

# LINEAR ALGEBRA

GABRIEL NAGY

Mathematics Department,  
Michigan State University,  
East Lansing, MI, 48824.

JULY 15, 2012

ABSTRACT. These are the lecture notes for the course MTH 415, Applied Linear Algebra, a one semester class taught in 2009-2012. These notes present a basic introduction to linear algebra with emphasis on few applications. Chapter 1 introduces systems of linear equations, the Gauss-Jordan method to find solutions of these systems which transforms the augmented matrix associated with a linear system into reduced echelon form, where the solutions of the linear system are simple to obtain. We end the Chapter with two applications of linear systems: First, to find approximate solutions to differential equations using the method of finite differences; second, to solve linear systems using floating-point numbers, as happens in a computer. Chapter 2 reviews matrix algebra, that is, we introduce the linear combination of matrices, the multiplication of appropriate matrices, and the inverse of a square matrix. We end the Chapter with the LU-factorization of a matrix. Chapter 3 reviews the determinant of a square matrix, the relation between a non-zero determinant and the existence of the inverse matrix, a formula for the inverse matrix using the matrix of cofactors, and the Cramer rule for the formula of the solution of a linear system with an invertible matrix of coefficients. The advanced part of the course really starts in Chapter 4 with the definition of vector spaces, subspaces, the linear dependence or independence of a set of vectors, bases and dimensions of vector spaces. Both finite and infinite dimensional vector spaces are presented, however finite dimensional vector spaces are the main interest in this notes. Chapter 5 presents linear transformations between vector spaces, the components of a linear transformation in a basis, and the formulas for the change of basis for both vector components and transformation components. Chapter 6 introduces a new structure on a vector space, called an inner product. The definition of an inner product is based on the properties of the dot product in  $\mathbb{R}^n$ . We study the notion of orthogonal vectors, orthogonal projections, best approximations of a vector on a subspace, and the Gram-Schmidt orthonormalization procedure. The central application of these ideas is the method of least-squares to find approximate solutions to inconsistent linear systems. One application is to find the best polynomial fit to a curve on a plane. Chapter 8 introduces the notion of a normed space, which is a vector space with a norm function which does not necessarily comes from an inner product. We study the main properties of the  $p$ -norms on  $\mathbb{R}^n$  or  $\mathbb{C}^n$ , which are useful norms in functional analysis. We briefly discuss induced operator norms. The last Section is an application of matrix norms. It discusses the condition number of a matrix and how to use this information to determine ill-conditioned linear systems. Finally, Chapter 9 introduces the notion of eigenvalue and eigenvector of a linear operator. We study diagonalizable operators, which are operators with diagonal matrix components in a basis of its eigenvectors. We also study functions of diagonalizable operators, with the exponential function as a main example. We also discuss how to apply these ideas to find solution of linear systems of ordinary differential equations.

## TABLE OF CONTENTS

<b>Overview</b>	1
Notation and conventions	1
Acknowledgments	2
<b>Chapter 1. Linear systems</b>	<b>3</b>
<b>1.1. Row and column pictures</b>	<b>3</b>
1.1.1. Row picture	3
1.1.2. Column picture	6
1.1.3. Exercises	10
<b>1.2. Gauss-Jordan method</b>	<b>11</b>
1.2.1. The augmented matrix	11
1.2.2. Gauss elimination operations	12
1.2.3. Square systems	13
1.2.4. Exercises	15
<b>1.3. Echelon forms</b>	<b>16</b>
1.3.1. Echelon and reduced echelon forms	16
1.3.2. The rank of a matrix	19
1.3.3. Inconsistent linear systems	21
1.3.4. Exercises	24
<b>1.4. Non-homogeneous equations</b>	<b>25</b>
1.4.1. Matrix-vector product	25
1.4.2. Linearity of matrix-vector product	27
1.4.3. Homogeneous linear systems	28
1.4.4. The span of vector sets	30
1.4.5. Non-homogeneous linear systems	31
1.4.6. Exercises	34
<b>1.5. Floating-point numbers</b>	<b>35</b>
1.5.1. Main definitions	35
1.5.2. The rounding function	37
1.5.3. Solving linear systems	38
1.5.4. Reducing rounding errors	40
1.5.5. Exercises	42
<b>Chapter 2. Matrix algebra</b>	<b>43</b>
<b>2.1. Linear transformations</b>	<b>43</b>
2.1.1. A matrix is a function	43
2.1.2. A matrix is a linear function	47
2.1.3. Exercises	50
<b>2.2. Linear combinations</b>	<b>51</b>
2.2.1. Linear combination of matrices	51
2.2.2. The transpose, adjoint, and trace of a matrix	52
2.2.3. Linear transformations on matrices	55
2.2.4. Exercises	56
<b>2.3. Matrix multiplication</b>	<b>57</b>
2.3.1. Algebraic definition	57
2.3.2. Matrix composition	59
2.3.3. Main properties	61
2.3.4. Block multiplication	62

2.3.5.	Matrix commutators	65
2.3.6.	Exercises	67
<b>2.4.</b>	<b>Inverse matrix</b>	68
2.4.1.	Main definition	68
2.4.2.	Properties of invertible matrices	71
2.4.3.	Computing the inverse matrix	72
2.4.4.	Exercises	74
<b>2.5.</b>	<b>Null and range spaces</b>	75
2.5.1.	Definition of the spaces	75
2.5.2.	Main properties	78
2.5.3.	Gauss operations	79
2.5.4.	Exercises	83
<b>2.6.</b>	<b>LU-factorization</b>	84
2.6.1.	Main definitions	84
2.6.2.	A sufficient condition	85
2.6.3.	Solving linear systems	89
2.6.4.	Exercises	90
<b>Chapter 3.</b>	<b>Determinants</b>	91
<b>3.1.</b>	<b>Definitions and properties</b>	91
3.1.1.	Determinant of $2 \times 2$ matrices	91
3.1.2.	Determinant of $3 \times 3$ matrices	95
3.1.3.	Determinant of $n \times n$ matrices	98
3.1.4.	Exercises	101
<b>3.2.</b>	<b>Applications</b>	102
3.2.1.	Inverse matrix formula	102
3.2.2.	Cramer's rule	104
3.2.3.	Determinants and Gauss operations	105
3.2.4.	Exercises	107
<b>Chapter 4.</b>	<b>Vector spaces</b>	108
<b>4.1.</b>	<b>Spaces and subspaces</b>	108
4.1.1.	Subspaces	110
4.1.2.	The span of finite sets	112
4.1.3.	Algebra of subspaces	113
4.1.4.	Internal direct sums	115
4.1.5.	Exercises	117
<b>4.2.</b>	<b>Linear dependence</b>	118
4.2.1.	Linearly dependent sets	118
4.2.2.	Main properties	119
4.2.3.	Exercises	121
<b>4.3.</b>	<b>Bases and dimension</b>	122
4.3.1.	Basis of a vector space	122
4.3.2.	Dimension of a vector space	125
4.3.3.	Extension of a set to a basis	127
4.3.4.	The dimension of subspace addition	128
4.3.5.	Exercises	130
<b>4.4.</b>	<b>Vector components</b>	131
4.4.1.	Ordered bases	131

4.4.2.	Vector components in a basis	131
4.4.3.	Exercises	136
<b>Chapter 5. Linear transformations</b>		<b>137</b>
<b>5.1.</b>	<b>Linear transformations</b>	<b>137</b>
5.1.1.	The null and range spaces	138
5.1.2.	Injections, surjections and bijections	139
5.1.3.	Nullity-Rank Theorem	141
5.1.4.	Exercises	143
<b>5.2.</b>	<b>Properties of linear transformations</b>	<b>144</b>
5.2.1.	The inverse transformation	144
5.2.2.	The vector space of linear transformations	147
5.2.3.	Linear functionals and the dual space	149
5.2.4.	Exercises	153
<b>5.3.</b>	<b>The algebra of linear operators</b>	<b>154</b>
5.3.1.	Polynomial functions of linear operators	156
5.3.2.	Functions of linear operators	157
5.3.3.	The commutator of linear operators	158
5.3.4.	Exercises	159
<b>5.4.</b>	<b>Transformation components</b>	<b>160</b>
5.4.1.	The matrix of a linear transformation	160
5.4.2.	Action as matrix-vector product	162
5.4.3.	Composition and matrix product	165
5.4.4.	Exercises	167
<b>5.5.</b>	<b>Change of basis</b>	<b>168</b>
5.5.1.	Vector components	168
5.5.2.	Transformation components	170
5.5.3.	Determinant and trace of linear operators	173
5.5.4.	Exercises	175
<b>Chapter 6. Inner product spaces</b>		<b>176</b>
<b>6.1.</b>	<b>Dot product</b>	<b>176</b>
6.1.1.	Dot product in $\mathbb{R}^2$	176
6.1.2.	Dot product in $\mathbb{F}^n$	179
6.1.3.	Exercises	184
<b>6.2.</b>	<b>Inner product</b>	<b>185</b>
6.2.1.	Inner product	185
6.2.2.	Inner product norm	188
6.2.3.	Norm distance	190
6.2.4.	Exercises	191
<b>6.3.</b>	<b>Orthogonal vectors</b>	<b>192</b>
6.3.1.	Definition and examples	192
6.3.2.	Orthonormal basis	194
6.3.3.	Vector components	196
6.3.4.	Exercises	198
<b>6.4.</b>	<b>Orthogonal projections</b>	<b>199</b>
6.4.1.	Orthogonal projection onto subspaces	199
6.4.2.	Orthogonal complement	202
6.4.3.	Exercises	205

<b>6.5. Gram-Schmidt method</b>	206
6.5.1. Exercises	210
<b>6.6. The adjoint operator</b>	211
6.6.1. The Riesz Representation Theorem	211
6.6.2. The adjoint operator	212
6.6.3. Normal operators	213
6.6.4. Bilinear forms	214
6.6.5. Exercises	216
<b>Chapter 7. Approximation methods</b>	217
<b>7.1. Best approximation</b>	217
7.1.1. Fourier expansions	217
7.1.2. Null and range spaces of a matrix	219
7.1.3. Exercises	222
<b>7.2. Least squares</b>	223
7.2.1. The normal equation	223
7.2.2. Least squares fit	226
7.2.3. Linear correlation	229
7.2.4. QR-factorization	230
7.2.5. Exercises	232
<b>7.3. Finite difference method</b>	233
7.3.1. Differential equations	233
7.3.2. Difference quotients	234
7.3.3. Method of finite differences	236
7.3.4. Exercises	241
<b>7.4. Finite element method</b>	242
7.4.1. Differential equations	242
7.4.2. The Galerkin method	243
7.4.3. Finite element method	243
7.4.4. Exercises	244
<b>Chapter 8. Normed spaces</b>	245
<b>8.1. The <math>p</math>-norm</b>	245
8.1.1. Not every norm is an inner product norm	249
8.1.2. Equivalent norms	252
8.1.3. Exercises	255
<b>8.2. Operator norms</b>	256
8.2.1. Exercises	262
<b>8.3. Condition numbers</b>	263
8.3.1. Exercises	265
<b>Chapter 9. Spectral decomposition</b>	266
<b>9.1. Eigenvalues and eigenvectors</b>	266
9.1.1. Main definitions	266
9.1.2. Characteristic polynomial	269
9.1.3. Eigenvalue multiplicities	270
9.1.4. Operators with distinct eigenvalues	272
9.1.5. Exercises	274
<b>9.2. Diagonalizable operators</b>	275

9.2.1.	Eigenvectors and diagonalization	275
9.2.2.	Functions of diagonalizable operators	277
9.2.3.	The exponential of diagonalizable operators	279
9.2.4.	Exercises	282
<b>9.3.</b>	<b>Differential equations</b>	<b>283</b>
9.3.1.	Non-repeated real eigenvalues	286
9.3.2.	Non-repeated complex eigenvalues	290
9.3.3.	Exercises	294
<b>9.4.</b>	<b>Normal operators</b>	<b>295</b>
9.4.1.	Exercises	300
<b>Chapter 10.</b>	<b>Appendix</b>	<b>301</b>
10.1.	Review exercises	301
10.2.	Practice Exams	310
10.3.	Answers to exercises	317
10.4.	Solutions to Practice Exams	336
	References	356

## OVERVIEW

Linear algebra is a collection of ideas involving algebraic systems of linear equations, vectors and vector spaces, and linear transformations between vector spaces.

Algebraic equations are called a system when there is more than one equation, and they are called linear when the unknown appears as a multiplicative factor with power zero or one. An example of a linear system of two equations in two unknowns is given in Eqs. (1.3)-(1.4) below. Systems of linear equations are the main subject of Chapter 1.

Examples of vectors are oriented segments on a line, plane, or space. An oriented segment is an ordered pair of points in these sets. Such ordered pair can be drawn as an arrow that starts on the first point and ends on the second point. Fix a preferred point in the line, plane or space, called the origin point, and then there exists a one-to-one correspondence between points in these sets and arrows that start at the origin point. The set of oriented segments with common origin in a line, plane, and space are called  $\mathbb{R}$ ,  $\mathbb{R}^2$  and  $\mathbb{R}^3$ , respectively. A sketch of vectors in these sets can be seen in Fig. 1. Two operations are defined on oriented segments with common origin point: An oriented segment can be stretched or compressed; and two oriented segments with the same origin point can be added using the parallelogram law. An addition of several stretched or compressed vectors is called a linear combination. The set of all oriented segments with common origin point together with this operation of linear combination is the essential structure called vector space. The origin of the word “space” in the term “vector space” originates precisely in these examples, which were associated with the physical space.

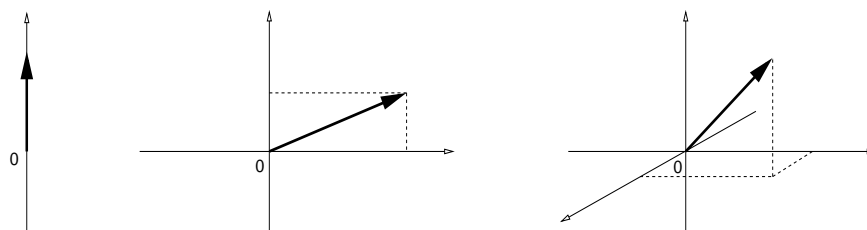


FIGURE 1. Example of vectors in the line, plane, and space, respectively.

Linear transformations are a particular type of functions between vector spaces that preserve the operation of linear combination. An example of a linear transformation is a  $2 \times 2$  matrix  $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$  together with a matrix-vector product that specifies how this matrix transforms a vector on the plane into another vector on the plane. The result is thus a function  $A : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ .

These notes try to be an elementary introduction to linear algebra with few applications.

**Notation and conventions.** We use the notation  $\mathbb{F} \in \{\mathbb{R}, \mathbb{C}\}$  to mean that  $\mathbb{F} = \mathbb{R}$  or  $\mathbb{F} = \mathbb{C}$ . Vectors will be denoted by boldface letters, like  $\mathbf{u}$  and  $\mathbf{v}$ . The exception are a column vectors in  $\mathbb{F}^n$  which are denoted in sanserif, like  $u$  and  $v$ . This notation permits to differentiate between a vector and its components on a basis. In a similar way, linear transformations between vector spaces are denoted by boldface capital letters, like  $\mathbf{T}$  and  $\mathbf{S}$ . The exception are matrices in  $\mathbb{F}^{m,n}$  which are denoted by capital sanserif letters like  $A$  and  $B$ . Again, this notation is useful to differentiate between a linear transformation and its

components on two bases. Below is a list of few mathematical symbols used in these notes:

$\mathbb{R}$	Set of real numbers,	$\mathbb{Q}$	Set of rational numbers,
$\mathbb{Z}$	Set of integer numbers,	$\mathbb{N}$	Set of positive integers,
$\{\mathbf{0}\}$	Zero set,	$\emptyset$	Empty set,
$\cup$	Union of sets,	$\cap$	Intersection of sets,
$:=$	Definition,	$\Rightarrow$	Implies,
$\forall$	For all,	$\exists$	There exists,
<b>Proof</b>	Beginning of a proof,	$\square$	End of a proof,
<b>Example</b>	Beginning of an example,	$\triangleleft$	End of an example.

**Acknowledgments.** I thanks all my students for pointing out several misprints and for helping make these notes more readable. I am specially grateful to Zhuo Wang and Wenning Feng.



## CHAPTER 1. LINEAR SYSTEMS

## 1.1. ROW AND COLUMN PICTURES

1.1.1. **Row picture.** A central problem in linear algebra is to find solutions of a system of linear equations. A  $2 \times 2$  *linear system* consists of two linear equations in two unknowns. More precisely, given the real numbers  $A_{11}$ ,  $A_{12}$ ,  $A_{21}$ ,  $A_{22}$ ,  $b_1$ , and  $b_2$ , find all numbers  $x$  and  $y$  solutions of both equations

$$A_{11}x + A_{12}y = b_1, \quad (1.1)$$

$$A_{21}x + A_{22}y = b_2. \quad (1.2)$$

These equations are called a system because there is more than one equation, and they are called linear because the unknowns,  $x$  and  $y$ , appear as multiplicative factors with power zero or one (for example, there is no term proportional to  $x^2$  or to  $y^3$ ). The *row picture* of a linear system is the method of finding solutions to this system as the intersection of all solutions to every single equation in the system. The individual equations are called *row equations*, or simply rows of the system.

**EXAMPLE 1.1.1:** Find all the numbers  $x$  and  $y$  solutions of the  $2 \times 2$  linear system

$$2x - y = 0, \quad (1.3)$$

$$-x + 2y = 3. \quad (1.4)$$

**SOLUTION:** The solution to each row of the system above is found geometrically in Fig. 2.

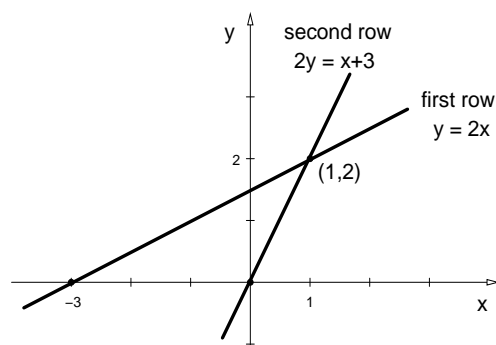


FIGURE 2. The solution of a  $2 \times 2$  linear system in the row picture is the intersection of the two lines, which are the solutions of each row equation.

Analytically, the solution can be found by substitution:

$$2x - y = 0 \quad \Rightarrow \quad y = 2x \quad \Rightarrow \quad -x + 4x = 3 \quad \Rightarrow \quad \begin{cases} x = 1, \\ y = 2. \end{cases}$$

◁

An interesting property of the solutions to any  $2 \times 2$  linear system is simple to prove using the row picture, and it is the following result.

**Theorem 1.1.1.** *Given any  $2 \times 2$  linear system, only one of the following statements holds:*

- (i) *There exists a unique solution;*
- (ii) *There exist infinitely many solutions;*
- (iii) *There exists no solution.*

It is interesting to remark what cannot happen, for example there is no  $2 \times 2$  linear system having only two solutions. Unlike the quadratic equation  $x^2 - 5x + 6 = 0$ , which has two solutions given by  $x = 2$  and  $x = 3$ , a  $2 \times 2$  linear system has only one solution, or infinitely many solutions, or no solution at all. Examples of these three cases, respectively, are given in Fig. 3.

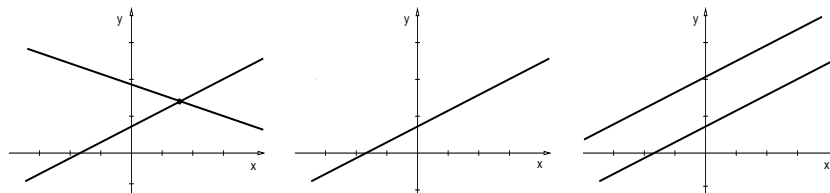


FIGURE 3. An example of the cases given in Theorem 1.1.1, cases (i)-(iii).

**Proof of Theorem 1.1.1:** The solutions of each equation in a  $2 \times 2$  linear system represents a line in  $\mathbb{R}^2$ . Two lines in  $\mathbb{R}^2$  can intersect at a point, or can be coincident, or can be parallel but not coincident. These are the cases given in (i)-(iii). This establishes the Theorem.  $\square$

We now generalize the definition of a  $2 \times 2$  linear system given in the Example 1.1.1 to  $m$  equations of  $n$  unknowns.

**Definition 1.1.2.** An  $m \times n$  **linear system** is a set of  $m \geq 1$  linear equations in  $n \geq 1$  unknowns is the following: Given the coefficients numbers  $A_{ij}$  and the source numbers  $b_i$ , with  $i = 1, \dots, m$  and  $j = 1, \dots, n$ , find the real numbers  $x_j$  solutions of

$$\begin{aligned} A_{11}x_1 + \dots + A_{1n}x_n &= b_1 \\ &\vdots \\ A_{m1}x_1 + \dots + A_{mn}x_n &= b_m. \end{aligned}$$

Furthermore, an  $m \times n$  linear system is called **consistent** iff it has a solution, and it is called **inconsistent** iff it has no solutions.

**EXAMPLE 1.1.2:** Find all numbers  $x_1$ ,  $x_2$  and  $x_3$  solutions of the  $2 \times 3$  linear system

$$\begin{aligned} x_1 + 2x_2 + x_3 &= 1 \\ -3x_1 - x_2 - 8x_3 &= 2 \end{aligned} \tag{1.5}$$

**SOLUTION:** Compute  $x_1$  from the first equation,  $x_1 = 1 - 2x_2 - x_3$ , and substitute it in the second equation,

$$-3(1 - 2x_2 - x_3) - x_2 - 8x_3 = 2 \Rightarrow 5x_2 - 5x_3 = 5 \Rightarrow x_2 = 1 + x_3.$$

Substitute the expression for  $x_2$  in the equation above for  $x_1$ , and we obtain

$$x_1 = 1 - 2(1 + x_3) - x_3 = 1 - 2 - 2x_3 - x_3 \Rightarrow x_1 = -1 - 3x_3.$$

Since there is no condition on  $x_3$ , the system above has infinitely many solutions parametrized by the number  $x_3$ . We conclude that  $x_1 = -1 - 3x_3$ , and  $x_2 = 1 + x_3$ , while  $x_3$  is free.  $\triangleleft$

**EXAMPLE 1.1.3:** Find all numbers  $x_1$ ,  $x_2$  and  $x_3$  solutions of the  $3 \times 3$  linear system

$$\begin{aligned} 2x_1 + x_2 + x_3 &= 2 \\ -x_1 + 2x_2 &= 1 \\ x_1 - x_2 + 2x_3 &= -2. \end{aligned} \tag{1.6}$$

**SOLUTION:** While the row picture is appropriate to solve small systems of linear equations, it becomes difficult to carry out on  $3 \times 3$  and bigger linear systems. The solution  $x_1, x_2, x_3$  of the system above can be found as follows: Substitute the second equation into the first,

$$x_1 = -1 + 2x_2 \Rightarrow x_3 = 2 - 2x_1 - x_2 = 2 + 2 - 4x_2 - x_2 \Rightarrow x_3 = 4 - 5x_2;$$

then, substitute the second equation and  $x_3 = 4 - 5x_2$  into the third equation,

$$(-1 + 2x_2) - x_2 + 2(4 - 5x_2) = -2 \Rightarrow x_2 = 1,$$

and then, substituting backwards,  $x_1 = 1$  and  $x_3 = -1$ . We conclude that the solution is a single point in space given by  $(1, 1, -1)$ .  $\triangleleft$

The solution of each separate equation in the examples above represents a plane in  $\mathbb{R}^3$ . A solution to the whole system is a point that belongs to the three planes. In the  $3 \times 3$  system in Example 1.1.3 above there is a unique solution, the point  $(1, 1, -1)$ , which means that the three planes intersect at a single point. In the general case, a  $3 \times 3$  system can have a unique solution, infinitely many solutions or no solutions at all, depending on how the three planes in space intersect among them. The case with unique solution was represented in Fig. 4, while two possible situations corresponding to no solution are given in Fig. 5. Finally, two cases of  $3 \times 3$  linear system having infinitely many solutions are pictured in Fig 6, where in the first case the solutions form a line, and in the second case the solutions form a plane because the three planes coincide.

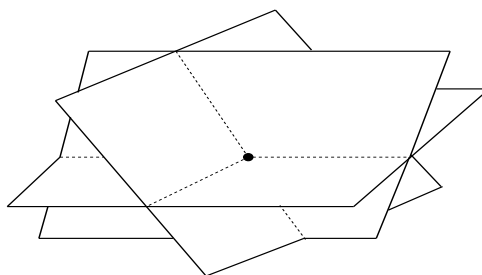


FIGURE 4. Planes representing the solutions of each row equation in a  $3 \times 3$  linear system having a unique solution.

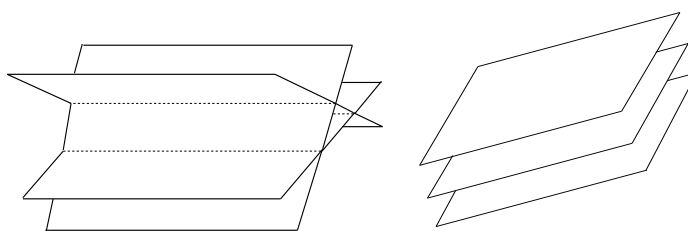


FIGURE 5. Two cases of planes representing the solutions of each row equation in  $3 \times 3$  linear systems having no solutions.

Solutions of linear systems with more than three unknowns can not be represented in the three dimensional space. Besides, the substitution method becomes more involved to solve. As a consequence, alternative ideas are needed to solve such systems. We now discuss one

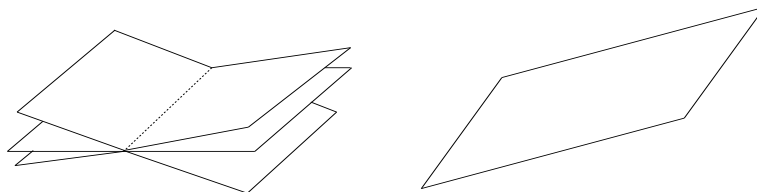


FIGURE 6. Two cases of planes representing the solutions of each row equation in  $3 \times 3$  linear systems having infinitely many solutions.

of such ideas, the use of vectors to interpret and find solutions of linear systems. In the next Section we introduce another idea, the use of matrices and vectors to solve linear systems following the Gauss-Jordan method. This latter procedure is suitable to solve large systems of linear equations in an efficient way.

**1.1.2. Column picture.** Consider again the linear system in Eqs. (1.3)-(1.4) and introduce a change in the names of the unknowns, calling them  $x_1$  and  $x_2$  instead of  $x$  and  $y$ . The problem is to find the numbers  $x_1$ , and  $x_2$  solutions of

$$2x_1 - x_2 = 0, \quad (1.7)$$

$$-x_1 + 2x_2 = 3. \quad (1.8)$$

We know that the answer is  $x_1 = 1$ ,  $x_2 = 2$ . The row picture consisted in solving each row separately. The main idea in the column picture is to interpret the  $2 \times 2$  linear system as an addition of new objects, column vectors, in the following way,

$$\begin{bmatrix} 2 \\ -1 \end{bmatrix} x_1 + \begin{bmatrix} -1 \\ 2 \end{bmatrix} x_2 = \begin{bmatrix} 0 \\ 3 \end{bmatrix}. \quad (1.9)$$

The new objects are called column vectors and they are denoted as follows,

$$A_1 = \begin{bmatrix} 2 \\ -1 \end{bmatrix}, \quad A_2 = \begin{bmatrix} -1 \\ 2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ 3 \end{bmatrix}.$$

We can represent these vectors in the plane, as it is shown in Fig. 7.

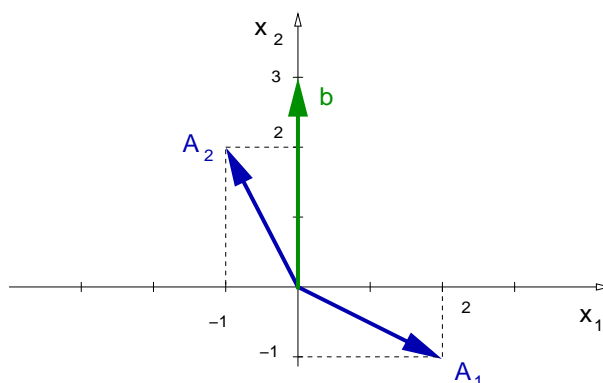


FIGURE 7. Graphical representation of column vectors in the plane.

The column vector interpretation of a  $2 \times 2$  linear system determines an addition law of vectors and a multiplication law of a vector by a number. In the example above, we

know that the solution is given by  $x_1 = 1$  and  $x_2 = 2$ , therefore in the column picture interpretation the following equation must hold

$$\begin{bmatrix} 2 \\ -1 \end{bmatrix} + \begin{bmatrix} -1 \\ 2 \end{bmatrix} 2 = \begin{bmatrix} 0 \\ 3 \end{bmatrix}.$$

The study of this example suggests that the multiplication law of a vector by numbers and the addition law of two vectors can be defined by the following equations, respectively,

$$\begin{bmatrix} -1 \\ 2 \end{bmatrix} 2 = \begin{bmatrix} (-1)2 \\ (2)2 \end{bmatrix}, \quad \begin{bmatrix} 2 \\ -1 \end{bmatrix} + \begin{bmatrix} -2 \\ 4 \end{bmatrix} = \begin{bmatrix} 2-2 \\ -1+4 \end{bmatrix}.$$

The study of several examples of  $2 \times 2$  linear systems in the column picture determines the following definition.

**Definition 1.1.3.** The *linear combination* of the 2-vectors  $\mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$  and  $\mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$ , with the real numbers  $a$  and  $b$ , is defined as follows,

$$a \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} + b \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} au_1 + bv_1 \\ au_2 + bv_2 \end{bmatrix}.$$

A linear combination includes the particular cases of addition ( $a = b = 1$ ), and multiplication of a vector by a number ( $b = 0$ ), respectively given by,

$$\begin{bmatrix} u_1 \\ u_2 \end{bmatrix} + \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} u_1 + v_1 \\ u_2 + v_2 \end{bmatrix}, \quad a \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} au_1 \\ au_2 \end{bmatrix}.$$

The addition law in terms of components is represented graphically by the parallelogram law, as it can be seen in Fig. 8. The multiplication of a vector by a number  $a$  affects the length and direction of the vector. The product  $a\mathbf{u}$  stretches the vector  $\mathbf{u}$  when  $a > 1$  and it compresses  $\mathbf{u}$  when  $0 < a < 1$ . If  $a < 0$  then it reverses the direction of  $\mathbf{u}$  and it stretches when  $a < -1$  and compresses when  $-1 < a < 0$ . Fig. 8 represents some of these possibilities. Notice that the difference of two vectors is a particular case of the parallelogram law, as it can be seen in Fig. 9.

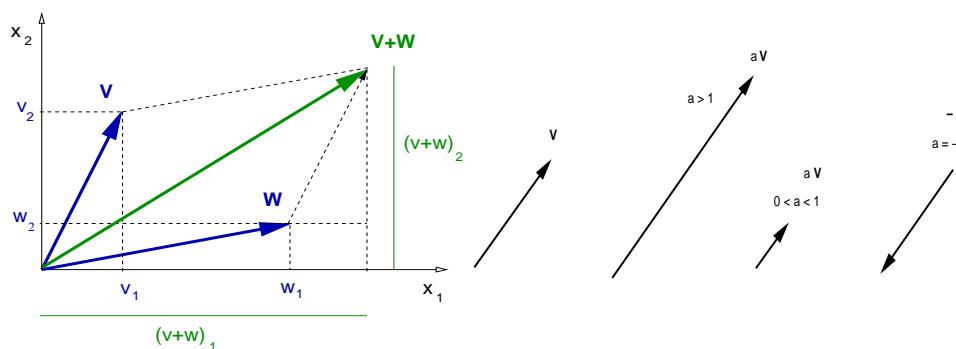


FIGURE 8. The addition of vectors can be computed with the parallelogram law. The multiplication of a vector by a number stretches or compresses the vector, and changes its direction in the case that the number is negative.

The *column picture* interpretation of a general  $2 \times 2$  linear system given in Eqs. (1.1)-(1.2) is the following: Introduce the coefficient and source column vectors

$$\mathbf{A}_1 = \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix}, \quad \mathbf{A}_2 = \begin{bmatrix} A_{12} \\ A_{22} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}, \quad (1.10)$$

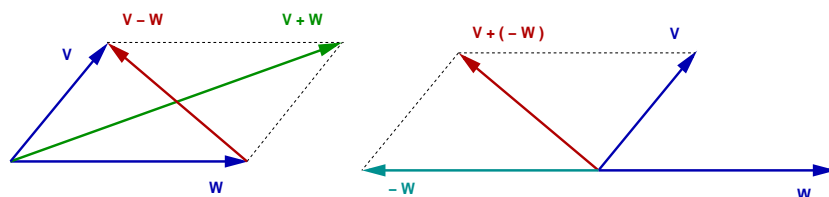


FIGURE 9. The difference of two vectors is a particular case of the parallelogram law of addition of vectors.

and then find the coefficients  $x_1$  and  $x_2$  that change the length of the coefficient column vectors  $A_1$  and  $A_2$  such that they add up to the source column vector  $\mathbf{b}$ , that is,

$$A_1 x_1 + A_2 x_2 = \mathbf{b}.$$

For example, the column picture of the linear system in Eqs. (1.7)-(1.8) is given in Eq. (1.9). The solution of this system are the numbers  $x_1 = 1$  and  $x_2 = 2$ , and this solution is represented in Fig. 10.

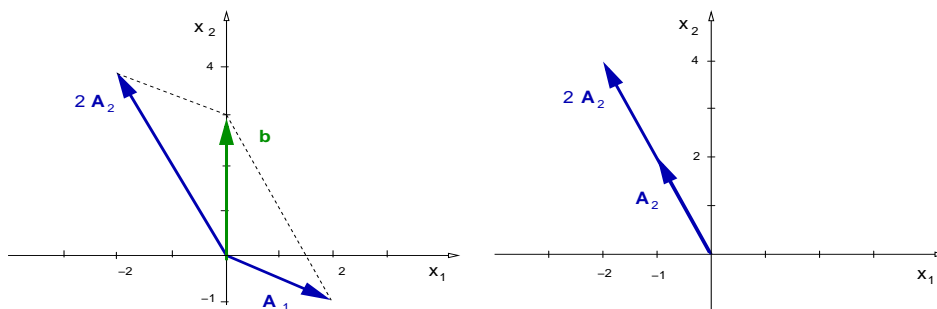


FIGURE 10. Representation of the solution of a  $2 \times 2$  linear system in the column picture.

The existence and uniqueness of solutions in the case of  $2 \times 2$  systems can be studied geometrically in the column picture as it was done in the row picture. In the latter case we have seen that all possible  $2 \times 2$  systems fall into one of these three cases, unique solution, infinitely many solutions and no solutions at all. The proof was to study all possible ways two lines can intersect on the plane. The same existence and uniqueness statement can be proved in the column picture. In Fig. 11 we present these three cases in both row and column pictures. In the latter case the proof is to study all possible relative positions of the column vectors  $A_1$ ,  $A_2$ , and  $\mathbf{b}$  on the plane.

We see in Fig. 11 that the first case corresponds to a system with unique solution. There is only one linear combination of the coefficient vectors  $A_1$  and  $A_2$  which adds up to  $\mathbf{b}$ . The reason is that the coefficient vectors are not proportional to each other. The second case corresponds to the case of infinitely many solutions. The coefficient vectors are proportional to each other and the source vector  $\mathbf{b}$  is also proportional to them. So, there are infinitely many linear combinations of the coefficient vectors that add up to the source vector. The last case corresponds to the case of no solutions. While the coefficient vectors are proportional to each other, the source vector is not proportional to them. So, there is no linear combination of the coefficient vectors that add up to the source vector.

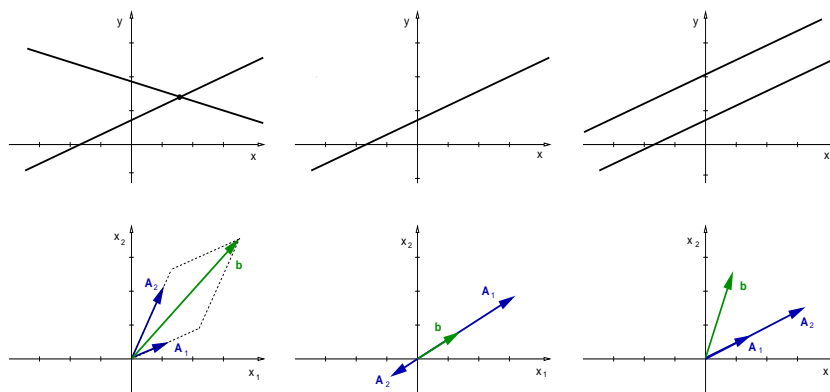


FIGURE 11. Examples of a solutions of general  $2 \times 2$  linear systems having a unique, infinitely many, and no solution, represented in the row picture and in the column picture.

The ideas in the column picture can be generalized to  $m \times n$  linear equations, which gives rise to the generalization to  $m$ -vectors of the definitions of linear combination presented above.

**Definition 1.1.4.** The *linear combination* of the  $m$ -vectors  $\mathbf{u} = \begin{bmatrix} u_1 \\ \vdots \\ u_m \end{bmatrix}$  and  $\mathbf{v} = \begin{bmatrix} v_1 \\ \vdots \\ v_m \end{bmatrix}$

with the real numbers  $a, b$  is defined as follows

$$a \begin{bmatrix} u_1 \\ \vdots \\ u_m \end{bmatrix} + b \begin{bmatrix} v_1 \\ \vdots \\ v_m \end{bmatrix} = \begin{bmatrix} au_1 + bv_1 \\ \vdots \\ au_m + bv_m \end{bmatrix}.$$

This definition can be generalized to an arbitrary number of vectors. Column vectors provide a new way to denote an  $m \times n$  system of linear equations.

**Definition 1.1.5.** An  $m \times n$  *linear system* of  $m \geq 1$  linear equations in  $n \geq 1$  unknowns is the following: Given the coefficient  $m$ -vectors  $\mathbf{A}_1, \dots, \mathbf{A}_n$  and the source  $m$ -vector  $\mathbf{b}$ , find the real numbers  $x_1, \dots, x_n$  solution of the linear combination

$$\mathbf{A}_1 x_1 + \dots + \mathbf{A}_n x_n = \mathbf{b}.$$

For example, recall the  $3 \times 3$  system given as the second system in Eq. (1.6). This system in the column picture is the following: Find numbers  $x_1, x_2$  and  $x_3$  such that

$$\begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix} x_1 + \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix} x_2 + \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix} x_3 = \begin{bmatrix} 2 \\ 1 \\ -2 \end{bmatrix}. \quad (1.11)$$

These are the main ideas in the column picture. We will see later that linear algebra emerges from the column picture. In the next Section we introduce the Gauss-Jordan method, which is a procedure to solve large systems of linear equations in an efficient way.

**Further reading.** For more details on the row picture see Section 1.1 in Lay's book [2]. There is a clear explanation of the column picture in Section 1.3 in Lay's book [2]. See also Section 1.2 in Strang's book [4] for a shorter summary of both the row and column pictures.

## 1.1.3. Exercises.

1.1.1.- Use the substitution method to find the solutions to the  $2 \times 2$  linear system

$$2x - y = 1,$$

$$x + y = 5.$$

1.1.2.- Sketch the three lines solution of each row in the system

$$x + 2y = 2$$

$$x - y = 2$$

$$y = 1.$$

Is this linear system consistent?

1.1.3.- Sketch a graph representing the solutions of each row in the following **non-linear** system, and decide whether it has solutions or not,

$$x^2 + y^2 = 4$$

$$x - y = 0.$$

1.1.4.- Graph on the plane the solution of each individual equation of the  $3 \times 2$  linear system system

$$3x - y = 0,$$

$$x + 2y = 4,$$

$$-x + y = -2,$$

and determine whether the system is consistent or inconsistent.

1.1.5.- Show that the  $3 \times 3$  linear system

$$x + y + z = 2,$$

$$x + 2y + 3z = 1,$$

$$y + 2z = 0,$$

is inconsistent, by finding a combination of the three equations that adds up to the equation  $0 = 1$ .

1.1.6.- Find all values of the constant  $k$  such that there exists infinitely many solutions to the  $2 \times 2$  linear system

$$kx + 2y = 0,$$

$$x + \frac{k}{2}y = 0.$$

1.1.7.- Sketch a graph of the vectors

$$A_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad A_2 = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ -1 \end{bmatrix}.$$

Is the linear system  $A_1x_1 + A_2x_2 = b$  consistent? If the answer is “yes,” find the solution.

1.1.8.- Consider the vectors

$$A_1 = \begin{bmatrix} 4 \\ 2 \end{bmatrix}, \quad A_2 = \begin{bmatrix} -2 \\ -1 \end{bmatrix}, \quad b = \begin{bmatrix} 2 \\ 0 \end{bmatrix}.$$

(a) Graph the vectors  $A_1$ ,  $A_2$  and  $b$  on the plane.

(b) Is the linear system  $A_1x_1 + A_2x_2 = b$  consistent?

(c) Given the vector  $c = \begin{bmatrix} 6 \\ 3 \end{bmatrix}$ , is the linear system  $A_1x_1 + A_2x_2 = c$  consistent? If the answer is “yes,” is the solution unique?

1.1.9.- Consider the vectors

$$A_1 = \begin{bmatrix} 4 \\ 2 \end{bmatrix}, \quad A_2 = \begin{bmatrix} -2 \\ -1 \end{bmatrix},$$

and given a real number  $h$ , set  $c = \begin{bmatrix} 2 \\ h \end{bmatrix}$ . Find all values of  $h$  such that the system  $A_1x_1 + A_2x_2 = c$  is consistent.

1.1.10.- Show that the three vectors below lie on the same plane, by expressing the third vector as a linear combination of the first two, where

$$A_1 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \quad A_2 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}, \quad A_3 = \begin{bmatrix} 1 \\ 3 \\ 2 \end{bmatrix}.$$

Is the linear system

$$A_1x_1 + A_2x_2 + A_3x_3 = 0$$

consistent? If the answer is “yes,” is the solution unique?



## 1.2. GAUSS-JORDAN METHOD

1.2.1. **The augmented matrix.** Solutions to  $m \times n$  linear systems can be obtained in an efficient way using the Gauss-Jordan method. Efficient here means to perform as few algebraic operations as possible either to find the solution or to show that the solution does not exist. Before introducing this method, we need several definitions.

**Definition 1.2.1.** The *coefficients matrix*, the *source vector*, and the *augmented matrix* of the  $m \times n$  linear system

$$\begin{aligned} A_{11}x_1 + \cdots + A_{1n}x_n &= b_1 \\ &\vdots \\ A_{m1}x_1 + \cdots + A_{mn}x_n &= b_m, \end{aligned}$$

are given by, respectively,

$$A = \begin{bmatrix} A_{11} & \cdots & A_{1n} \\ \vdots & & \vdots \\ A_{m1} & \cdots & A_{mn} \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}, \quad [A|\mathbf{b}] = \left[ \begin{array}{ccc|c} A_{11} & \cdots & A_{1n} & b_1 \\ \vdots & & \vdots & \vdots \\ A_{m1} & \cdots & A_{mn} & b_m \end{array} \right].$$

We call  $A$  an  $m \times n$  matrix, with  $m$  the number of rows and  $n$  the number of columns. The source vector  $\mathbf{b}$  is a particular case of an  $m \times 1$  matrix. The augmented matrix of an  $m \times n$  linear system is given by the coefficient matrix and the source vector together, hence it is an  $m \times (n + 1)$  matrix.

**EXAMPLE 1.2.1:** Find the coefficient matrix, the source vector and the augmented matrix of the  $2 \times 2$  linear system

$$2x_1 - x_2 = 0, \tag{1.12}$$

$$-x_1 + 2x_2 = 3. \tag{1.13}$$

**SOLUTION:** The coefficient matrix is  $2 \times 2$ , the source vector is  $2 \times 1$ , and the augmented matrix is  $2 \times 3$ , given respectively by

$$A = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ 3 \end{bmatrix}, \quad [A|\mathbf{b}] = \left[ \begin{array}{cc|c} 2 & -1 & 0 \\ -1 & 2 & 3 \end{array} \right]. \tag{1.14}$$

◁

**EXAMPLE 1.2.2:** Find the coefficient matrix, the source vector and the augmented matrix of the  $2 \times 3$  linear system

$$2x_1 - x_2 = 0,$$

$$-x_1 + 2x_2 + 3x_3 = 0.$$

**SOLUTION:** The coefficient matrix is  $2 \times 3$ , the source vector is  $2 \times 1$ , and the augmented matrix is  $2 \times 4$ , given respectively by

$$A = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & 3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad [A|\mathbf{b}] = \left[ \begin{array}{ccc|c} 2 & -1 & 0 & 0 \\ -1 & 2 & 3 & 0 \end{array} \right].$$

Notice that the coefficient matrix in this example is equal to the augmented matrix in the previous Example. This means that the vertical separator in the definition of the augmented

matrix is important. If one does not write down the vertical separator in Example 1.2.1, then one is actually working on the system in Example 1.2.2.  $\triangleleft$

We also use the alternative notation  $\mathbf{A} = [A_{ij}]$  to denote a matrix with components  $A_{ij}$ , and  $\mathbf{b} = [b_i]$  to denote a vector with components  $b_i$ , where  $i = 1, \dots, m$  and  $j = 1, \dots, n$ .

**Definition 1.2.2.** The *diagonal* of an  $m \times n$  matrix  $\mathbf{A} = [A_{ij}]$  is the set of all coefficients  $A_{ii}$  for  $i = 1, \dots, m$ . A matrix is *upper triangular* iff every coefficient below the diagonal vanishes, and *lower triangular* iff every coefficient above the diagonal vanishes.

**EXAMPLE 1.2.3:** Highlight the diagonal coefficients in  $3 \times 3$ ,  $2 \times 3$  and  $3 \times 2$  matrices.

**SOLUTION:** The diagonal coefficients in the  $3 \times 3$ ,  $2 \times 3$  and  $3 \times 2$  matrices below are highlighted in green, where the coefficients with \* denote non-diagonal coefficients:

$$\begin{bmatrix} A_{11} & * & * \\ * & A_{22} & * \\ * & * & A_{33} \end{bmatrix}, \quad \begin{bmatrix} A_{11} & * & * \\ * & A_{22} & * \end{bmatrix}, \quad \begin{bmatrix} A_{11} & * \\ * & A_{22} \\ * & * \end{bmatrix}.$$

$\triangleleft$

**EXAMPLE 1.2.4:** Write down the most general  $3 \times 3$ ,  $2 \times 3$  and  $3 \times 2$  upper triangular matrices, highlighting the diagonal elements.

**SOLUTION:** The following matrices are upper triangular:

$$\begin{bmatrix} A_{11} & * & * \\ 0 & A_{22} & * \\ 0 & 0 & A_{33} \end{bmatrix}, \quad \begin{bmatrix} A_{11} & * & * \\ 0 & A_{22} & * \end{bmatrix}, \quad \begin{bmatrix} A_{11} & * \\ 0 & A_{22} \\ 0 & 0 \end{bmatrix}.$$

$\triangleleft$

**1.2.2. Gauss elimination operations.** The Gauss-Jordan Method is a procedure performed on the augmented matrix of a linear system. It consists on a sequence of operations, called Gauss elimination operations, that change the augmented matrix of the system but they do not change the solutions of the system. The Gauss elimination operations were already known in China around 200 BC. We call them after Carl Friedrich Gauss, since he made them very popular around 1810, when he used them to study the orbit of the asteroid Pallas, giving a systematic method to solve a  $6 \times 6$  algebraic linear system.

**Definition 1.2.3.** The *Gauss elimination operations* are three operations on a matrix:

- (i) To add a multiple of one row to another row;
- (ii) To interchange two rows;
- (iii) To multiply a row by a non-zero number.

These operations are respectively represented by the symbols given in Fig. 12.

$$\left[ \begin{array}{c} \text{---} \\ \text{---} \end{array} \right] \overset{a}{\curvearrowright} \quad \left[ \begin{array}{c} \text{---} \\ \text{---} \end{array} \right] \curvearrowright \quad \left[ \text{---} \right] \xleftarrow{a \neq 0}$$

FIGURE 12. A representation of the Gauss elimination operations (i), (ii) and (iii), respectively.

**REMARK:** If the factor in operation (iii) is allowed to be zero, then multiplying an equation by zero modifies the solutions of the original system, since such operation on an equation is equivalent to erase that equation from the system. This explains why the factor in operation (iii) is constrained to be non-zero.

The Gauss elimination operations change the coefficients of the augmented matrix of a system but do not change its solution. Two systems of linear equations having the same solutions are called *equivalent*. The Gauss-Jordan Method is an algorithm using these operations that transforms any linear system into an equivalent system where the solutions are given explicitly. We describe the Gauss-Jordan Method only through examples.

**EXAMPLE 1.2.5:** Find the solution of the  $2 \times 2$  linear system with augmented matrix given in Eq. (1.14) using Gauss operations.

**SOLUTION:**

$$\begin{aligned} \left[ \begin{array}{cc|c} 2 & -1 & 0 \\ -1 & 2 & 3 \end{array} \right] &\rightarrow \left[ \begin{array}{cc|c} 2 & -1 & 0 \\ -2 & 4 & 6 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 2 & -1 & 0 \\ 0 & 3 & 6 \end{array} \right] \rightarrow \\ &\left[ \begin{array}{cc|c} 2 & -1 & 0 \\ 0 & 1 & 2 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 2 & 0 & 2 \\ 0 & 1 & 2 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 0 & 1 \\ 0 & 1 & 2 \end{array} \right]. \end{aligned}$$

The Gauss operations have changed the augmented matrix of the original system as follows:

$$\left[ \begin{array}{cc|c} 2 & -1 & 0 \\ -1 & 2 & 3 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 0 & 1 \\ 0 & 1 & 2 \end{array} \right].$$

Since the Gauss operations do not change the solution of the associated linear systems, the augmented matrices above imply that the following two linear systems have the same solutions,

$$\left. \begin{array}{l} 2x_1 - x_2 = 0, \\ -x_1 + 2x_2 = 3. \end{array} \right\} \Leftrightarrow \left\{ \begin{array}{l} x_1 = 1, \\ x_2 = 2. \end{array} \right.$$

On the last system the solution is given explicitly,  $x_1 = 1$ ,  $x_2 = 2$ , since no additional algebraic operations are needed to find the solution.  $\triangleleft$

A precise way to define the notion of an augmented matrix corresponding to a linear system with solutions “easy to read” is captured in the notions of echelon form of a matrix, and reduced echelon form of a matrix. Next Section is dedicated to present these notions in some detail.

**1.2.3. Square systems.** In the rest of this Section we study a particular type of linear systems having the same number of equations and unknowns.

**Definition 1.2.4.** An  $m \times n$  linear system is called a *square system* iff holds that  $m = n$ .

An example of a square system is the  $2 \times 2$  linear system in Eqs. (1.12)-(1.13). In the rest of this Section we introduce the Gauss operations and back substitution method to find solutions to square linear systems. We later on compare this method with the Gauss-Jordan Method restricted to square linear systems.

We start using *Gauss operations and back substitution* to find solutions to square  $n \times n$  linear systems. This method has two main parts: First, use Gauss operations to transform the augmented matrix of the system into an upper triangular form. Second, use back substitution to compute the solution to the system.

**EXAMPLE 1.2.6:** Use the Gauss operations and back substitution to solve the  $3 \times 3$  system

$$\begin{aligned} 2x_1 + x_2 + x_3 &= 2 \\ -x_1 + 2x_2 &= 1 \\ x_1 - x_2 + 2x_3 &= -2. \end{aligned}$$

**SOLUTION:** We have already seen in Example 1.1.3 that the solution of this system is given by  $x_1 = x_2 = 1$ ,  $x_3 = -1$ . Let us find that solution using Gauss operations with back

substitution. First transform the augmented matrix of the system into upper triangular form:

$$\begin{aligned} \left[ \begin{array}{ccc|c} 2 & 1 & 1 & 2 \\ -1 & 2 & 0 & 1 \\ 1 & -1 & 2 & -2 \end{array} \right] &\rightarrow \left[ \begin{array}{ccc|c} 1 & -1 & 2 & -2 \\ -1 & 2 & 0 & 1 \\ 2 & 1 & 1 & 2 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & -1 & 2 & -2 \\ 0 & 1 & 2 & -1 \\ 0 & 3 & -3 & 6 \end{array} \right] \\ &\rightarrow \left[ \begin{array}{ccc|c} 1 & -1 & 2 & -2 \\ 0 & 1 & 2 & -1 \\ 0 & 0 & -9 & 9 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & -1 & 2 & -2 \\ 0 & 1 & 2 & -1 \\ 0 & 0 & 1 & -1 \end{array} \right]. \end{aligned}$$

We now write the linear system corresponding to the last augmented matrix,

$$\begin{aligned} 2x_1 - x_2 + 2x_3 &= -2 \\ x_2 + 2x_3 &= -1 \\ x_3 &= -1. \end{aligned}$$

We now use the back substitution to obtain the solution, that is, introduce  $x_3 = -1$  into the second equation, which gives us  $x_2 = 1$ . Then, substitute  $x_3 = -1$  and  $x_2 = 1$  into the first equation, which gives us  $x_1 = 1$ .  $\triangleleft$

We now use the *Gauss-Jordan Method* to find solutions to square  $n \times n$  linear systems. This is a minor modification of the Gauss operations and back substitution method. The difference is that now we do not stop doing Gauss operations when the augmented matrix becomes upper triangular. We keep doing Gauss operations in order to make zeros above the diagonal. Then, back substitution will no be needed to find the solution of the linear system. The solution will be given explicitly at the end of the procedure.

**EXAMPLE 1.2.7:** Use the Gauss-Jordan method to solve the same  $3 \times 3$  linear system as in Example 1.2.6, that is,

$$\begin{aligned} 2x_1 + x_2 + x_3 &= 2 \\ -x_1 + 2x_2 &= 1 \\ x_1 - x_2 + 2x_3 &= -2. \end{aligned}$$

**SOLUTION:** In Example 1.2.6 we performed Gauss operations on the augmented matrix of the system above until we obtained an upper triangular matrix, that is,

$$\left[ \begin{array}{ccc|c} 2 & 1 & 1 & 2 \\ -1 & 2 & 0 & 1 \\ 1 & -1 & 2 & -2 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & -1 & 2 & -2 \\ 0 & 1 & 2 & -1 \\ 0 & 0 & 1 & -1 \end{array} \right].$$

The idea now is to continue with Gauss operations, as follows:

$$\left[ \begin{array}{ccc|c} 1 & -1 & 2 & -2 \\ 0 & 1 & 2 & -1 \\ 0 & 0 & 1 & -1 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & 0 & 4 & -3 \\ 0 & 1 & 2 & -1 \\ 0 & 0 & 1 & -1 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & -1 \end{array} \right] \Rightarrow \begin{cases} x_1 = 1, \\ x_2 = 1, \\ x_3 = -1. \end{cases}$$

In the last step we do not need to do back substitution to compute the solution. It is obtained without doing any further algebraic operation.  $\triangleleft$

**Further reading.** Almost every linear algebra book describes the Gauss-Jordan method. See Section 1.2 in Lay's book [2] for a summary of echelon forms and the Gauss-Jordan method. See Sections 1.2 and 1.3 in Meyer's book [3] for more details on Gauss elimination operations and back substitution. Also see Section 1.3 in Strang's book [4].

## 1.2.4. Exercises.

**1.2.1.-** Use Gauss operations and back substitution to find the solution of the  $3 \times 3$  linear system

$$\begin{aligned}x_1 + x_2 + x_3 &= 1, \\x_1 + 2x_2 + 2x_3 &= 1, \\x_1 + 2x_2 + 3x_3 &= 1.\end{aligned}$$

**1.2.2.-** Use Gauss operations and back substitution to find the solution of the  $3 \times 3$  linear system

$$\begin{aligned}2x_1 - 3x_2 &= 3, \\4x_1 - 5x_2 + x_3 &= 7, \\2x_1 - x_2 - 3x_3 &= 5.\end{aligned}$$

**1.2.3.-** Find the solution of the following linear system with Gauss-Jordan's method,

$$\begin{aligned}4x_2 - 3x_3 &= 3, \\-x_1 + 7x_2 - 5x_3 &= 4, \\-x_1 + 8x_2 - 6x_3 &= 5.\end{aligned}$$

**1.2.4.-** Find the solutions to the following two linear systems, which have the same matrix of coefficient  $A$  but different source vectors  $\mathbf{b}_1$  and  $\mathbf{b}_2$ , given respectively by,

$$\begin{aligned}4x_1 - 8x_2 + 5x_3 &= 1, & 4x_1 - 8x_2 + 5x_3 &= 0, \\4x_1 - 7x_2 + 4x_3 &= 0, & 4x_1 - 7x_2 + 4x_3 &= 1, \\3x_1 - 4x_2 + 2x_3 &= 0, & 3x_1 - 4x_2 + 2x_3 &= 0.\end{aligned}$$

Solve these two systems at one time using the Gauss-Jordan method on an augmented matrix of the form  $[A|\mathbf{b}_1|\mathbf{b}_2]$ .

**1.2.5.-** Use the Gauss-Jordan method to solve the following three linear systems at the same time,

$$\begin{aligned}2x_1 - x_2 &= 1, & = 0, & = 0, \\-x_1 + 2x_2 - x_3 &= 0, & = 1, & = 0, \\-x_2 + x_3 &= 0, & = 0, & = 1.\end{aligned}$$

## 1.3. ECHELON FORMS

**1.3.1. Echelon and reduced echelon forms.** The Gauss-Jordan method is a procedure that uses Gauss operations to transform the augmented matrix of an  $m \times n$  linear system into the augmented matrix of an equivalent system whose solutions can be found without performing further algebraic operations. A precise way to define the notion of a linear system with solutions that can be found without doing further algebraic operations is captured in the notions of echelon form and reduced echelon form of its augmented matrix.

**Definition 1.3.1.** An  $m \times n$  matrix is in **echelon form** iff the following conditions hold:

- (i) The zero rows are located at the bottom rows of the matrix;
- (ii) The first non-zero coefficient on a row is always to the right of the first non-zero coefficient of the row above it.

The **pivot** coefficient is the first non-zero coefficient on every non-zero row in a matrix in echelon form.

**EXAMPLE 1.3.1:** The  $6 \times 8$ ,  $3 \times 5$  and  $3 \times 3$  matrices given below are in echelon form, where pivots are highlighted and the \* means any non-zero number.

$$\begin{bmatrix} * & * & * & * & * & * & * & * \\ 0 & 0 & * & * & * & * & * & * \\ 0 & 0 & 0 & * & * & * & * & * \\ 0 & 0 & 0 & 0 & 0 & 0 & * & * \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}, \quad \begin{bmatrix} * & * & * & * & * \\ 0 & 0 & * & * & * \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}, \quad \begin{bmatrix} * & * & * \\ 0 & * & * \\ 0 & 0 & * \end{bmatrix}.$$

&lt;

**EXAMPLE 1.3.2:** The following matrices are in echelon form, with pivots highlighted:

$$\begin{bmatrix} \mathbf{1} & 3 \\ 0 & \mathbf{1} \end{bmatrix}, \quad \begin{bmatrix} \mathbf{2} & 3 & 2 \\ 0 & \mathbf{4} & -2 \end{bmatrix}, \quad \begin{bmatrix} \mathbf{2} & 1 & 1 \\ 0 & \mathbf{3} & 4 \\ 0 & 0 & 0 \end{bmatrix}.$$

&lt;

**Definition 1.3.2.** An  $m \times n$  matrix is in **reduced echelon form** iff the matrix is in echelon form and the following two conditions hold:

- (i) The pivot coefficient is equal to 1;
- (ii) The pivot coefficient is the only non-zero coefficient in that column.

We denote by  $E_A$  a reduced echelon form of a matrix  $A$ .

**EXAMPLE 1.3.3:** The  $6 \times 8$ ,  $3 \times 5$  and  $3 \times 3$  matrices given below are in echelon form, where pivots are highlighted and the \* means any non-zero number.

$$\begin{bmatrix} \mathbf{1} & * & 0 & 0 & * & * & 0 & * \\ 0 & 0 & \mathbf{1} & 0 & * & * & 0 & * \\ 0 & 0 & 0 & \mathbf{1} & * & * & 0 & * \\ 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{1} & * \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}, \quad \begin{bmatrix} \mathbf{1} & * & 0 & * & * \\ 0 & 0 & \mathbf{1} & * & * \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}, \quad \begin{bmatrix} \mathbf{1} & 0 & 0 \\ 0 & \mathbf{1} & 0 \\ 0 & 0 & \mathbf{1} \end{bmatrix}.$$

&lt;

**EXAMPLE 1.3.4:** And the following matrices are not only in echelon form but also in reduced echelon form; again, pivot coefficients are highlighted:

$$\begin{bmatrix} \mathbf{1} & 0 \\ 0 & \mathbf{1} \end{bmatrix}, \quad \begin{bmatrix} \mathbf{1} & 0 & 4 \\ 0 & \mathbf{1} & 5 \end{bmatrix}, \quad \begin{bmatrix} \mathbf{1} & 0 & 0 \\ 0 & \mathbf{1} & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

◁

Summarizing, the Gauss-Jordan method uses Gauss operations to transform the augmented matrix of a linear system into reduced echelon form. Then, the solutions of the linear system can be obtained without doing further algebraic operations.

**EXAMPLE 1.3.5:** Consider a  $3 \times 3$  linear system with augmented matrix having a reduced echelon form given below. Then, the solution of this linear system is simple to obtain:

$$\left[ \begin{array}{ccc|c} 1 & 0 & 0 & -1 \\ 0 & 1 & 0 & 3 \\ 0 & 0 & 1 & 2 \end{array} \right] \Rightarrow \begin{cases} x_1 = -1 \\ x_2 = 3 \\ x_3 = 2. \end{cases}$$

◁

We use the notation  $E_A$  for a reduced echelon form of an  $m \times n$  matrix  $A$ , since a reduced echelon form of a matrix has an important property: It is unique. We state this property in Proposition 1.3.3 below. Since the proof is somehow involved, after the proof we show a different proof for the particular case of  $2 \times 2$  matrices.

**Proposition 1.3.3.** *The reduced echelon  $E_A$  form of an  $m \times n$  matrix  $A$  is unique.*

**Proof of Proposition 1.3.3:** We assume that there exists an  $m \times n$  matrix  $A$  having two reduced echelon forms  $E_A$  and  $\hat{E}_A$ , and we will show that  $E_A = \hat{E}_A$ . Introduce the column vector notation for matrices (see Sect. 1.4) and for vectors (see Sect. 1.4), respectively,

$$A = [A_1, \dots, A_n], \quad E_A = [E_1, \dots, E_n], \quad \hat{E}_A = [\hat{E}_1, \dots, \hat{E}_n], \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}.$$

If the reduced echelon forms  $E_A$  and  $\hat{E}_A$  were different, then they should start differing at a particular column, say column  $i$ , with  $1 \leq i \leq n$ . That is, there would be vectors  $E_i$  and  $\hat{E}_i$  satisfying  $E_i \neq \hat{E}_i$  and  $E_j = \hat{E}_j$  for  $j = 1, \dots, (i-1)$ . The proof of Proposition 1.3.3 reduces to show that this is not the case, by studying the following two cases:

- (a) The vector of one of the reduced echelon form matrices, say the column vector  $E_i$  of  $E_A$ , does not contain a pivot coefficient.
- (b) Both vectors  $E_i$  and  $\hat{E}_i$  do contain a pivot coefficient.

Consider the case in (a). The definition of reduced echelon form implies that the vector  $E_i = [e_j]$  has non-zero components among the first  $k$  components, with  $0 \leq k < i$ , hence,  $e_j = 0$  for  $j = (k+1), \dots, n$ . Furthermore, the definition of reduced echelon form also implies that there are  $k$  columns to the left of vector  $E_i$  with pivot coefficients. Let us denote these columns by  $E_{p_j}$ , for  $j = 1, \dots, k$ . We use a new subindex  $p_j$ , since the  $k$  columns need not be the first  $k$  columns. This information about the reduced echelon form matrix  $E_A$  can be translated into a vector  $x$  solution of the linear system with augmented matrix  $[E_A|0]$ . This solution is the vector  $x = [x_j]$  with non-zero components  $x_{p_j} = e_j$  for  $j = 1, \dots, k$ , and  $x_i = -1$ , while the rest of the components vanish, since the following equations hold,

$$e_1 E_{p_1} + \dots + e_k E_{p_k} + (-1)E_i = 0 \quad \Leftrightarrow \quad x_{p_1} E_{p_1} + \dots + x_{p_k} E_{p_k} + x_i E_i = 0.$$

Recall now that the solution of a linear system is not changed when Gauss operations are performed on the coefficient matrix. Since Gauss operations change  $E_A \rightarrow A \rightarrow \hat{E}_A$ , the same vector  $x$  above is solution of the linear system with augmented matrix  $[\hat{E}_A|0]$ , that is,

$$x_{p_1}\hat{E}_{p_1} + \cdots + x_{p_k}\hat{E}_{p_k} + x_i\hat{E}_i = 0 \quad \Leftrightarrow \quad e_1\hat{E}_{p_1} + \cdots + e_k\hat{E}_{p_k} + (-1)\hat{E}_i = 0.$$

Since  $E_j = \hat{E}_j$  for  $j = 1, \dots, k$ , the second equation above implies that  $E_i = \hat{E}_i$  also holds. We conclude that the first vector in  $E_A$  that differs from a vector in  $\hat{E}_A$  cannot be a vector without pivot coefficients.

Consider the case in (b), that is, if  $E_i$  and  $\hat{E}_i$  are the first vectors from  $E_A$  and  $\hat{E}_A$ , respectively, which are different, then both vectors have pivot coefficients. (If only one of these vectors had a pivot coefficient, then the argument in case (a) on the vector without pivot coefficient would imply that the first vector could not have a pivot coefficient.) Suppose that the vector  $E_i = [e_j]$  has the pivot coefficient at the position  $k_0$ , that is,  $e_j = 0$  for  $j \neq k_0$  and  $e_{k_0} = 1$ . Similarly, suppose that the vector  $\hat{E}_i = [\hat{e}_j]$  has the pivot coefficient at the position  $k_1$ , that is,  $\hat{e}_j = 0$  for  $j \neq k_1$  and  $\hat{e}_{k_1} = 1$ . By definition of reduced echelon form all rows in  $E_A$  above the  $k_0$  row have pivots columns to the left of  $E_i$ . This statement also applies to matrix  $\hat{E}_A$ , since  $\hat{E}_i$  is the first column that differs from  $E_A$ . Therefore  $k_0 \leq k_1$ . The exactly same argument also applies to matrix  $\hat{E}_A$  concluding that  $k_1 \leq k_0$ . Therefore,  $k_0 = k_1$ , and so,  $E_i = \hat{E}_i$ . We conclude that the first vector in  $E_A$  that differs from a vector in  $\hat{E}_A$  cannot be a vector with pivot a coefficient. Therefore, parts (a) and (b) establish the Proposition.  $\square$

The proof above is not straightforward to understand, so it might be a good idea to present a simpler proof for the particular case of  $2 \times 2$  matrices.

**Alternative proof of Proposition 1.3.3 in the case of  $2 \times 2$  matrices:** We start recalling once again the following observation that holds for any  $m \times n$  matrix  $A$ : Since Gauss operations do not change the solutions of the homogeneous system  $[A|0]$ , if a matrix  $A$  has two different reduced echelon forms,  $E_A$  and  $\tilde{E}_A$ , then the set of solutions of the systems  $[E_A|0]$  and  $[\tilde{E}_A|0]$  must coincide with the solutions of the system  $[A|0]$ . What we are going to show now is the following: All possible  $2 \times 2$  reduced echelon form matrices  $E$  have different solutions to the homogeneous equation  $[E|0]$ . This property then establishes that every  $2 \times 2$  matrix  $A$  has a unique reduced echelon form.

Given any  $2 \times 2$  matrix  $A$ , all possible reduced echelon forms are the following:

$$E_A = \begin{bmatrix} 1 & c \\ 0 & 0 \end{bmatrix}, \quad c \in \mathbb{R}, \quad E_A = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}, \quad E_A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

We claim that all these matrices determine different sets of solutions to the equation  $[E_A|0]$ . In the first case, the set of solutions are lines given by

$$x_1 = -cx_2, \quad c \in \mathbb{R}. \tag{1.15}$$

In the second case, the set of solutions is a single line given by

$$x_2 = 0, \quad x_1 \in \mathbb{R}.$$

Notice that this line does not belong to the set given in Eq. (1.15). In the third case, the solution is a single point

$$x_1 = 0, \quad x_2 = 0.$$

Since all these sets of solutions are different, and only one corresponds to the equation  $[A|0]$ , then there is only one reduced echelon form  $E_A$  for matrix  $A$ . This establishes the Proposition for  $2 \times 2$  matrices.  $\square$



While the reduced echelon form of a matrix is unique, a matrix may have many different echelon forms. Also notice that given a matrix  $A$ , there are many different sequences of Gauss operations that produce the reduced echelon form  $E_A$ , as it is sketched in Fig. 13.

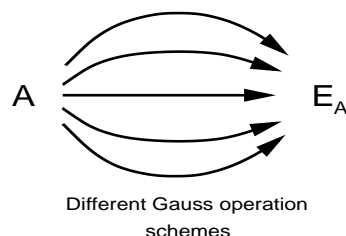


FIGURE 13. The reduced echelon form of a matrix can be obtained in many different ways.

**EXAMPLE 1.3.6:** Use two different sequences of Gauss operations to find the reduced echelon form of matrix

$$A = \begin{bmatrix} 2 & 4 & 10 \\ 1 & 3 & 7 \end{bmatrix}$$

**SOLUTION:** Here are two different sequences of Gauss operations to find  $E_A$ :

$$A = \begin{bmatrix} 2 & 4 & 10 \\ 1 & 3 & 7 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 5 \\ 1 & 3 & 7 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 5 \\ 0 & 1 & 2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 2 \end{bmatrix} = E_A,$$

$$A = \begin{bmatrix} 2 & 4 & 10 \\ 1 & 3 & 7 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 3 & 7 \\ 2 & 4 & 10 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 5 \\ 0 & -2 & -4 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 5 \\ 0 & 1 & 2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 2 \end{bmatrix} = E_A.$$

The matrix  $E_A$  is the same using the two Gauss operation sequences above.  $\triangleleft$

**1.3.2. The rank of a matrix.** Since the reduced echelon form  $E_A$  of an  $m \times n$  matrix  $A$  is unique, any property of matrix  $E_A$  is indeed a property of matrix  $A$ . In particular, the uniqueness of  $E_A$  implies that the number of pivots of  $E_A$  is also unique. This number is a property of matrix  $A$ , and we give that number a name, since it will be important later on.

**Definition 1.3.4.** The **rank** of an  $m \times n$  matrix  $A$ , denoted as  $\text{rank}(A)$ , is the number of pivots in its reduced echelon form  $E_A$ .

**EXAMPLE 1.3.7:** Find the rank of the coefficient matrix and all solutions of the linear system

$$\begin{aligned} x_1 + 3x_3 &= -1 \\ x_2 - 2x_3 &= 3 \\ 2x_1 + x_2 + 4x_3 &= 1. \end{aligned}$$

**SOLUTION:** The  $3 \times 3$  system has an augmented matrix with reduced echelon form

$$\left[ \begin{array}{ccc|c} 1 & 0 & 3 & -1 \\ 0 & 1 & -2 & 3 \\ 2 & 1 & 4 & 1 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & 0 & 3 & -1 \\ 0 & 1 & -2 & 3 \\ 0 & 0 & 0 & 0 \end{array} \right]$$

Since the reduced echelon form of the coefficient matrix has two pivots, we obtain that **the coefficient matrix has rank two**. The solutions of the linear system can be obtained without

any further calculations, since the coefficient matrix is already in reduced echelon form. These solutions are the following,

$$\left[ \begin{array}{ccc|c} 1 & 0 & 3 & -1 \\ 0 & 1 & -2 & 3 \\ 0 & 0 & 0 & 0 \end{array} \right] \Rightarrow \begin{cases} x_1 = -1 - 3x_3, \\ x_2 = 3 + 2x_3, \\ x_3 : \text{free variable.} \end{cases}$$

We call  $x_3$  a free variable, since there is a solution for every value of this variable.  $\triangleleft$

**Definition 1.3.5.** A variable in a linear system is called a **free variable** iff there is a solution of the linear system for every value of that variable.

In Example 1.3.7 the free variable is  $x_3$ , since for every value of  $x_3$  the other two variables  $x_1$  and  $x_2$  are fixed by the system. Notice that *only the number of free variables is relevant*, but not which particular variable is the free one. In the following example we express the solutions in Example 1.3.7 in terms of  $x_1$  as a free variable.

**EXAMPLE 1.3.8:** Express the solutions in Example 1.3.7 in terms of  $x_1$  as a free variable.

**SOLUTION:** The solutions given in Example 1.3.7 can also be expressed as

$$\begin{aligned} x_1 &: \text{free variable,} \\ x_2 &= 3 + \frac{2}{3}(-1 - x_1), \\ x_3 &= \frac{1}{3}(-1 - x_1). \end{aligned}$$

$\triangleleft$

As the reader may have noticed there is a relation between the rank of the coefficient matrix and the number of free variables in the solution of the linear system.

**Theorem 1.3.6.** The number of free variables in the solutions of an  $m \times n$  consistent linear system with augmented matrix  $[A|b]$  is given by  $(n - \text{rank}(A))$ .

**Proof of Theorem 1.3.6:** The number of non-pivots columns in  $E_A$  is actually the number of variables not fixed by the linear system with augmented matrix  $[A|b]$ . The number of non-pivot columns is the total number of columns minus the pivot columns, that is,  $(n - \text{rank}(A))$ . This establishes the Theorem.  $\square$

We saw in Example 1.3.7 that the  $3 \times 3$  coefficient matrix has  $\text{rank}(A) = 2$ . Since the system has  $n = 3$  variables, we conclude, without actually computing the solution, that the solution has  $n - \text{rank}(A) = 3 - 2 = 1$  free variable. The free variable can be any variable in the system. The relevant information is that there is one free variable.

**EXAMPLE 1.3.9:** We now present three examples of consistent linear systems.

- (i) Show that the  $2 \times 2$  linear system below is consistent with coefficient matrix having rank one and solutions having one free variable,

$$\begin{aligned} 2x_1 - x_2 &= 1 \\ -\frac{1}{2}x_1 + \frac{1}{4}x_2 &= -\frac{1}{4}. \end{aligned}$$

**SOLUTION:** Gauss operations transform the augmented matrix of the system as follows,

$$\left[ \begin{array}{cc|c} 2 & -1 & 1 \\ -1/2 & 1/4 & -1/4 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 2 & -1 & 1 \\ 0 & 0 & 0 \end{array} \right].$$

The system is consistent, has rank one, and has a free variable, hence the system has infinitely many solutions.

- (ii) Show that the  $2 \times 3$  linear system below is consistent with coefficient matrix having rank two and solution having one free variable,

$$\begin{aligned}x_1 + 2x_2 + 3x_3 &= 1, \\3x_1 + 5x_2 + 2x_3 &= 2.\end{aligned}$$

**SOLUTION:** Perform the following Gauss operations on the augmented matrix,

$$\left[ \begin{array}{ccc|c} 1 & 2 & 3 & 1 \\ 3 & 5 & 2 & 2 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & 2 & 3 & 1 \\ 0 & -1 & -7 & -1 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & 0 & -11 & -1 \\ 0 & 1 & 7 & 1 \end{array} \right].$$

We see that the  $2 \times 3$  system has rank two, hence one free variable.

- (iii) Show that the  $3 \times 2$  linear system below is consistent with coefficient matrix having rank two and solution having no free variables,

$$\begin{aligned}x_1 + 3x_2 &= 2, \\2x_1 + 2x_2 &= 0, \\3x_1 + x_2 &= -2.\end{aligned}$$

**SOLUTION:** Perform the following Gauss operations on the augmented matrix,

$$\left[ \begin{array}{cc|c} 1 & 3 & 2 \\ 2 & 2 & 0 \\ 3 & 1 & -2 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 3 & 2 \\ 0 & -4 & -4 \\ 0 & -8 & -8 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 0 & -1 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{array} \right].$$

We see that the  $3 \times 2$  system has rank two, hence no free variables. ◁

**1.3.3. Inconsistent linear systems.** Not every linear system has solutions. We now use both the reduced echelon form of the coefficient matrix and of the augmented matrix to characterize inconsistent linear systems. We start with an example.

**EXAMPLE 1.3.10:** Show that the following  $2 \times 2$  linear system has no solutions,

$$2x_1 - x_2 = 0 \tag{1.16}$$

$$-\frac{1}{2}x_1 + \frac{1}{4}x_2 = -\frac{1}{4}. \tag{1.17}$$

**SOLUTION:** One way to see that there is no solution is the following: Multiplying the second equation by  $-4$  one obtains the equation

$$2x_1 - x_2 = 1,$$

whose solutions form a parallel line to the line given in Eq. (1.16). Therefore, the system in Eqs. (1.16)-(1.17) has no solution. A second way to see that the system above has no solution is using Gauss operations. The system above has augmented matrix

$$\left[ \begin{array}{cc|c} 2 & -1 & 0 \\ -\frac{1}{2} & \frac{1}{4} & -\frac{1}{4} \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 2 & -1 & 0 \\ 0 & 0 & 1 \end{array} \right].$$

The echelon form above corresponds to the linear system

$$\begin{aligned}2x_1 - x_2 &= 0 \\0 &= 1.\end{aligned}$$

The solutions of this second system coincide with the solutions of Eqs. (1.16)-(1.17). Since the second system has no solution, the system in Eqs. (1.16)-(1.17) has no solution. ◁

This example is a particular cases of the following result.

**Theorem 1.3.7.** *An  $m \times n$  linear system with augmented matrix  $[A|b]$  is inconsistent iff the reduced echelon form of its augmented matrix contains a row of the form  $[0, \dots, 0 | 1]$ ; equivalently  $\text{rank}(A) < \text{rank}(A|b)$ . Furthermore, a consistent system contains:*

- (i) *A unique solution iff it has no free variables; equivalently  $\text{rank}(A) = n$ .*
- (ii) *Infinitely many solutions iff it has at least one free variable; equivalently  $\text{rank}(A) < n$ .*

This Theorem says that a system is inconsistent iff its augmented matrix has a pivot in the source vector column, which is the last column in an augmented matrix. In that case the system includes an equation of the form  $0 = 1$ , which has no solution.

The idea of the proof of this Theorem is to study all possible reduced echelon forms  $E_A$  of an arbitrary matrix  $A$ , and then to study all possible augmented matrices  $[E_A|c]$ . One then concludes that there are three main cases, no solutions, unique solutions, or infinitely many solutions.

**Proof of Theorem 1.3.7:** We only give the proof in the case of  $3 \times 3$  linear systems. The reduced echelon form  $E_A$  of a  $3 \times 3$  matrix  $A$  determines a  $3 \times 4$  augmented matrix  $[E_A|c]$ . There are 14 possible forms for this matrix. We start with the case of three pivots:

$$\left[ \begin{array}{ccc|c} 1 & 0 & 0 & * \\ 0 & 1 & 0 & * \\ 0 & 0 & 1 & * \end{array} \right], \quad \left[ \begin{array}{ccc|c} 1 & 0 & * & * \\ 0 & 1 & * & * \\ 0 & 0 & 0 & 1 \end{array} \right], \quad \left[ \begin{array}{ccc|c} 1 & * & 0 & * \\ 0 & 0 & 1 & * \\ 0 & 0 & 0 & 1 \end{array} \right], \quad \left[ \begin{array}{ccc|c} 0 & 1 & 0 & * \\ 0 & 0 & 1 & * \\ 0 & 0 & 0 & 1 \end{array} \right].$$

In the first case we have a unique solution, and  $\text{rank}(A) = 3$ . The other three cases correspond to no solutions. We now continue with the case of two pivots, which contains six possibilities, with the first three given by

$$\left[ \begin{array}{ccc|c} 1 & 0 & * & * \\ 0 & 1 & * & * \\ 0 & 0 & 0 & 0 \end{array} \right], \quad \left[ \begin{array}{ccc|c} 1 & * & 0 & * \\ 0 & 0 & 1 & * \\ 0 & 0 & 0 & 0 \end{array} \right], \quad \left[ \begin{array}{ccc|c} 0 & 1 & 0 & * \\ 0 & 0 & 1 & * \\ 0 & 0 & 0 & 0 \end{array} \right],$$

which correspond to infinitely many solutions, with one free variable,  $\text{rank}(A) = 2$ . The other three possibilities are given by

$$\left[ \begin{array}{ccc|c} 1 & * & * & * \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{array} \right], \quad \left[ \begin{array}{ccc|c} 0 & 1 & * & * \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{array} \right], \quad \left[ \begin{array}{ccc|c} 0 & 0 & 1 & * \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{array} \right];$$

which correspond to no solution. Finally, we have the case of one pivot. It contains four possibilities, the first three of them are given by

$$\left[ \begin{array}{ccc|c} 1 & * & * & * \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right], \quad \left[ \begin{array}{ccc|c} 0 & 1 & * & * \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right], \quad \left[ \begin{array}{ccc|c} 0 & 0 & 1 & * \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right];$$

which correspond to infinitely many solutions, with two free variables and  $\text{rank}(A) = 1$ . The last possibility is the trivial case

$$\left[ \begin{array}{ccc|c} 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \right],$$

which has no solutions. This establishes the Theorem for  $3 \times 3$  linear systems. □

**EXAMPLE 1.3.11:** Show that the  $2 \times 2$  linear system below is inconsistent:

$$\begin{aligned} 2x_1 - x_2 &= 1 \\ -\frac{1}{2}x_1 + \frac{1}{4}x_2 &= -\frac{1}{4}. \end{aligned}$$

**SOLUTION:** The proof of the statement above is the following: Gauss-Jordan operations transform the augmented matrix of the system as follows,

$$\left[ \begin{array}{cc|c} 2 & -1 & 0 \\ -\frac{1}{2} & \frac{1}{4} & -\frac{1}{4} \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 2 & -1 & 0 \\ 0 & 0 & -\frac{1}{4} \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 2 & -1 & 0 \\ 0 & 0 & 1 \end{array} \right].$$

Since there is a line of the form  $[0, 0|1]$ , the system is inconsistent. We can also say that the coefficient matrix has rank one, but the definition of free variables does not apply, since the system is inconsistent.  $\triangleleft$

**EXAMPLE 1.3.12:** Find all numbers  $h$  and  $k$  such that the system below has only one, many, or no solutions,

$$\begin{aligned} x_1 + hx_2 &= 1 \\ x_1 + 2x_2 &= k. \end{aligned}$$

**SOLUTION:** Start finding the associated augmented matrix and reducing it into echelon form,

$$\left[ \begin{array}{cc|c} 1 & h & 1 \\ 1 & 2 & k \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & h & 1 \\ 0 & 2-h & k-1 \end{array} \right].$$

Suppose  $h \neq 2$ , for example set  $h = 1$ , then

$$\left[ \begin{array}{cc|c} 1 & 1 & 1 \\ 0 & 1 & k-1 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 0 & 2-k \\ 0 & 1 & k-1 \end{array} \right],$$

so the system has a unique solution for all values of  $k$ . (The same conclusion holds if one sets  $h$  to any number different of 2.) Suppose now that  $h = 2$ , then,

$$\left[ \begin{array}{cc|c} 1 & 2 & 1 \\ 0 & 0 & k-1 \end{array} \right].$$

If  $k = 1$  then

$$\left[ \begin{array}{cc|c} 1 & 2 & 1 \\ 0 & 0 & 0 \end{array} \right] \Rightarrow \begin{cases} x_1 = 1 - 2x_2, \\ x_2 : \text{free variable.} \end{cases}$$

so there are infinitely many solutions. If  $k \neq 1$ , the system is inconsistent, since

$$\left[ \begin{array}{cc|c} 1 & 2 & 1, \\ 0 & 0 & k-1 \neq 0. \end{array} \right]$$

Summarizing, for  $h \neq 2$  the system has a unique solution for every  $k$ . If  $h = 2$  and  $k = 1$  the system has infinitely many solutions, if  $h = 2$  and  $k \neq 1$  the system has no solution.  $\triangleleft$

**Further reading.** See Sections 2.1, 2.2 and 2.3 in Meyer's book [3] for a detailed discussion on echelon forms and rank, reduced echelon forms and inconsistent systems, respectively. Again, see Section 1.2 in Lay's book [2] for a summary of echelon forms and the Gauss-Jordan method. Section 1.3 in Strang's book [4] also helps.

## 1.3.4. Exercises.

- 1.3.1.- Find the rank and the pivot columns of the matrix

$$A = \begin{bmatrix} 1 & 2 & 1 & 1 \\ 2 & 4 & 2 & 2 \\ 3 & 6 & 3 & 4 \end{bmatrix}.$$

- 1.3.2.- Find all the solutions  $\mathbf{x}$  of the linear system  $A\mathbf{x} = \mathbf{0}$ , where the matrix  $A$  is given by

$$A = \begin{bmatrix} 1 & 2 & 1 & 3 \\ 2 & 1 & -1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix}$$

- 1.3.3.- Find all the solutions  $\mathbf{x}$  of the linear system  $A\mathbf{x} = \mathbf{b}$ , where the matrix  $A$  and the vector  $\mathbf{b}$  are given by

$$A = \begin{bmatrix} 1 & 2 & 4 \\ 2 & -1 & -7 \\ 3 & 2 & 0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 3 \\ -4 \\ 1 \end{bmatrix}.$$

- 1.3.4.- Construct a  $3 \times 4$  matrix  $A$  and 3-vectors  $\mathbf{b}$ ,  $\mathbf{c}$ , such that  $[A|\mathbf{b}]$  is the augmented matrix of a consistent system and  $[A|\mathbf{c}]$  is the augmented matrix of an inconsistent system.

- 1.3.5.- Let  $A$  be an  $m \times n$  matrix having  $\text{rank}(A) = m$ . Explain why the system with augmented matrix  $[A|\mathbf{b}]$  is consistent for every  $m$ -vector  $\mathbf{b}$ .

- 1.3.6.- Consider the following system of linear equations, where  $k$  represents any real number,

$$\begin{aligned} 2x_1 + kx_2 &= 1, \\ -4x_1 + 2x_2 &= 4. \end{aligned}$$

- (a) Find all possible values of the number  $k$  such that the system above is **inconsistent**.
- (b) Set  $k = 2$ . In this case, find the solution  $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ .

## 1.4. NON-HOMOGENEOUS EQUATIONS

In this Section we introduce the definitions of homogeneous and non-homogeneous linear systems and discuss few relations between their solutions. We start, however, introducing new notation. We first define a vector form for the solution of a linear system and we then introduce a matrix-vector product in order to express a linear system in a compact way. One advantage of this notation appears in this Section: It is simple to realize that solutions of a non-homogeneous linear system are translations of solutions of the associated homogeneous system. A somehow deeper insight on the relation between matrices and vectors can be obtained from the matrix-vector product, which will be discussed in the next Chapter. Here we can summarize this insight as follows: A matrix can be thought as a function acting on the space of vectors.

**1.4.1. Matrix-vector product.** We start generalizing the vector notation to include the unknowns of a linear system. First, recall that given an  $m \times n$  linear system

$$A_{11}x_1 + \cdots + A_{1n}x_n = b_1, \quad (1.18)$$

$$\vdots$$

$$A_{m1}x_1 + \cdots + A_{mn}x_n = b_m, \quad (1.19)$$

we have introduced the  $m \times n$  coefficients matrix  $\mathbf{A} = [A_{ij}]$  and the source  $m$ -vector  $\mathbf{b} = [b_i]$ , where  $i = 1, \dots, m$  and  $j = 1, \dots, n$ . Now, introduce the *unknown  $n$ -vector*  $\mathbf{x} = [x_j]$ . The matrix  $\mathbf{A}$  and the vectors  $\mathbf{b}$ ,  $\mathbf{x}$  can be used to express a linear system if we introduce an operation between a matrix and a vector. The result of this operation must be the left-hand side in Eqs. (1.18)-(1.19).

**Definition 1.4.1.** The *matrix-vector product* of an  $m \times n$  matrix  $\mathbf{A} = [A_{ij}]$  and an  $n$ -vector  $\mathbf{x} = [x_j]$ , where  $i = 1, \dots, m$  and  $j = 1, \dots, n$ , is the  $m$ -vector

$$\mathbf{Ax} = \begin{bmatrix} A_{11} & \cdots & A_{1n} \\ \vdots & & \vdots \\ A_{n1} & \cdots & A_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + \cdots + A_{1n}x_n \\ \vdots \\ A_{m1}x_1 + \cdots + A_{mn}x_n \end{bmatrix}.$$

We now use the matrix, vectors above together with the matrix-vector product to express a linear system in a compact notation.

**Theorem 1.4.2.** Given the algebraic linear system in Eqs. (1.18)-(1.19), introduce the coefficient matrix  $\mathbf{A}$ , the unknown vector  $\mathbf{x}$ , and the source vector  $\mathbf{b}$ , as follows,

$$\mathbf{A} = \begin{bmatrix} A_{11} & \cdots & A_{1n} \\ \vdots & & \vdots \\ A_{n1} & \cdots & A_{nn} \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}.$$

Then, the algebraic linear system can be written as

$$\mathbf{Ax} = \mathbf{b}.$$

**Proof of Theorem 1.4.2:** From the definition of the matrix-vector product we have that

$$\mathbf{Ax} = \begin{bmatrix} A_{11} & \cdots & A_{1n} \\ \vdots & & \vdots \\ A_{n1} & \cdots & A_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + \cdots + A_{1n}x_n \\ \vdots \\ A_{n1}x_1 + \cdots + A_{nn}x_n \end{bmatrix}.$$

Then, we conclude that

$$\left. \begin{array}{l} A_{11}x_1 + \cdots + A_{1n}x_n = b_1, \\ \vdots \\ A_{n1}x_1 + \cdots + A_{nn}x_n = b_n, \end{array} \right\} \Leftrightarrow \begin{bmatrix} A_{11}x_1 + \cdots + a_{1n}x_n \\ \vdots \\ A_{n1}x_1 + \cdots + A_{1n}x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} \Leftrightarrow \mathbf{Ax} = \mathbf{b}.$$

□

Therefore, an  $m \times n$  linear system with augmented matrix  $[\mathbf{A}|\mathbf{b}]$  can now be expressed as using the matrix-vector product as  $\mathbf{Ax} = \mathbf{b}$ , where  $\mathbf{x}$  is the variable  $n$ -vector and  $\mathbf{b}$  is the usual source  $m$ -vector.

**EXAMPLE 1.4.1:** Use the matrix-vector product to express the linear system

$$\begin{aligned} 2x_1 - x_2 &= 3, \\ -x_1 + 2x_2 &= 0. \end{aligned}$$

**SOLUTION:** The matrix of coefficient  $\mathbf{A}$ , the variable vector  $\mathbf{x}$  and the source vector  $\mathbf{b}$  are respectively given by

$$\mathbf{A} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 3 \\ 0 \end{bmatrix}.$$

The coefficient matrix can be written as

$$\mathbf{A} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} = [\mathbf{A}_{:1}, \mathbf{A}_{:2}], \quad \mathbf{A}_{:1} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}, \quad \mathbf{A}_{:2} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}.$$

The linear system above can be written in the compact way  $\mathbf{Ax} = \mathbf{b}$ , that is,

$$\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 3 \\ 0 \end{bmatrix}.$$

We now verify that the notation above actually represents the original linear system:

$$\begin{bmatrix} 3 \\ 0 \end{bmatrix} = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2 \\ -1 \end{bmatrix} x_1 + \begin{bmatrix} -1 \\ 2 \end{bmatrix} x_2 = \begin{bmatrix} 2x_1 \\ -x_1 \end{bmatrix} + \begin{bmatrix} -x_2 \\ 2x_2 \end{bmatrix} = \begin{bmatrix} 2x_1 - x_2 \\ -x_1 + 2x_2 \end{bmatrix},$$

which indeed is the original linear system. ◀

**EXAMPLE 1.4.2:** Use the matrix-vector product to express the  $2 \times 3$  linear system

$$\begin{aligned} 2x_1 - 2x_2 + 4x_3 &= 6, \\ x_1 + 3x_2 + 2x_3 &= 10. \end{aligned}$$

**SOLUTION:** The matrix of coefficients  $\mathbf{A}$ , variable vector  $\mathbf{x}$  and source vector  $\mathbf{b}$  are given by

$$\mathbf{A} = \begin{bmatrix} 2 & -2 & 4 \\ 1 & 3 & 2 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 6 \\ 10 \end{bmatrix}.$$

Therefore, the linear system above can be written as  $\mathbf{Ax} = \mathbf{b}$ , that is,

$$\begin{bmatrix} 2 & -2 & 4 \\ 1 & 3 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 6 \\ 10 \end{bmatrix}.$$

We now verify that this notation reproduces the linear system above:

$$\begin{bmatrix} 6 \\ 10 \end{bmatrix} = \begin{bmatrix} 2 & -2 & 4 \\ 1 & 3 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix} x_1 + \begin{bmatrix} -2 \\ 3 \end{bmatrix} x_2 + \begin{bmatrix} 4 \\ 2 \end{bmatrix} x_3 = \begin{bmatrix} 2x_1 - 2x_2 + 4x_3 \\ x_1 + 3x_2 + 2x_3 \end{bmatrix},$$



which indeed is the linear system above.  $\triangleleft$

**EXAMPLE 1.4.3:** Use the matrix-vector product to express the  $3 \times 2$  linear system

$$\begin{aligned}x_1 - x_2 &= 0, \\ -x_1 + x_2 &= 2, \\ x_1 + x_2 &= 0.\end{aligned}$$

**SOLUTION:** Using the matrix-vector product we get

$$\begin{bmatrix} 1 & -1 \\ -1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 2 \\ 0 \end{bmatrix}.$$

We now verify that the notation above actually represents the original linear system:

$$\begin{bmatrix} 0 \\ 2 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} x_1 + \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} x_2 = \begin{bmatrix} x_1 - x_2 \\ -x_1 + x_2 \\ x_1 + x_2 \end{bmatrix},$$

which is indeed the linear system above.  $\triangleleft$

**1.4.2. Linearity of matrix-vector product.** Introduce a column-vector notation for the coefficient matrix, that is,

$$\mathbf{A} = \begin{bmatrix} A_{11} & \cdots & A_{1n} \\ \vdots & & \vdots \\ A_{m1} & \cdots & A_{mn} \end{bmatrix} = [\mathbf{A}_{:1}, \cdots, \mathbf{A}_{:n}],$$

where  $\mathbf{A}_{:j}$  is the  $j$ -th column of the coefficient matrix  $\mathbf{A}$ . Using this notation we can rewrite a matrix-vector product as a linear combination of the matrix column vectors  $\mathbf{A}_{:j}$ . This is our first result.

**Theorem 1.4.3.** *The matrix-vector product of an  $m \times n$  matrix  $\mathbf{A} = [\mathbf{A}_{:1}, \cdots, \mathbf{A}_{:n}]$  and an  $n$ -vector  $\mathbf{x} = [x_i]$  can be written as follows,*

$$\mathbf{Ax} = \mathbf{A}_{:1}x_1 + \cdots + \mathbf{A}_{:n}x_n.$$

**Proof of Theorem 1.4.3:** Write down the matrix-vector product of the  $m \times n$  matrix  $\mathbf{A} = [A_{ij}]$  and the  $n$ -vector  $\mathbf{x} = [x_j]$ , that is,

$$\mathbf{Ax} = \begin{bmatrix} A_{11} & \cdots & A_{1n} \\ \vdots & & \vdots \\ A_{m1} & \cdots & A_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} A_{11}x_1 + \cdots + A_{1n}x_n \\ \vdots \\ A_{m1}x_1 + \cdots + A_{mn}x_n \end{bmatrix} = \begin{bmatrix} A_{11}x_1 \\ \vdots \\ A_{m1}x_1 \end{bmatrix} + \cdots + \begin{bmatrix} A_{1n}x_n \\ \vdots \\ A_{mn}x_n \end{bmatrix}.$$

The expression on the far right can be rewritten as

$$\mathbf{Ax} = \begin{bmatrix} A_{11} \\ \vdots \\ A_{m1} \end{bmatrix} x_1 + \cdots + \begin{bmatrix} A_{1n} \\ \vdots \\ A_{mn} \end{bmatrix} x_n \quad \Rightarrow \quad \mathbf{Ax} = \mathbf{A}_{:1}x_1 + \cdots + \mathbf{A}_{:n}x_n.$$

This establishes the Theorem.  $\square$

We are now ready to show that the matrix-vector product has an important property: It preserves the linear combination of vectors. We say that the matrix-vector product is a *linear operation*. This property is summarized below.

**Theorem 1.4.4.** For every  $m \times n$  matrix  $A$ , every  $n$ -vectors  $x$ ,  $y$  and for every numbers  $a$ ,  $b$ , the matrix-vector product satisfies that

$$A(ax + by) = aAx + bAy.$$

In words, the Theorem says that the matrix-vector product of a linear combination of vectors is the linear combination of the matrix-vector products. The expression above contains the particular cases  $a = b = 1$  and  $b = 0$ , which are respectively given by

$$A(x + y) = Ax + Ay, \quad A(ax) = aAx.$$

**Proof of Theorem 1.4.4:** From the definition of the matrix-vector product we see that:

$$A(ax + by) = [A_{:1}, \dots, A_{:n}] \begin{bmatrix} ax_1 + by_1 \\ \vdots \\ ax_n + by_n \end{bmatrix} = A_{:1}(ax_1 + by_1) + \dots + A_{:n}(ax_n + by_n).$$

Reorder terms in the expression above to get,

$$A(ax + by) = a(A_{:1}x_1 + \dots + A_{:n}x_n) + b(A_{:1}y_1 + \dots + A_{:n}y_n) = aAx + bAy.$$

This establishes the Theorem.  $\square$

**1.4.3. Homogeneous linear systems.** All possible linear systems can be classified into two main classes, homogeneous and non-homogeneous, depending whether the source vector vanishes or not. This classification will be useful to express solutions of non-homogeneous systems in terms of solutions of the associated homogeneous system.

**Definition 1.4.5.** The  $m \times n$  linear system  $Ax = b$  is called **homogeneous** iff it holds that  $b = 0$ ; and is called **non-homogeneous** iff it holds that  $b \neq 0$ .

Every  $m \times n$  homogeneous linear system has at least one solution, given by  $x = 0$ , called the *trivial* solution of the homogeneous system. The following example shows that homogeneous linear systems can also have non-trivial solutions.

**EXAMPLE 1.4.4:** Find all solutions of the system  $Ax = 0$ , with coefficient matrix

$$A = \begin{bmatrix} -2 & 4 \\ 1 & -2 \end{bmatrix}.$$

**SOLUTION:** The linear system above can be written in the usual way as follows,

$$-2x_1 + 4x_2 = 0, \tag{1.20}$$

$$x_1 - 2x_2 = 0. \tag{1.21}$$

The solution of this system can be found performing the Gauss elimination operations

$$\begin{bmatrix} -2 & 4 \\ 1 & -2 \end{bmatrix} \rightarrow \begin{bmatrix} -2 & 4 \\ 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -2 \\ 0 & 0 \end{bmatrix} \Rightarrow \begin{cases} x_1 = 2x_2, \\ x_2 : \text{free variable.} \end{cases}$$

We see above that the coefficient matrix of this system has rank one, so the solutions have one free variable. The set of all solutions of the linear system above is given by

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2x_2 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix} x_2, \quad x_2 \in \mathbb{R}.$$

We conclude that the set of all solutions of this linear system can be identified with the set of points that belong to the line shown in Fig. 14.  $\triangleleft$

It is also possible that an homogeneous linear system has only the trivial solution, like the following example shows.

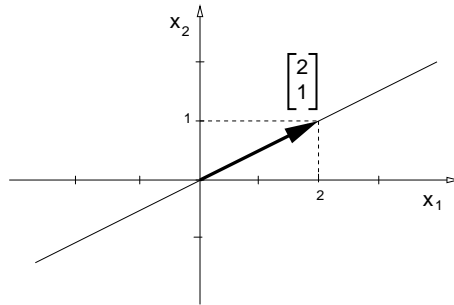


FIGURE 14. We plot the solutions of the homogeneous linear system given in Eqs. (1.20)-(1.21).

**EXAMPLE 1.4.5:** Find all solutions of the system  $Ax = 0$  with coefficient matrix

$$A = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}.$$

**SOLUTION:** The linear system above can be written in the usual way as follows,

$$\begin{aligned} 2x_1 - x_2 &= 0, \\ -x_1 + 2x_2 &= 0. \end{aligned}$$

The solutions of this system can be found performing the Gauss elimination operations

$$\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -2 \\ 2 & -1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -2 \\ 0 & 3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \Rightarrow \begin{cases} x_1 = 0, \\ x_2 = 0. \end{cases}$$

We see that the coefficient matrix of this system has rank two, so the solutions have no free variable. The solution is unique and is the trivial solution  $x = 0$ .  $\triangleleft$

Examples 1.4.4 and 1.4.5 are particular cases of the following statement: An  $m \times n$  homogeneous linear system has non-trivial solutions iff the system has at least one free variable. We show more examples of this statement.

**EXAMPLE 1.4.6:** Find all solutions of the  $2 \times 3$  homogeneous linear system  $Ax = 0$  with coefficient matrix

$$A = \begin{bmatrix} 2 & -2 & 4 \\ 1 & 3 & 2 \end{bmatrix}.$$

**SOLUTION:** The linear system above can be written in the usual way as follows,

$$2x_1 - 2x_2 + 4x_3 = 0, \tag{1.22}$$

$$x_1 + 3x_2 + 2x_3 = 0. \tag{1.23}$$

The solutions of system above can be obtained through the Gauss elimination operations

$$\begin{bmatrix} 1 & 3 & 2 \\ 2 & -2 & 4 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 3 & 2 \\ 0 & -8 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 2 \\ 0 & 1 & 0 \end{bmatrix} \Rightarrow \begin{cases} x_1 = -2x_3, \\ x_2 = 0, \\ x_3 : \text{free variable.} \end{cases}$$

The set of all solutions of the linear system above can be written in vector notation,

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -2x_3 \\ 0 \\ x_3 \end{bmatrix} = \begin{bmatrix} -2 \\ 0 \\ 1 \end{bmatrix} x_3, \quad x_3 \in \mathbb{R}.$$

In Fig. 15 we emphasize that the solution vector  $\mathbf{x}$  belongs to the space  $\mathbb{R}^3$ , while the column vectors of the coefficient matrix of this same system belong to the space  $\mathbb{R}^2$ .  $\triangleleft$

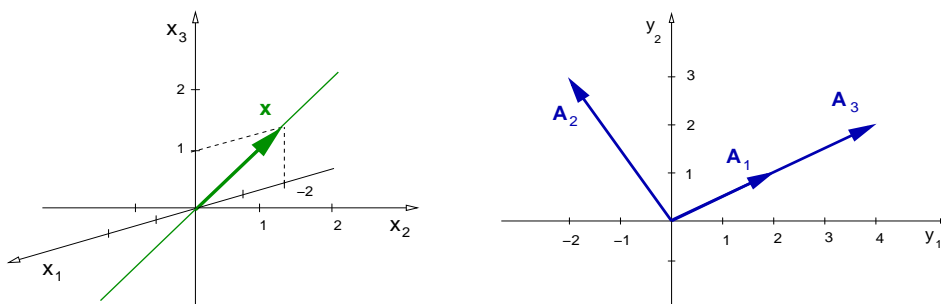


FIGURE 15. The picture on the left represents the solutions of the homogeneous linear system given in Eq. (1.22)-(1.23), which are 3-vectors, elements in  $\mathbb{R}^3$ . The picture on the right represents the column vectors of the coefficient matrix in this system which are 2-vectors, elements in  $\mathbb{R}^2$ .

**EXAMPLE 1.4.7:** Find all solutions to the linear system  $A\mathbf{x} = \mathbf{0}$ , with coefficient matrix

$$A = \begin{bmatrix} 1 & 3 & 4 \\ 2 & 6 & 8 \end{bmatrix}.$$

**SOLUTION:** We only need to use the Gauss-Jordan method on the coefficient matrix

$$\begin{bmatrix} 1 & 3 & 4 \\ 2 & 6 & 8 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 3 & 4 \\ 0 & 0 & 0 \end{bmatrix} \Rightarrow \begin{cases} x_1 = -3x_2 - 4x_3 \\ x_2, x_3 : \text{free variables.} \end{cases}$$

In this case the solution can be expressed in vector notation as follows

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -3x_2 - 4x_3 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -3 \\ 1 \\ 0 \end{bmatrix} x_2 + \begin{bmatrix} -4 \\ 0 \\ 1 \end{bmatrix} x_3.$$

We conclude that all solutions of the homogeneous linear system above are all possible linear combinations of the vectors

$$\mathbf{u}_1 = \begin{bmatrix} -3 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{u}_2 = \begin{bmatrix} -4 \\ 0 \\ 1 \end{bmatrix}.$$

$\triangleleft$

**1.4.4. The span of vector sets.** From Examples 1.4.4-1.4.7 above we see that solutions of homogeneous linear systems can be expressed as the sets of all possible linear combinations of particular vectors. Since this type of set will appear very often in our studies, it is convenient to give such sets a name. We introduce the span of a finite set of vectors as the infinite set of all possible linear combinations of this finite set of vectors.

**Definition 1.4.6.** The *span* of a finite set  $U = \{\mathbf{u}_1, \dots, \mathbf{u}_k\} \subset \mathbb{R}^n$ , with  $k \geq 1$ , denoted as  $\text{Span}(U)$ , is the set given by

$$\text{Span}(U) = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} = x_1\mathbf{u}_1 + \dots + x_n\mathbf{u}_n, \quad \forall x_1, \dots, x_n \in \mathbb{R}\}$$

Recall that the symbol “ $\forall$ ” means “for all”. Using this definition we express the solutions of  $Ax = 0$  in Example 1.4.7 as

$$x \in \text{Span}\left(\left\{u_1 = \begin{bmatrix} -3 \\ 1 \\ 0 \end{bmatrix}, u_2 = \begin{bmatrix} -4 \\ 0 \\ 1 \end{bmatrix}\right\}\right).$$

In this case, the set of all solutions forms a plane in  $\mathbb{R}^3$ , which contains the vectors  $u_1$  and  $u_2$ . In the case of Example 1.4.4 the solutions  $x$  belong to a line in  $\mathbb{R}^2$  given by

$$x \in \text{Span}\left(\left\{\begin{bmatrix} 2 \\ 1 \end{bmatrix}\right\}\right).$$

**EXAMPLE 1.4.8:** Find the set  $S = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 \\ 4 \end{bmatrix}\right\}\right) \subset \mathbb{R}^2$ .

**SOLUTION:** Since the vectors are not proportional to each other, the set of all linear combinations of these vectors is the whole plane. We conclude that  $S = \mathbb{R}^2$ .  $\triangleleft$

**EXAMPLE 1.4.9:** Find the set  $S = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} -1 \\ -2 \end{bmatrix}\right\}\right) \subset \mathbb{R}^2$ .

**SOLUTION:** Since the vectors lay on a line, the set of all linear combinations of these vectors also belongs to the same line. We conclude that  $S = \left\{a \begin{bmatrix} 1 \\ 2 \end{bmatrix}, a \in \mathbb{R}\right\}$ .  $\triangleleft$

**1.4.5. Non-homogeneous linear systems.** Spans of vector sets are not only useful to express solution sets to homogeneous equations, they are also useful to characterize when a non-homogeneous linear system is consistent.

**Theorem 1.4.7.** An  $m \times n$  linear system  $Ax = b$ , with coefficient matrix  $A = [A_{:1}, \dots, A_{:n}]$ , is consistent iff  $b \in \text{Span}(\{A_{:1}, \dots, A_{:n}\})$ .

In words, a non-homogeneous linear system is consistent iff the source vector belongs to the Span of the coefficient matrix column vectors.

**Proof of Theorem 1.4.7:**

( $\Rightarrow$ ) Suppose that  $x = [x_i]$  is any solution of the linear system  $Ax = b$ . Using the column vector notation for matrix  $A$  we get that

$$b = Ax = A_{:1}x_1 + \dots + A_{:n}x_n.$$

This last expression says that  $b \in \text{Span}(\{A_{:1}, \dots, A_{:n}\})$ .

( $\Leftarrow$ ) If  $b \in \text{Span}(\{A_{:1}, \dots, A_{:n}\})$ , then there exists constants  $c_1, \dots, c_n$  such that

$$b = A_{:1}c_1 + \dots + A_{:n}c_n = Ac, \quad c = [c_i].$$

This last equation says that the vector  $x = c$  is a solution of the linear system  $Ax = b$ , so the system is consistent. This establishes the Theorem.  $\square$

Knowing the solutions of an homogeneous linear system provides important information about the solutions of an inhomogeneous linear system with the same coefficient matrix. The next result establishes this relation in a precise way.

**Theorem 1.4.8.** If the  $m \times n$  linear system  $Ax = b$  is consistent and the vector  $x_p$  is one particular solution of this linear system, then any solution  $x$  to this system can be decomposed as  $x = x_p + x_h$ , where the vector  $x_h$  is a solution of the homogeneous linear system  $Ax_h = 0$ .

**Proof of Theorem 1.4.8:** We know that the vector  $\mathbf{x}_p$  is a solution of the linear system, that is,  $A\mathbf{x}_p = \mathbf{b}$ . Suppose that there is any other solution  $\mathbf{x}$  of the same system,  $A\mathbf{x} = \mathbf{b}$ . Then, their difference  $\mathbf{x}_h = \mathbf{x} - \mathbf{x}_p$  satisfies the homogeneous equation  $A\mathbf{x}_h = \mathbf{0}$ , since

$$A\mathbf{x}_h = A(\mathbf{x} - \mathbf{x}_p) = A\mathbf{x} - A\mathbf{x}_p = \mathbf{b} - \mathbf{b} = \mathbf{0},$$

where the linearity of the matrix-vector product, proven in Theorem 1.4.4, is used in the step  $A(\mathbf{x} - \mathbf{x}_p) = A\mathbf{x} - A\mathbf{x}_p$ . We have shown that  $\mathbf{x}_h = \mathbf{x} - \mathbf{x}_p$ , that is, any solution  $\mathbf{x}$  of the non-homogeneous system can be written as  $\mathbf{x} = \mathbf{x}_p + \mathbf{x}_h$ . This establishes the Theorem.  $\square$

We say that the solution to a non-homogeneous linear system is written in *vector form* or in *parametric form* when it is expressed as in Theorem 1.4.8, that is, as  $\mathbf{x} = \mathbf{x}_p + \mathbf{x}_h$ , where the vector  $\mathbf{x}_h$  is a solution of the homogeneous linear system, and the vector  $\mathbf{x}_p$  is any solution of the non-homogeneous linear system.

**EXAMPLE 1.4.10:** Find all solutions of the linear system and write them in parametric form,

$$\begin{bmatrix} 2 & -4 \\ 1 & -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 3 \end{bmatrix}. \quad (1.24)$$

**SOLUTION:** We first find the solutions of this non-homogeneous linear system using Gauss elimination operations,

$$\left[ \begin{array}{cc|c} 1 & -2 & 3 \\ 2 & -4 & 6 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & -2 & 3 \\ 0 & 0 & 0 \end{array} \right] \Rightarrow \begin{cases} x_1 = 2x_2 + 3, \\ x_2 : \text{free variable.} \end{cases}$$

Therefore, the set of all solutions of the linear system above is given by

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2x_2 + 3 \\ x_2 \end{bmatrix} \Rightarrow \mathbf{x} = \begin{bmatrix} 2 \\ 1 \end{bmatrix} x_2 + \begin{bmatrix} 3 \\ 0 \end{bmatrix}.$$

In this case we see that

$$\mathbf{x}_p = \begin{bmatrix} 3 \\ 0 \end{bmatrix} \quad (\text{by choosing } x_2 = 0), \quad \text{and} \quad \mathbf{x}_h = \begin{bmatrix} 2 \\ 1 \end{bmatrix} x_2.$$

The vector  $\mathbf{x}_p$  is the particular solution to the non-homogeneous system given by  $x_2 = 0$ , while it is not difficult to check that  $\mathbf{x}_h$  above is solution of the homogeneous equation  $A\mathbf{x}_h = \mathbf{0}$ . In Fig. 16 we plot these solutions on the plane. The solution of the non-homogeneous system is the translation by  $\mathbf{x}_p$  of the solutions of the homogeneous system.  $\triangleleft$

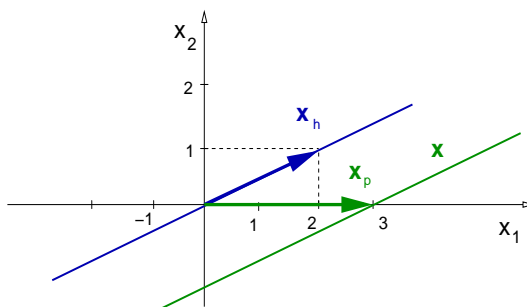


FIGURE 16. The blue line represents solutions to the homogeneous system in Eq. (1.24). The green line represents the solutions to the non-homogeneous system in Eq. (1.24), which is the translation by  $\mathbf{x}_p$  of the line passing by the origin.

**Further reading.** See Sections 2.4 and 2.5 in Meyer's book [3] for a detailed discussion on homogeneous and non-homogeneous linear systems, respectively. See Section 1.4 in Lay's book [2] for a detailed discussion of the matrix-vector product, and Section 1.5 for detailed discussions on homogeneous and non-homogeneous linear systems.

## 1.4.6. Exercises.

- 1.4.1.- Find the general solution of the homogeneous linear system

$$\begin{aligned}x_1 + 2x_2 + x_3 + 2x_4 &= 0 \\2x_1 + 4x_2 + x_3 + 3x_4 &= 0 \\3x_1 + 6x_2 + x_3 + 4x_4 &= 0.\end{aligned}$$

- 1.4.2.- Find all the solutions  $\mathbf{x}$  of the linear system  $\mathbf{Ax} = \mathbf{b}$ , where  $\mathbf{A}$  and  $\mathbf{b}$  are given by

$$\mathbf{A} = \begin{bmatrix} 1 & -2 & -1 \\ 2 & 1 & 8 \\ 1 & -1 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix},$$

and write these solutions in parametric form, that is, in terms of column vectors.

- 1.4.3.- Prove the following statement: If the vectors  $\mathbf{c}$  and  $\mathbf{d}$  are solutions of the homogeneous linear system  $\mathbf{Ax} = \mathbf{0}$ , then  $\mathbf{c} + \mathbf{d}$  is also a solution.

- 1.4.4.- Find the general solution of the non-homogeneous linear system

$$\begin{aligned}x_1 + 2x_2 + x_3 + 2x_4 &= 3 \\2x_1 + 4x_2 + x_3 + 3x_4 &= 4 \\3x_1 + 6x_2 + x_3 + 4x_4 &= 5.\end{aligned}$$

- 1.4.5.- Suppose that the solution to a system of linear equation is given by

$$\begin{aligned}x_1 &= 5 + 4x_3 \\x_2 &= -2 - 7x_3 \\x_3 &\text{ free.}\end{aligned}$$

Use column vectors to describe this set as a line in  $\mathbb{R}^3$ .

- 1.4.6.- Suppose that the solution to a system of linear equation is given by

$$\begin{aligned}x_1 &= 3x_4 \\x_2 &= 8 + 4x_4 \\x_3 &= 2 - 5x_4 \\x_4 &\text{ free.}\end{aligned}$$

Use column vectors to describe this set as a line in  $\mathbb{R}^4$ .

- 1.4.7.- Consider the following system of linear equations, where  $k$  represents any real number,

$$\begin{bmatrix} 2 & 2 & 3 \\ 4 & 8 & 12 \\ 6 & 2 & k \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ -4 \\ 4 \end{bmatrix}.$$

- (a) Find all possible values of the number  $k$  such that the system above has a unique solution.
- (b) Find all possible values of the number  $k$  such that the system above has infinitely many solutions, and express those solutions in parametric form.



## 1.5. FLOATING-POINT NUMBERS

**1.5.1. Main definitions.** Floating-point numbers are a finite subset of the rational numbers. Many different types of floating-point numbers exist, all of them are characterized by having a finite number of digits when written in a particular base. Digital computers use floating-point numbers to carry out almost every arithmetic operation. When an  $m \times n$  algebraic linear system is solved using a computer, every Gauss operations is performed in a particular set of floating-point numbers. In this Section we study what type of approximations occur in this process.

**Definition 1.5.1.** A non-zero rational number  $x$  is a **floating-point number** in base  $b \in \mathbb{N}$ , of precision  $p \in \mathbb{N}$ , with exponent  $n \in \mathbb{Z}$  in the range  $-N \leq n \leq N \in \mathbb{N}$ , iff there exist integers  $d_i$ , for  $i = 1, \dots, p$ , satisfying  $0 \leq d_i \leq b-1$  and  $d_1 \neq 0$ , such that the number  $x$  has the form

$$x = \pm 0.d_1 \cdots d_p \times b^n. \quad (1.25)$$

We call  $p$  the **precision**,  $b$  the **base** and  $N$  the **exponent range** of the floating point number  $x$ . We denote by  $\mathbb{F}_{p,b,N}$  the set of all floating-point numbers of fixed precision  $p$ , base  $b$  and exponent range  $N$ .

In this notes we always work in base  $b = 10$ . Computers usually work with base  $b = 2$ , but also with base  $b = 16$ , and they present their results with base  $b = 10$ .

**EXAMPLE 1.5.1:** The following numbers belong to the set  $\mathbb{F}_{3,10,3}$ ,

$$210 = 0.210 \times 10^3, \quad 1 = 0.100 \times 10, \quad -0.02 = -0.200 \times 10^{-1}, \quad \frac{215}{10^6} = 0.215 \times 10^{-3}.$$

The set  $\mathbb{F}_{3,10,3}$  is a finite subset of the rational numbers. The biggest number and the smallest number in absolute value are the following, respectively,

$$0.999 \times 10^3 = 999 \quad 0.100 \times 10^{-3} = 0.0001.$$

Any number bigger 999 or closer to 0 than 0.0001 does not belong to  $\mathbb{F}_{3,10,3}$ . Here are other examples of numbers that do not belong to  $\mathbb{F}_{3,10,3}$ ,

$$1000 = 0.100 \times 10^4 \quad 210.3 = 0.2103 \times 10^3, \quad 1.001 = 0.1001 \times 10, \quad 0.000027 = 0.270 \times 10^{-4}.$$

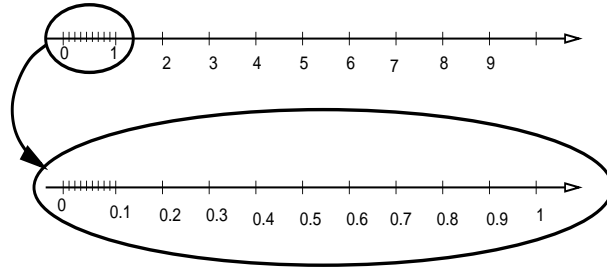
In the first case the number is too big, we need an exponent  $n = 4$  to specify it. In the second and third cases we need a precision  $p = 4$  to specify those numbers. In the last case the number is too close to zero, we need an exponent  $n = -4$  to specify it.  $\triangleleft$

**EXAMPLE 1.5.2:** The set of floating-point numbers  $\mathbb{F}_{1,10,1}$  is small enough to picture it on the real line. The set of all positive elements in  $\mathbb{F}_{1,10,1}$  is shown on Fig. 17, and this is the union of the following three sets,

$$\begin{aligned} \{0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09\} &= \{0.i \times 10^{-1}\}_{i=1}^9; \\ \{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9\} &= \{0.i \times 10^0\}_{i=1}^9; \\ \{1, 2, 3, 4, 5, 6, 7, 8, 9\} &= \{0.i \times 10^1\}_{i=1}^9. \end{aligned}$$

One can see in this example that the elements on  $\mathbb{F}_{1,10,1}$  are not homogeneously distributed on the interval  $[0, 10]$ . The irregular distribution of the floating-point numbers plays an important role when one computes the addition of a small number to a big number.  $\triangleleft$

**EXAMPLE 1.5.3:** In Table 1 we show the main set of floating-point numbers used nowadays in computers. For example, the format called Binary64 represents all floating-point numbers

FIGURE 17. All the positive elements in  $\mathbb{F}_{1,10,1}$ .

in the set  $\mathbb{F}_{53,2,1024}$ . One of the bigger numbers in this set is  $2^{1023}$ . To have an idea how big is this number, let us rewrite it as  $10^x$ , that is,

$$2^{1023} = 10^x \Leftrightarrow 1023 \ln(2) = x \ln(10) \Leftrightarrow x = 307.95\dots$$

So the biggest number in  $\mathbb{F}_{53,2,1024}$  is close to  $10^{308}$ . ◀

Format name	Base $b$	Digits $p$	Max. exp. $N$
Binary 32	2	24	128
Binary 64	2	53	1024
Binary 128	2	113	16384
Decimal 164	10	16	385
Decimal 128	10	34	6145

TABLE 1. List of the parameters  $b$ ,  $p$  and  $N$  that determine the floating-point sets  $\mathbb{F}_{p,b,N}$ , which are most used in computers. The first column presents a standard name given in the scientific computing community to these floating-point formats.

We say that a set  $A \subset \mathbb{R}$  is *closed under addition* iff for every two elements in  $A$  the sum of these two numbers also belongs to  $A$ . The definition of a set  $A \subset \mathbb{R}$  being closed under multiplication is similar. The sets of floating-point numbers  $\mathbb{F}_{p,b,N} \subset \mathbb{R}$  are not closed under addition or multiplication. This means that the sum of two numbers in  $\mathbb{F}_{p,b,N}$  might not belong to  $\mathbb{F}_{p,b,N}$ . And the multiplication of two numbers in  $\mathbb{F}_{p,b,N}$  might not belong to  $\mathbb{F}_{p,b,N}$ . Here are some examples.

**EXAMPLE 1.5.4:** Consider the set  $\mathbb{F}_{2,10,2}$ . It is not difficult to see that  $\mathbb{F}_{2,10,2}$  is not closed under multiplication, as the first line below shows. The rest of the example below shows that the sum of two numbers in  $\mathbb{F}_{2,10,2}$  does not belong to that set.

$$x = 10^{-3} = 0.10 \times 10^{-2} \in \mathbb{F}_{2,10,2} \Rightarrow x^2 = 10^{-6} = 0.0001 \times 10^{-2} \notin \mathbb{F}_{2,10,2},$$

$$\left. \begin{array}{l} x = 10^{-3} = 0.10 \times 10^{-2} \in \mathbb{F}_{2,10,2}, \\ y = 1 = 0.10 \times 10^1 \in \mathbb{F}_{2,10,2}, \end{array} \right\} \Rightarrow x + y = 0.001 + 1,$$

$$x + y = 1.001 \quad \Rightarrow \quad x + y = 0.1001 \times 10 \notin \mathbb{F}_{2,10,2}.$$

◁

**1.5.2. The rounding function.** Since the set  $\mathbb{F}_{p,b,N}$  is not closed under addition or multiplication, not every arithmetic calculation involving real numbers can be performed in  $\mathbb{F}_{p,b,N}$ . A way to perform a sequence of arithmetic operations in  $\mathbb{F}_{p,b,N}$  is first, to project the real numbers into the floating-point numbers, and then to perform the calculation. Since the result might not be in the set  $\mathbb{F}_{p,b,N}$ , one must project again the result onto  $\mathbb{F}_{p,b,N}$ . The action to project a real number into the floating-point set is called to *round-off* the real number. There are many different ways to do this. We now present a common round-off function.

**Definition 1.5.2.** Given the floating-point number set  $\mathbb{F}_{p,b,N}$ , let  $x_N = 0.9 \cdots 9 \times b^N$  be the biggest number in  $\mathbb{F}_{p,b,N}$ . The **rounding function**  $f_\ell : \mathbb{R} \cap [-x_N, x_N] \rightarrow \mathbb{F}_{p,b,N}$  is defined as follow: Given  $x \in \mathbb{R} \cap [-x_N, x_N]$ , with  $x = \pm 0.d_1 \cdots d_p d_{p+1} \cdots \times b^n$  and  $-N \leq n \leq N$ , holds

$$f_\ell(x) = \begin{cases} \pm 0.d_1 \cdots d_p \times b^n & \text{if } d_{p+1} < 5, \\ \pm (0.d_1 \cdots d_p + b^{-p}) \times b^n & \text{if } d_{p+1} \geq 5. \end{cases}$$

**EXAMPLE 1.5.5:** We now present few numbers not in  $\mathbb{F}_{3,10,3}$  and their respective round-offs

$$\begin{aligned} x = 0.2103 \times 10^3, & & f_\ell(x) = 0.210 \times 10^3, \\ x = 0.21037 \times 10^3, & & f_\ell(x) = 0.210 \times 10^3, \\ x = 0.2105 \times 10^3, & & f_\ell(x) = 0.211 \times 10^3, \\ x = 0.21051 \times 10^3, & & f_\ell(x) = 0.211 \times 10^3. \end{aligned}$$

◁

The rounding function has the following properties:

**Proposition 1.5.3.** Given  $\mathbb{F}_{p,b,N}$  there always exist  $x, y \in \mathbb{R}$  such that

$$f_\ell[f_\ell(x) + f_\ell(y)] \neq f_\ell(x + y), \quad f_\ell[f_\ell(x)f_\ell(y)] \neq f_\ell(xy).$$

We do not prove this Proposition, we only provide a particular example, in the case that the floating-point number is  $\mathbb{F}_{2,10,2}$ .

**EXAMPLE 1.5.6:** The real numbers  $x = 21/2$  and  $y = 11/2$  add up to  $x + y = 16$ . Only one of them belongs to  $\mathbb{F}_{2,10,2}$ , since

$$\begin{aligned} x = 0.105 \times 10^2 \notin \mathbb{F}_{2,10,2} & \Rightarrow & f_\ell(x) = 0.11 \times 10^2, \\ y = 0.55 \times 10 \in \mathbb{F}_{2,10,2} & \Rightarrow & f_\ell(y) = 0.55 \times 10 = y. \end{aligned}$$

We now verify that for these numbers holds that  $f_\ell[f_\ell(x) + f_\ell(y)] \neq f_\ell(x + y)$ , since

$$x + y = 0.16 \times 10^2 \in \mathbb{F}_{2,10,2} \quad \Rightarrow \quad f_\ell(x + y) = 0.16 \times 10^2 = x + y,$$

$$f_\ell[f_\ell(x) + f_\ell(y)] = f_\ell(0.11 \times 10^2 + 0.55 \times 10) = f_\ell(0.165 \times 10^2) = 0.17 \times 10^2.$$

Therefore, we conclude that

$$17 = f_\ell[f_\ell(x) + f_\ell(y)] \neq f_\ell(x + y) = 16.$$

◁

**1.5.3. Solving linear systems.** Arithmetic operations are, in general, not possible in  $\mathbb{F}_{p,b,N}$ , since this set is not closed under these operations. Rounding after each arithmetic operation is a procedure that determines numbers in  $\mathbb{F}_{p,b,N}$  which are close to the result of the arithmetic operations. The difference between these two numbers is the error of the procedure, is the error of making calculations in a finite set of numbers. One could think that the Gauss-Jordan method could be carried out in the set  $\mathbb{F}_{p,b,N}$ , if one round-off each intermediate calculation. This Subsection presents an example that this is not the case. The Gauss-Jordan method is, in general, *not possible* in any set  $\mathbb{F}_{p,b,N}$  using rounding after each arithmetic operation.

**EXAMPLE 1.5.7:** Use the floating-point set  $\mathbb{F}_{3,10,3}$  and Gauss operations to find the solution of the  $2 \times 2$  linear system

$$\begin{aligned} 5x_1 + x_2 &= 6, \\ 9.43x_1 + 1.57x_2 &= 11. \end{aligned}$$

**SOLUTION:** Note that the solution is  $x_1 = 1$ ,  $x_2 = 1$ . The calculation we do now shows that the Gauss-Jordan method is not possible in the set  $\mathbb{F}_{3,10,3}$  using round-off functions; the Gauss-Jordan method does not work in this example. The first step in the Gauss-Jordan method is to compute the augmented matrix of the system and perform a Gauss operation to make the coefficient in the position of  $a_{21}$  vanish, that is,

$$\left[ \begin{array}{cc|c} 5 & 1 & 6 \\ 9.43 & 1.57 & 11 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 5 & 1 & 6 \\ 0 & -0.316 & -0.316 \end{array} \right].$$

This operation cannot be performed in the set  $\mathbb{F}_{3,10,3}$  using rounding. In order to understand this, let us review what we did in the calculation above. We multiplied the first row by  $-9.43/5$  and add that result to the second row. The result using real numbers is that the new coefficient obtained after this calculation is  $\tilde{a}_{21} = 0$ . If we do this calculation in the set  $\mathbb{F}_{3,10,3}$  using rounding, we have to do the following calculation:

$$\tilde{a}_{21} = f_\ell \left( 9.43 - f_\ell \left[ f_\ell(5) f_\ell \left( \frac{9.43}{5} \right) \right] \right),$$

that is, we round the quotient  $-9.43/5$ , then we multiply by 5, we round again, then we subtract that from 9.43, and we finally round the result:

$$\begin{aligned} \tilde{a}_{21} &= f_\ell(9.43 - f_\ell[5 f_\ell(1.886)]) \\ &= f_\ell(9.43 - f_\ell[5(1.89)]) \\ &= f_\ell[9.43 - f_\ell(9.45)] \\ &= f_\ell[9.43 - 9.45] \\ &= f_\ell(-0.02) \\ &= -0.02. \end{aligned}$$

Therefore, with this Gauss operation on  $\mathbb{F}_{3,10,3}$  using rounding one obtains  $\tilde{a}_{21} = -0.02 \neq 0$ . The same type of calculation on the other coefficients  $\tilde{a}_{22}$ ,  $\tilde{b}_2$ , produces the following new augmented matrix

$$\left[ \begin{array}{cc|c} 5 & 1 & 6 \\ 9.43 & 1.57 & 11 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 5 & 1 & 6 \\ -0.02 & -0.32 & -0.3 \end{array} \right].$$

The Gauss-Jordan method cannot follow unless the coefficient  $\tilde{a}_{21} = 0$ . But this is not possible in our example. A usual procedure used in scientific computation is to modify the Gauss-Jordan method. The modification introduces further approximation errors in

the calculation. The modification in our example is the following: Replace the coefficient  $\tilde{a}_{21} = -0.02$  by  $\tilde{a}_{21} = 0$ . The modified Gauss-Jordan method in our example is given by

$$\left[ \begin{array}{cc|c} 5 & 1 & 6 \\ 9.43 & 1.57 & 11 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 5 & 1 & 6 \\ 0 & -0.32 & -0.3 \end{array} \right]. \quad (1.26)$$

What we have done here is not a rounding error. It is a modification of the Gauss-Jordan method to find an approximate solution of the linear system in the set  $\mathbb{F}_{3,10,3}$ . The rest of the calculation to find the solution of the linear system is the following:

$$\left[ \begin{array}{cc|c} 5 & 1 & 6 \\ 0 & -0.32 & -0.3 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 5 & 1 & 6 \\ 0 & 1 & f_\ell\left(\frac{0.3}{0.32}\right) \end{array} \right];$$

since

$$f_\ell\left(\frac{0.3}{0.32}\right) = f_\ell(0.9375) = 0.938,$$

we have that

$$\left[ \begin{array}{cc|c} 5 & 1 & 6 \\ 0 & 1 & f_\ell\left(\frac{0.3}{0.32}\right) \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 5 & 1 & 6 \\ 0 & 1 & 0.938 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 5 & 0 & f_\ell(6 - 0.938) \\ 0 & 1 & 0.938 \end{array} \right];$$

since

$$f_\ell(6 - 0.938) = f_\ell(5.062) = 5.06,$$

we also have that

$$\left[ \begin{array}{cc|c} 5 & 0 & 5.06 \\ 0 & 1 & 0.938 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 0 & f_\ell\left(\frac{5.06}{5}\right) \\ 0 & 1 & 0.938 \end{array} \right];$$

since

$$f_\ell\left(\frac{5.06}{5}\right) = f_\ell(1.012) = 1.01,$$

we conclude that

$$\left[ \begin{array}{cc|c} 1 & 0 & 1.01 \\ 0 & 1 & 0.938 \end{array} \right] \Rightarrow \begin{array}{l} x_1 = 0.101 \times 10, \\ x_2 = 0.938. \end{array}.$$

We conclude that the solution in the set  $\mathbb{F}_{3,10,3}$  differs from the exact solution  $x_1 = 1$ ,  $x_2 = 1$ . The errors in the result are produced by rounding errors and by the modification of the Gauss-Jordan method discussed in Eq. (1.26).  $\triangleleft$

We finally comment that the round-off error becomes important when adding a small number to a big number, or when dividing by a small number. Here is an example of the former case.

**EXAMPLE 1.5.8:** Add together the numbers  $x = 10^3$  and  $y = 4$  in the set  $\mathbb{F}_{3,10,4}$ .

**SOLUTION:** Since  $x = 0.100 \times 10^4$  and  $y = 0.400 \times 10$ , both numbers belong to  $\mathbb{F}_{3,10,4}$  and so  $f_\ell(x) = x$ ,  $f_\ell(y) = y$ . Therefore, their addition is the following,

$$f_\ell(x + y) = f_\ell(1000 + 4) = f_\ell(0.1004 \times 10^3) = 1 \times 10^3 = x.$$

That is,  $f_\ell(x + y) = x$ , and the information of  $y$  is completely lost.  $\triangleleft$

**1.5.4. Reducing rounding errors.** We have seen that solving an algebraic linear system in a floating-point number set  $\mathbb{F}_{p,b,N}$  introduces rounding errors in the solution. There are several techniques that help keep these errors from becoming an important part of the solution. We comment here on four of these techniques. We present these techniques without any proof.

The first two techniques are implemented before one starts to solve the linear system. They are called column scaling and row scaling. The *column scaling* consists in multiplying by a same scale factor a whole column in the linear system. One performs a column scaling when one column of a linear system has coefficients that are far bigger or far smaller than the rest of the matrix. The factor in the column scaling is chosen in order that all the coefficients in the linear system are similar in size. One can interpret the column scaling on the column  $i$  of a linear system as changing the physical units of the unknown  $x_i$ . The *row scaling* consists in multiplying by the same scale factor a whole equation of the system. One performs a row scaling when one row of the linear system has coefficients that are far bigger or far smaller than the rest of the system. Like in the column scaling, one chooses the scaling factor in order that all the coefficients of the linear system are similar in size.

The last two techniques refer to what sequence of Gauss operations is more efficient in reducing rounding errors. It is well-known that there are many different ways to solve a linear system using Gauss operations in the set of the real numbers  $\mathbb{R}$ . For example, we now solve the linear system below using two different sequences of Gauss operations:

$$\left[ \begin{array}{cc|c} 2 & 4 & 10 \\ 1 & 3 & 7 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 2 & 5 \\ 1 & 3 & 7 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 2 & 5 \\ 0 & 1 & 2 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 0 & 1 \\ 0 & 1 & 2 \end{array} \right],$$

$$\left[ \begin{array}{cc|c} 2 & 4 & 10 \\ 1 & 3 & 7 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 3 & 7 \\ 2 & 4 & 10 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 2 & 5 \\ 0 & -2 & -4 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 2 & 5 \\ 0 & 1 & 2 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 0 & 1 \\ 0 & 1 & 2 \end{array} \right].$$

The solution obtained is independent of the sequences of Gauss operations used to find them. This property does not hold in the floating-point number set  $\mathbb{F}_{p,b,N}$ . Two different sequences of Gauss operations on the same augmented matrix might produce different approximate solutions when they are performed in  $\mathbb{F}_{p,b,N}$ . The main idea behind the last two techniques we now present to solve linear systems using floating-point numbers is the following: To find the sequence of Gauss operations that minimize the rounding errors in the approximate solution. The last two techniques are called partial pivoting and complete pivoting.

The *partial pivoting* is the row interchange in a matrix in order to have the biggest coefficient in that column as pivot. Here below we do a partial pivoting for the pivot on the first column:

$$\left[ \begin{array}{ccc} 10^{-2} & 2 & 1 \\ 1 & 10^3 & 10^2 \\ 10 & 4 & 1 \end{array} \right] \rightarrow \left[ \begin{array}{ccc} 10 & 4 & 1 \\ 1 & 10^3 & 10^2 \\ 10^{-2} & 2 & 1 \end{array} \right],$$

that is, use as pivot for the first column the coefficient with 10 instead of the coefficient with  $10^{-2}$  or the coefficient with 1. Now proceed in the usual way:

$$\left[ \begin{array}{ccc} 10 & 4 & 1 \\ 1 & 10^3 & 10^2 \\ 10^{-2} & 2 & 1 \end{array} \right] \rightarrow \left[ \begin{array}{ccc} 1 & 0.4 & 0.1 \\ 1 & 10^3 & 10^2 \\ 10^{-2} & 2 & 1 \end{array} \right] \rightarrow \left[ \begin{array}{ccc} 1 & 0.4 & 0.1 \\ 0 & 999.6 & 99.9 \\ 0 & 1.996 & 0.999 \end{array} \right].$$

The next step is again to use as pivot the biggest coefficient in the second column. In this case it is the coefficient 999.6, so no further row interchanges are necessary. Repeat this procedure till the last column.

The *complete pivoting* is the row and column interchange in a matrix in order to have as pivot the biggest coefficient in the lower-right block from the pivot position. Here below we

do a complete pivoting for the pivot on the first column:

$$\begin{bmatrix} 10^{-2} & 2 & 1 \\ 1 & 10^3 & 10^2 \\ 10 & 4 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 10^3 & 10^2 \\ 10^{-2} & 2 & 1 \\ 10 & 4 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 10^3 & 1 & 10^2 \\ 2 & 10^{-2} & 1 \\ 4 & 10 & 1 \end{bmatrix}.$$

For the first pivot, the lower-right part of the matrix is the whole matrix. In the example above we used as pivot for the first column the coefficient  $10^3$  in position  $(2, 2)$ . We needed to do a row interchange and then a column interchange. Now proceed in the usual way:

$$\begin{bmatrix} 10^3 & 1 & 10^2 \\ 2 & 10^{-2} & 1 \\ 4 & 10 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 10^{-3} & 10^{-1} \\ 2 & 10^{-2} & 1 \\ 4 & 10 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0.001 & 0.1 \\ 0 & 0.008 & 0.8 \\ 0 & 9.996 & 0.6 \end{bmatrix}.$$

The next step in complete pivoting is to choose the coefficient for the pivot position  $(2, 2)$ . We have to look in the lower-right block from the pivot position, that is, in the block

$$\begin{bmatrix} 0.008 & 0.8 \\ 9.996 & 0.6 \end{bmatrix}.$$

The biggest coefficient in that block is 9.996, so we have to do a row interchange but we do not need to do column interchanges, that is,

$$\begin{bmatrix} 1 & 0.001 & 0.1 \\ 0 & 0.008 & 0.8 \\ 0 & 9.996 & 0.6 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0.001 & 0.1 \\ 0 & 9.996 & 0.6 \\ 0 & 0.008 & 0.8 \end{bmatrix}.$$

Repeat this procedure till the last column.

**Further reading.** See Section 1.5 in Meyer's book [3] for a detailed discussion on solving linear systems using floating point numbers.

## 1.5.5. Exercises.

1.5.1.- Consider the following system:

$$\begin{aligned} 10^{-3}x_1 - x_2 &= 1, \\ x_1 + x_2 &= 0. \end{aligned}$$

- (a) Solve this system in  $\mathbb{F}_{3,10,6}$  with rounding, but without partial or complete pivoting.
- (b) Find the system that is exactly satisfied by your solution in (a), and note how close is this system to the original system.
- (c) Use partial pivoting to solve this system in  $\mathbb{F}_{3,10,5}$  with rounding.
- (d) Find the system that is exactly satisfied by your solution in (c), and note how close is this system to the original system.
- (e) Solve this system in  $\mathbb{R}$  without partial or complete pivoting, and compare this exact solution with the solutions in (a) and (c).
- (f) Round the exact solution up to three digits, and compare it with the results from (a) and (c).

1.5.2.- Consider the following system:

$$\begin{aligned} x_1 + x_2 &= 3, \\ -10x_1 + 10^5x_2 &= 10^5. \end{aligned}$$

- (a) Solve this system in  $\mathbb{F}_{4,10,6}$  with partial pivoting but no scaling.
- (b) Solve this system in  $\mathbb{F}_{4,10,6}$  with complete pivoting but no scaling.
- (c) Use partial pivoting to solve this system in  $\mathbb{F}_{3,10,5}$  with rounding.
- (d) This time row scale the original system, and then solve it in  $\mathbb{F}_{4,10,6}$  with partial pivoting.
- (e) Solve this system in  $\mathbb{R}$  and compare this exact solution with the solutions in (a)-(d).

1.5.3.- Consider the linear system

$$\begin{aligned} -3x_1 + x_2 &= -2, \\ 10x_1 - 3x_2 &= 7. \end{aligned}$$

Solve this system in  $\mathbb{F}_{3,10,6}$  without partial pivoting, and then solve it again with partial pivoting. Compare your results with the exact solution.



## CHAPTER 2. MATRIX ALGEBRA

## 2.1. LINEAR TRANSFORMATIONS

In Sect. 1.4 we introduced the matrix-vector product because it provided a convenient notation to express a system of linear equations in a compact way. Here we study a deeper meaning of the matrix-vector product. It is the key to interpret a matrix as a function acting on vectors. We will see that a matrix is a particular type of function, called linear function or linear transformation. We also show examples of functions determined by matrices.

From now on we consider both real-valued and complex-valued matrices and vectors. We use the notation  $\mathbb{F} \in \{\mathbb{R}, \mathbb{C}\}$  to mean that  $\mathbb{F}$  can be either  $\mathbb{F} = \mathbb{R}$  or  $\mathbb{F} = \mathbb{C}$ . Elements in  $\mathbb{F}$  are called *scalars*. So a scalar is a real number or a complex number. We denote by  $\mathbb{F}^n$  the set of all  $n$ -vectors  $\mathbf{x} = [x_i]$  with components  $x_i \in \mathbb{F}$ , where  $i = 1, \dots, n$ . Finally we denote by  $\mathbb{F}^{m,n}$  the set of all  $m \times n$  matrices  $\mathbf{A} = [A_{ij}]$  with components  $A_{ij} \in \mathbb{F}$ , where  $i = 1 \dots m$  and  $j = 1, \dots, n$ .

**2.1.1. A matrix is a function.** The matrix-vector product provides a new interpretation for a matrix. A matrix is not only an artifact for a compact notation, it defines a function on the set of vectors. Here is a precise definition.

**Definition 2.1.1.** Every  $m \times n$  matrix  $\mathbf{A} \in \mathbb{F}^{m,n}$  defines the function  $\mathbf{A} : \mathbb{F}^n \rightarrow \mathbb{F}^m$ , where the image of an  $n$ -vector  $\mathbf{x}$  is the  $m$ -vector  $\mathbf{y} = \mathbf{A}\mathbf{x}$ , the matrix-vector product of  $\mathbf{A}$  and  $\mathbf{x}$ .

**EXAMPLE 2.1.1:** Compute the function defined by the  $2 \times 3$  matrix  $\mathbf{A} = \begin{bmatrix} 2 & -2 & 4 \\ 1 & 3 & 2 \end{bmatrix}$ .

**SOLUTION:** The matrix  $\mathbf{A}$  defines a function  $\mathbf{A} : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ , since for every  $\mathbf{x} \in \mathbb{R}^3$  holds,

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}, \quad \mathbf{A}\mathbf{x} = \begin{bmatrix} 2 & -2 & 4 \\ 1 & 3 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \Rightarrow \mathbf{y} = \mathbf{A}\mathbf{x} = \begin{bmatrix} 2x_1 - 2x_2 + 4x_3 \\ x_1 + 3x_2 + 2x_3 \end{bmatrix} \in \mathbb{R}^2.$$

For example, given  $\mathbf{x} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \in \mathbb{R}^3$ , then  $\mathbf{y} = \mathbf{A}\mathbf{x} = \begin{bmatrix} 2 & -2 & 4 \\ 1 & 3 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 4 \\ 6 \end{bmatrix} \in \mathbb{R}^2$ . ◁

**EXAMPLE 2.1.2:** Describe the function defined by the  $2 \times 2$  matrix  $\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$ .

**SOLUTION:** The matrix  $\mathbf{A}$  defines a function  $\mathbf{A} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , which can be interpreted as a *reflection along the horizontal line*. Indeed, the action of the matrix  $\mathbf{A}$  on an arbitrary element in  $\mathbf{x} \in \mathbb{R}^2$  is the following,

$$\mathbf{A}\mathbf{x} = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_1 \\ -x_2 \end{bmatrix}.$$

Here are particular cases, represented in Fig. 18,

$$\mathbf{A} \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}, \quad \mathbf{A} \begin{bmatrix} -1 \\ -3 \end{bmatrix} = \begin{bmatrix} -1 \\ 3 \end{bmatrix}, \quad \mathbf{A} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}. \quad (2.1)$$

◁

**EXAMPLE 2.1.3:** Describe the function defined by the  $2 \times 2$  matrix  $\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$ .

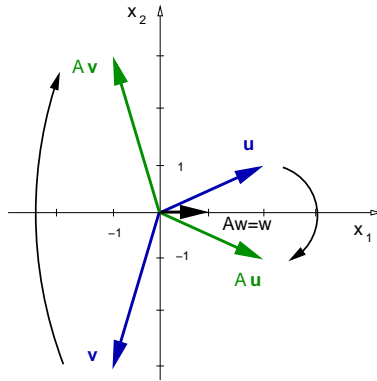


FIGURE 18. We sketch the action of matrix  $A$  in Example 2.1.2 on the vectors given in Eq. (2.1), which we called  $u$ ,  $v$ , and  $w$ , respectively. Since  $A : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , we used the same plane to plot both  $u$  and  $Au$ .

**SOLUTION:** The matrix  $A$  defines a function  $A : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , which can be interpreted as a *reflection along the line  $x_1 = x_2$* , see Fig. 19. Indeed, the action of the matrix  $A$  on an arbitrary element in  $x \in \mathbb{R}^2$  is the following,

$$Ax = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ x_1 \end{bmatrix}.$$

Here are particular cases, represented in Fig. 19,

$$A \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad A \begin{bmatrix} -3 \\ -1 \end{bmatrix} = \begin{bmatrix} -1 \\ -3 \end{bmatrix}, \quad A \begin{bmatrix} 2 \\ 2 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}. \quad (2.2)$$

◁

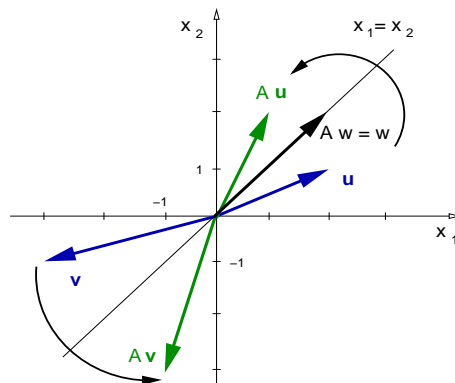


FIGURE 19. We sketch the action of matrix  $A$  in Example 2.1.3 on the vectors given in Eq. (2.2), which we called  $u$ ,  $v$ , and  $w$ , respectively. Since  $A : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , we used the same plane to plot both  $u$  and  $Au$ .

**EXAMPLE 2.1.4:** Describe the function defined by the  $2 \times 2$  matrix  $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$ .

**SOLUTION:** The matrix  $A$  defines a function  $A : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , that can be interpreted as a *rotation by an angle  $\pi/2$  counterclockwise*, see Fig. 20. Indeed, the action of the matrix  $A$  on an arbitrary element in  $x \in \mathbb{R}^2$  is the following,

$$Ax = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -x_2 \\ x_1 \end{bmatrix}.$$

Here are particular cases, represented in Fig. 20,

$$A \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}, \quad A \begin{bmatrix} -3 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ -3 \end{bmatrix}, \quad A \begin{bmatrix} 2 \\ 2 \end{bmatrix} = \begin{bmatrix} -2 \\ 2 \end{bmatrix}. \quad (2.3)$$

◀

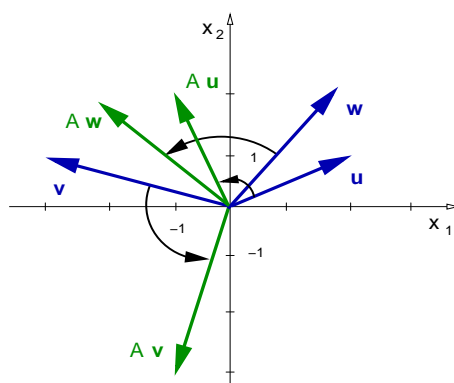


FIGURE 20. We sketch the action of matrix  $A$  in Example 2.1.3 on the vectors given in Eq. (2.2), which we called  $u$ ,  $v$ , and  $w$ , respectively. Since  $A : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , we used the same plane to plot both  $u$  and  $Au$ .

**EXAMPLE 2.1.5:** Describe the function defined by the  $2 \times 2$  matrix  $A = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}$ , where  $\theta$  is a fixed real number.

**SOLUTION:** The matrix  $A$  defines a function  $A : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , that can be interpreted as a *rotation by an angle  $\theta$  counterclockwise*. To verify this, first compute  $y = Ax$  for an arbitrary vector  $x$ , and then check that vector  $y$  is the counterclockwise rotation by  $\theta$  of vector  $x$ .

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_1 \cos(\theta) - x_2 \sin(\theta) \\ x_1 \sin(\theta) + x_2 \cos(\theta) \end{bmatrix} \Rightarrow \begin{cases} y_1 = x_1 \cos(\theta) - x_2 \sin(\theta), \\ y_2 = x_1 \sin(\theta) + x_2 \cos(\theta). \end{cases}$$

We now show that the components  $y_1$  and  $y_2$  above are precisely the counterclockwise rotation by  $\theta$  of the vector  $x$ . From Fig. 21 we see that the following relation holds:

$$y_1 = \cos(\theta + \phi) \|y\|, \quad y_2 = \sin(\theta + \phi) \|y\|,$$

where  $\|y\|$  is the magnitude of the vector  $y$ . Since a rotation does not change the magnitude of the vector, then  $\|y\| = \|x\|$  and so,

$$y_1 = \cos(\theta + \phi) \|x\|, \quad y_2 = \sin(\theta + \phi) \|x\|.$$

Recalling now the formulas for the cosine and the sine of a sum of two angles,

$$\begin{aligned} \cos(\theta + \phi) &= \cos(\theta) \cos(\phi) - \sin(\theta) \sin(\phi), \\ \sin(\theta + \phi) &= \sin(\theta) \cos(\phi) + \cos(\theta) \sin(\phi), \end{aligned}$$

we obtain that

$$\begin{aligned} y_1 &= \cos(\theta) \cos(\phi) \|\mathbf{x}\| - \sin(\theta) \sin(\phi) \|\mathbf{x}\|, \\ y_2 &= \sin(\theta) \cos(\phi) \|\mathbf{x}\| + \cos(\theta) \sin(\phi) \|\mathbf{x}\|. \end{aligned}$$

Recalling that

$$x_1 = \cos(\phi) \|\mathbf{x}\|, \quad x_2 = \sin(\phi) \|\mathbf{x}\|,$$

we obtain the formula

$$\begin{aligned} y_1 &= \cos(\theta) x_1 - \sin(\theta) x_2, \\ y_2 &= \sin(\theta) x_1 + \cos(\theta) x_2. \end{aligned}$$

This is precisely the expression that defines matrix  $A$ . Therefore, the action of the matrix  $A$  above is a rotation by  $\theta$  counterclockwise.  $\triangleleft$

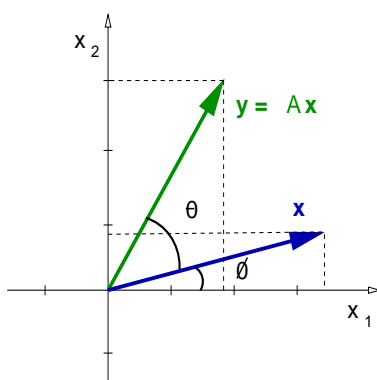


FIGURE 21. The vector  $y$  is the rotation by an angle  $\theta$  counterclockwise of the vector  $x$ .

**EXAMPLE 2.1.6:** Describe the function defined by the  $2 \times 2$  matrix  $A = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$ .

**SOLUTION:** The matrix  $A$  defines a function  $A : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , which can be interpreted as a *dilation*, see Fig. 22. Indeed, the action of the matrix  $A$  on an arbitrary element in  $x \in \mathbb{R}^2$  is the following,  $Ax = 2x$ .  $\triangleleft$

**EXAMPLE 2.1.7:** Describe the function defined by the  $2 \times 2$  matrix  $A = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$ .

**SOLUTION:** The matrix  $A$  defines a function  $A : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , which can be interpreted as a *shear*, see Fig. 23. Indeed, the action of the matrix  $A$  on an arbitrary element in  $x \in \mathbb{R}^2$  is the following,

$$Ax = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2x_1 \\ x_2 \end{bmatrix}.$$

$\triangleleft$

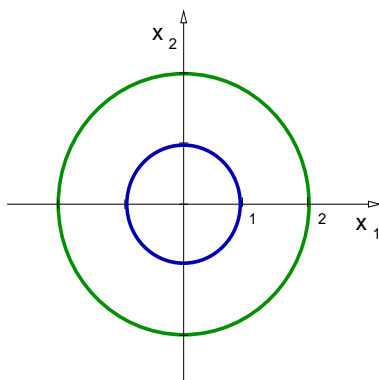


FIGURE 22. We sketch the action of matrix  $A$  in Example 2.1.6. Given any vector  $x$  with end point on the circle of radius one, the vector  $Ax = 2x$  is a vector parallel to  $x$  and with end point on the dashed circle of radius two.

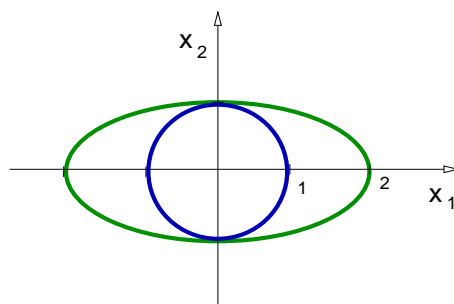


FIGURE 23. We sketch the action of matrix  $A$  in Example 2.1.7. Given any vector  $x$  with end point on the circle of radius one, the vector  $Ax = 2x$  is a vector parallel to  $x$  and with end point on the dashed curve.

**2.1.2. A matrix is a linear function.** We have seen several examples of functions given by matrices. All these functions are defined using the matrix-vector product. We have seen in Theorem 1.4.4 that the matrix-vector product is a linear operation. That is, given any  $m \times n$  matrix  $A$ , for all vectors  $x, y \in \mathbb{F}^n$  and all scalars  $a, b \in \mathbb{F}$  holds that

$$A(ax + by) = aAx + bAy$$

This property of the function defined by a matrix will be important later on, so any function with this property will be given a particular name.

**Definition 2.1.2.** A function  $T : \mathbb{F}^n \rightarrow \mathbb{F}^m$  is called a **linear transformation** iff for all vectors  $x, y \in \mathbb{F}^n$  and for all scalars  $a, b \in \mathbb{F}$  holds

$$T(ax + by) = aT(x) + bT(y).$$

The expression above contains the particular cases  $a = b = 1$  and  $b = 0$ , which are respectively given by

$$T(x + y) = T(x) + T(y), \quad T(ax) = aT(x).$$

We will also use the name **linear function** for a linear transformation. At the end of Section 2.2 we generalize the definition of a linear transformation to include functions on the

space of matrices. We will present two examples of linear functions on the space of matrices, the transpose and the trace functions. In this Section we present simpler examples of linear functions.

**EXAMPLE 2.1.8:** Show that the only linear transformations  $T : \mathbb{R} \rightarrow \mathbb{R}$  are straight lines through the origin.

**SOLUTION:** Since  $y = T(x)$  is linear and  $x \in \mathbb{R}$ , we have that

$$y = T(x) = T(x \times 1) = xT(1).$$

If we denote  $T(1) = m$ , we then conclude that a linear transformation  $T : \mathbb{R} \rightarrow \mathbb{R}$  must be

$$y = mx,$$

which is a straight line through the origin with slope  $m$ . ◁

**EXAMPLE 2.1.9:** Any  $m \times n$  matrix  $A \in \mathbb{F}^{m,n}$  defines a linear transformation  $T : \mathbb{F}^n \rightarrow \mathbb{F}^m$  by the equation  $T(\mathbf{x}) = A\mathbf{x}$ . In particular, all the functions defined in Examples 2.1.1-2.1.7 are linear transformations. Consider the most general  $2 \times 2$  matrix,

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \in \mathbb{F}^{2,2}$$

and explicitly show that the function  $T(\mathbf{x}) = A\mathbf{x}$  is a linear transformation.

**SOLUTION:** The explicit form of the function  $T : \mathbb{F}^2 \rightarrow \mathbb{F}^2$  given by  $T(\mathbf{x}) = A\mathbf{x}$  is

$$T\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Rightarrow T\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 \\ A_{21}x_1 + A_{22}x_2 \end{bmatrix}.$$

This function is linear, as it can be seen from the following explicit computation,

$$\begin{aligned} T(c\mathbf{x} + d\mathbf{y}) &= T\left(c \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + d \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}\right) \\ &= T\left(\begin{bmatrix} cx_1 + dy_1 \\ cx_2 + dy_2 \end{bmatrix}\right) \\ &= \begin{bmatrix} A_{11}(cx_1 + dy_1) + A_{12}(cx_2 + dy_2) \\ A_{21}(cx_1 + dy_1) + A_{22}(cx_2 + dy_2) \end{bmatrix} \\ &= \begin{bmatrix} A_{11}cx_1 + A_{12}cx_2 \\ A_{21}cx_1 + A_{22}cx_2 \end{bmatrix} + \begin{bmatrix} A_{11}dy_1 + A_{12}dy_2 \\ A_{21}dy_1 + A_{22}dy_2 \end{bmatrix} \\ &= c \begin{bmatrix} A_{11}x_1 + A_{12}x_2 \\ A_{21}x_1 + A_{22}x_2 \end{bmatrix} + d \begin{bmatrix} A_{11}y_1 + A_{12}y_2 \\ A_{21}y_1 + A_{22}y_2 \end{bmatrix} \\ &= cT(\mathbf{x}) + dT(\mathbf{y}). \end{aligned}$$

This establishes that  $T$  is a linear transformation. Notice that this proof is just a repetition of the Theorem 1.4.4 proof in the case of  $2 \times 2$  matrices. ◁

**EXAMPLE 2.1.10:** Find a function  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  that projects a vector onto the line  $x_1 = x_2$ , see Fig. 24. Show that this function is linear. Finally, find a matrix  $A$  such that  $T(\mathbf{x}) = A\mathbf{x}$ .

**SOLUTION:** From Fig. 24 one can see that a possible way to compute the projection of a vector  $\mathbf{x}$  onto the line  $x_1 = x_2$  is the following: Add to the vector  $\mathbf{x}$  its reflection along the line  $x_1 = x_2$ , and divide the result by two, that is,

$$T\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \frac{1}{2}\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} x_2 \\ x_1 \end{bmatrix}\right) \Rightarrow T\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \frac{(x_1 + x_2)}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix}.$$

We have obtained the projection function  $T$ . We now show that this function is linear: Indeed

$$\begin{aligned} T(ax + by) &= T\left(a \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + b \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}\right) \\ &= T\left(\begin{bmatrix} ax_1 + by_1 \\ ax_2 + by_2 \end{bmatrix}\right) \\ &= \frac{(ax_1 + by_1 + ax_2 + by_2)}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \\ &= \frac{(ax_1 + ax_2)}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \frac{(by_1 + by_2)}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \\ &= aT(x) + bT(y). \end{aligned}$$

This shows that  $T$  is linear. We now find a matrix  $A$  such that  $T(x) = Ax$ , as follows

$$\begin{aligned} T\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) &= \frac{(x_1 + x_2)}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \\ &= \frac{1}{2} \begin{bmatrix} x_1 + x_2 \\ x_1 + x_2 \end{bmatrix} \\ &= \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Rightarrow A = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}. \end{aligned}$$

This matrix projects vectors onto the line  $x_1 = x_2$ . ◁

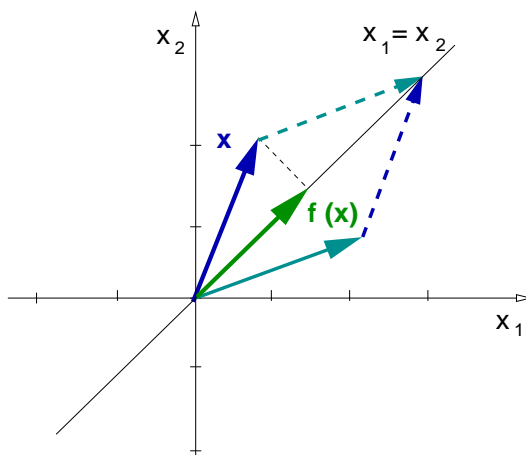


FIGURE 24. The function  $T$  projects the vector  $x$  onto the line  $x_1 = x_2$ .

**Further reading.** See Sections 1.8 and 1.9 in Lay's book [2]. Also Sections 3.3 and 3.4 in Meyer's book [3].

## 2.1.3. Exercises.

2.1.1.- Determine which of the following functions  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  is linear:

(a)

$$T\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} 3x_2 \\ 2 + x_1 \end{bmatrix}.$$

(b)

$$T\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1 + x_2 \\ x_1 - x_2 \end{bmatrix}.$$

(c)

$$T\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1 x_2 \\ 0 \end{bmatrix}.$$

2.1.2.- Let  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be the linear transformation given by  $T(\mathbf{x}) = \mathbf{A}\mathbf{x}$ , where

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix}.$$

Find  $\mathbf{x} \in \mathbb{R}^2$  such that  $T(\mathbf{x}) = \begin{bmatrix} 5 \\ 7 \end{bmatrix}$ .

2.1.3.- Given the matrix and vector

$$\mathbf{A} = \begin{bmatrix} 1 & -5 & -7 \\ -3 & 7 & 5 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -2 \\ -2 \end{bmatrix},$$

define the function  $T : \mathbb{R}^3 \rightarrow \mathbb{R}^2$  as  $T(\mathbf{x}) = \mathbf{A}\mathbf{x}$ , and then find all vectors  $\mathbf{x}$  such that  $T(\mathbf{x}) = \mathbf{b}$ .

2.1.4.- Let  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be a linear transformation such that

$$T\left(\begin{bmatrix} 1 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} 3 \\ 1 \end{bmatrix}, \quad T\left(\begin{bmatrix} 0 \\ 1 \end{bmatrix}\right) = \begin{bmatrix} 1 \\ 3 \end{bmatrix}.$$

Find the values of  $T\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right)$  for any vector  $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = x_1 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + x_2 \begin{bmatrix} 0 \\ 1 \end{bmatrix} \in \mathbb{R}^2$ .

2.1.5.- Describe geometrically what is the action of  $T$  over a vector  $\mathbf{x} \in \mathbb{R}^2$ , where

$$(a) \quad T(\mathbf{x}) = \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix};$$

$$(b) \quad T(\mathbf{x}) = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix};$$

$$(c) \quad T(\mathbf{x}) = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix};$$

$$(d) \quad T(\mathbf{x}) = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}.$$

2.1.6.- Let  $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ ,  $\mathbf{v} = \begin{bmatrix} -2 \\ 5 \end{bmatrix}$ ,  $\mathbf{u} = \begin{bmatrix} 7 \\ -3 \end{bmatrix}$ , and let  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be the linear transformation  $T(\mathbf{x}) = x_1\mathbf{v} + x_2\mathbf{u}$ . Find a matrix  $\mathbf{A}$  such that  $T(\mathbf{x}) = \mathbf{A}\mathbf{x}$  for all  $\mathbf{x} \in \mathbb{R}^2$ .



## 2.2. LINEAR COMBINATIONS

One could say that the idea to introduce matrix operations originates from the interpretation of a matrix as a function. A matrix  $A \in \mathbb{F}^{m,n}$  determines a function  $A : \mathbb{F}^n \rightarrow \mathbb{F}^m$ . Such functions are generalizations of scalar valued functions of a single variable,  $f : \mathbb{R} \rightarrow \mathbb{R}$ . It is well-known how to compute the linear combination of two functions  $f, g : \mathbb{R} \rightarrow \mathbb{R}$  and, when possible, how to compute their composition and their inverses. Matrices determine a particular generalizations of scalar functions where the operations mentioned above on scalar functions can be defined on matrices. The result is called the linear combination of matrices, and when possible, the product of matrices and the inverse of a matrix. Since matrices are generalizations of scalar valued functions, there are few operations on matrices that reduce to the identity operation in the the case of scalar functions. Among such operations belong the transpose of a matrix and the trace of a matrix.

**2.2.1. Linear combination of matrices.** The addition of two matrices and the multiplication of a matrix by scalar are defined component by component.

**Definition 2.2.1.** Let  $A = [A_{ij}]$  and  $B = [B_{ij}]$  be  $m \times n$  matrices in  $\mathbb{F}^{m,n}$  and  $a, b \in \mathbb{F}$  be any scalars. The **linear combination** of  $A$  and  $B$  is also an  $m \times n$  matrix in  $\mathbb{F}^{m,n}$ , denoted as  $aA + bB$ , and given by

$$aA + bB = [aA_{ij} + bB_{ij}].$$

Recall that the notation  $aA + bB = [(aA + bB)_{ij}]$  means that the numbers  $(aA + bB)_{ij}$  are the components of the matrix  $aA + bB$ . Using this notation the definition above can be expressed in terms of matrix components as follows

$$(aA + bB)_{ij} = aA_{ij} + bB_{ij}.$$

This definition contains the particular cases  $a = b = 1$  and  $b = 0$ , given by, respectively,

$$(A + B)_{ij} = A_{ij} + B_{ij}, \quad (aA)_{ij} = aA_{ij}.$$

**EXAMPLE 2.2.1:** Consider the matrices

$$A = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \quad B = \begin{bmatrix} 3 & 0 \\ 2 & -1 \end{bmatrix}.$$

- Find the matrix  $A + B$  and  $3A$ .
- Find a matrix  $C$  such that  $2C + 6A = 4B$ .

**SOLUTION:**

**Part (a):** The definition above gives,

$$A + B = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} + \begin{bmatrix} 3 & 0 \\ 2 & -1 \end{bmatrix} \Rightarrow A + B = \begin{bmatrix} 5 & -1 \\ 1 & 1 \end{bmatrix}, \quad 3A = \begin{bmatrix} 6 & -3 \\ -3 & 6 \end{bmatrix}.$$

**Part (b):** Matrix  $C$  is given by

$$C = \frac{1}{2}(4B - 6A)$$

The definition above implies that

$$C = 2B - 3A = 2 \begin{bmatrix} 3 & 0 \\ 2 & -1 \end{bmatrix} - 3 \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} = \begin{bmatrix} 6 & 0 \\ 4 & -2 \end{bmatrix} - 3 \begin{bmatrix} 6 & -3 \\ -3 & 6 \end{bmatrix},$$

therefore, we conclude that

$$C = \begin{bmatrix} 0 & 3 \\ 7 & -8 \end{bmatrix}.$$

◁

We now summarize the main properties of the matrix linear combination.

**Theorem 2.2.2.** For all matrices  $A, B \in \mathbb{F}^{m,n}$  and all scalars  $a, b \in \mathbb{F}$ , hold:

- (a)  $(ab)A = a(bA)$ , (associativity);
- (b)  $a(A + B) = aA + aB$ , (distributivity);
- (c)  $(a + b)A = aA + bA$ , (distributivity);
- (d)  $1A = A$ , ( $1 \in \mathbb{R}$  is the identity).

The definition of linear combination of matrices is defined as the linear combination of their components, which are real or complex numbers. Therefore, all properties in Theorem 2.2.2 on linear combinations of matrices are obtained from the analogous properties on the linear combination of real or complex numbers.

**Proof of Theorem 2.2.2:** We use components notation.

(a):

$$[(ab)A]_{ij} = (ab)A_{ij} = a(bA_{ij}) = a(bA)_{ij} = [a(bA)]_{ij}.$$

(b):

$$[a(A + B)]_{ij} = a(A + B)_{ij} = a(A_{ij} + B_{ij}) = aA_{ij} + aB_{ij} = (aA)_{ij} + (aB)_{ij} = [aA + aB]_{ij}.$$

(c):

$$[(a + b)A]_{ij} = (a + b)A_{ij} = aA_{ij} + bA_{ij} = (aA)_{ij} + (bA)_{ij} = [aA + bA]_{ij}.$$

(d):

$$(1A)_{ij} = 1A_{ij} = A_{ij}.$$

□

**2.2.2. The transpose, adjoint, and trace of a matrix.** Since matrices are generalizations of scalar-valued functions, one can define operations on matrices that, unlike linear combinations, have no analogue on scalar-valued functions. One of such operations is the transpose of a matrix, which is a new matrix with the rows and columns interchanged.

**Definition 2.2.3.** The *transpose* of a matrix  $A = [A_{ij}] \in \mathbb{F}^{m,n}$  is the matrix denoted as  $A^T = [(A^T)_{kl}] \in \mathbb{F}^{n,m}$ , with its components given by

$$(A^T)_{kl} = A_{lk}.$$

**EXAMPLE 2.2.2:** Find the transpose of the  $2 \times 3$  matrix  $A = \begin{bmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \end{bmatrix}$ .

**SOLUTION:** Matrix  $A$  has components  $A_{ij}$  with  $i = 1, 2$  and  $j = 1, 2, 3$ . Therefore, its transpose has components  $(A^T)_{ji} = A_{ij}$ , that is,  $A^T$  has three rows and two columns,

$$A^T = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}.$$

◁

**EXAMPLE 2.2.3:** Show that the transpose operation satisfies  $(A^T)^T = A$ .

**SOLUTION:** The proof is:

$$[(A^T)^T]_{ij} = (A^T)_{ji} = A_{ij}.$$

An example of this property is the following: In Example 2.2.2 we showed that

$$\begin{bmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \end{bmatrix}^T = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}.$$

Therefore,

$$\left( \begin{bmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \end{bmatrix}^T \right)^T = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}^T = \begin{bmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \end{bmatrix}.$$

◁

If a matrix has complex-valued coefficients, then the complex conjugate of a matrix, also called the conjugate, is defined as the conjugate of each component.

**Definition 2.2.4.** The *complex conjugate* of a matrix  $\mathbf{A} = [A_{ij}] \in \mathbb{F}^{m,n}$  is the matrix

$$\bar{\mathbf{A}} = [\bar{A}_{ij}] \in \mathbb{F}^{m,n}.$$

**EXAMPLE 2.2.4:** Find the conjugate if matrix  $\mathbf{A} = \begin{bmatrix} 1 & 2+i \\ -i & 3-4i \end{bmatrix}$ .

**SOLUTION:** We have to conjugate each component of matrix  $\mathbf{A}$ , that is,

$$\bar{\mathbf{A}} = \begin{bmatrix} 1 & 2-i \\ i & 3+4i \end{bmatrix}.$$

◁

**EXAMPLE 2.2.5:** Show that a matrix  $\mathbf{A}$  has real coefficients iff  $\mathbf{A} = \bar{\mathbf{A}}$ ; and a matrix has purely imaginary coefficients iff  $\mathbf{A} = -\bar{\mathbf{A}}$ . Show one example of each case.

**SOLUTION:** The first condition is  $\mathbf{A} = \bar{\mathbf{A}}$ , that is,  $A_{ij} = \bar{A}_{ij}$  holds for every matrix component, which implies  $2i \operatorname{Im}(A_{ij}) = A_{ij} - \bar{A}_{ij} = 0$ . Therefore, every matrix component  $A_{ij}$  is real. The second condition is  $\mathbf{A} = -\bar{\mathbf{A}}$ , that is,  $A_{ij} = -\bar{A}_{ij}$  holds for every matrix component, which implies  $2\operatorname{Re}(A_{ij}) = A_{ij} + \bar{A}_{ij} = 0$ . Therefore, every matrix component  $A_{ij}$  is purely imaginary. Here are examples of these two situations:

$$\begin{aligned} \mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} & \Rightarrow \bar{\mathbf{A}} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = \mathbf{A}; \\ \mathbf{A} = \begin{bmatrix} i & 2i \\ 3i & 4i \end{bmatrix} & \Rightarrow \bar{\mathbf{A}} = \begin{bmatrix} -i & -2i \\ -3i & -4i \end{bmatrix} = -\mathbf{A}. \end{aligned}$$

◁

**Definition 2.2.5.** The *adjoint* of a matrix  $\mathbf{A} \in \mathbb{F}^{m,n}$  is the matrix  $\mathbf{A}^* = \bar{\mathbf{A}}^T \in \mathbb{F}^{n,m}$ .

It is not difficult to show, using components, that  $(\bar{\mathbf{A}})^T = \overline{(\mathbf{A}^T)}$ , that is, the order of the transpose and conjugate operation does not change the resulting matrix. This property is the reason why there is no parenthesis in the definition of  $\mathbf{A}^*$ .

**EXAMPLE 2.2.6:** Find the adjoint of matrix  $\mathbf{A} = \begin{bmatrix} 1 & 2+i \\ -i & 3-4i \end{bmatrix}$ .

**SOLUTION:** We need to switch rows with columns and complex conjugate the result, that is,

$$\mathbf{A}^* = \begin{bmatrix} 1 & i \\ 2-i & 3+4i \end{bmatrix}.$$

◁

The transpose, conjugate and adjoint operations are useful to specify different matrix classes having particular symmetries. These matrix classes are the symmetric, the skew-symmetric, the Hermitian, and the skew-Hermitian matrices. Here is a precise definition.

**Definition 2.2.6.** An  $n \times n$  matrix  $A$  is called:

- (a) **symmetric** iff holds  $A = A^T$ ;
- (b) **skew-symmetric** iff holds  $A = -A^T$ ;
- (c) **Hermitian** iff holds  $A = A^*$ ;
- (d) **skew-Hermitian** iff holds  $A = -A^*$ .

**EXAMPLE 2.2.7:** We present examples of each of the classes introduced in Def. 2.2.6.

**Part (a):** Matrices  $A$  and  $B$  are symmetric. Notice that  $A$  is also Hermitian, while  $B$  is not Hermitian,

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 7 & 4 \\ 3 & 4 & 8 \end{bmatrix} = A^T, \quad B = \begin{bmatrix} 1 & 2+3i & 3 \\ 2+3i & 7 & 4i \\ 3 & 4i & 8 \end{bmatrix} = B^T.$$

**Part (b):** Matrix  $C$  is skew-symmetric,

$$C = \begin{bmatrix} 0 & -2 & 3 \\ 2 & 0 & -4 \\ -3 & 4 & 0 \end{bmatrix} \Rightarrow C^T = \begin{bmatrix} 0 & 2 & -3 \\ -2 & 0 & 4 \\ 3 & -4 & 0 \end{bmatrix} = -C.$$

Notice that the diagonal elements in a skew-symmetric matrix must vanish, since  $C_{ij} = -C_{ji}$  in the case  $i = j$  means  $C_{ii} = -C_{ii}$ , that is,  $C_{ii} = 0$ .

**Part (c):** Matrix  $D$  is Hermitian but is not symmetric:

$$D = \begin{bmatrix} 1 & 2+i & 3 \\ 2-i & 7 & 4+i \\ 3 & 4-i & 8 \end{bmatrix} \Rightarrow D^T = \begin{bmatrix} 1 & 2-i & 3 \\ 2+i & 7 & 4-i \\ 3 & 4+i & 8 \end{bmatrix} \neq D,$$

however,

$$D^* = \overline{D}^T = \begin{bmatrix} 1 & 2+i & 3 \\ 2-i & 7 & 4+i \\ 3 & 4-i & 8 \end{bmatrix} = D.$$

Notice that the diagonal elements in a Hermitian matrix must be real numbers, since the condition  $A_{ij} = \overline{A_{ji}}$  in the case  $i = j$  implies  $A_{ii} = \overline{A_{ii}}$ , that is,  $2i\text{Im}(A_{ii}) = A_{ii} - \overline{A_{ii}} = 0$ . We can also verify what we said in part (a), matrix  $A$  is Hermitian since  $A^* = \overline{A}^T = A^T = A$ .

**Part (d):** The following matrix  $E$  is skew-Hermitian:

$$E = \begin{bmatrix} i & 2+i & -3 \\ -2+i & 7i & 4+i \\ 3 & -4+i & 8i \end{bmatrix} \Rightarrow E^T = \begin{bmatrix} i & -2+i & 3 \\ 2+i & 7i & -4+i \\ -3 & 4+i & 8i \end{bmatrix}$$

therefore,

$$E^* = \overline{E}^T = \begin{bmatrix} -i & -2-i & 3 \\ 2-i & -7i & -4-i \\ -3 & 4-i & -8i \end{bmatrix} = -E.$$

A skew-Hermitian matrix has purely imaginary elements in its diagonal, and the off diagonal elements have skew-symmetric real parts with symmetric imaginary parts.  $\triangleleft$

The trace of a square matrix is a number, the sum of all the diagonal matrix coefficients.

**Definition 2.2.7.** The **trace** of a square matrix  $A = [A_{ij}] \in \mathbb{F}^{n,n}$ , denoted as  $\text{tr}(A) \in \mathbb{F}$ , is the scalar given by

$$\text{tr}(A) = A_{11} + \cdots + A_{nn}.$$

**EXAMPLE 2.2.8:** Find the trace of the matrix  $A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$ .

**SOLUTION:** We add up the diagonal elements:  $\text{tr}(A) = 1 + 5 + 9$ , that is,  $\text{tr}(A) = 15$ .  $\triangleleft$

**2.2.3. Linear transformations on matrices.** The operation of computing the transpose matrix can be thought as function on the set of all  $m \times n$  matrices. The transpose operation is in fact a function  $T : \mathbb{F}^{m,n} \rightarrow \mathbb{F}^{n,m}$  given by  $T(A) = A^T$ . In a similar way, the operation of computing the trace of a square matrix can be thought as a function on the set of all square matrices. The trace is a function  $\text{tr} : \mathbb{F}^{n,n} \rightarrow \mathbb{F}$ , where  $\text{tr}(A) = A_{11} + \cdots + A_{nn}$ . One can verify that transpose function  $T$  and the trace function  $\text{tr}$  are indeed linear functions.

**Theorem 2.2.8.** *Both the transpose function  $T : \mathbb{F}^{m,n} \rightarrow \mathbb{F}^{n,m}$ , given by  $T(A) = A^T$ , and the trace function  $\text{tr} : \mathbb{F}^{n,n} \rightarrow \mathbb{F}$ , given by  $\text{tr}(A) = A_{11} + \cdots + A_{nn}$ , are linear functions.*

**Proof of Theorem 2.2.8:** We start with the transpose function  $T(A) = A^T$ , which in matrix components has the form  $(T(A))_{ij} = A_{ji}$ . Therefore, given two matrices  $A, B \in \mathbb{F}^{m,n}$  and arbitrary scalars  $a, b \in \mathbb{F}$ , holds.

$$\begin{aligned} (T(aA + bB))_{ij} &= (aA + bB)_{ji} \\ &= aA_{ji} + bB_{ji} \\ &= a(T(A))_{ij} + b(T(B))_{ij}, \end{aligned}$$

so we conclude that

$$T(aA + bB) = aT(A) + bT(B),$$

showing that the transpose operation is a linear function. We now consider the trace function. Given any two matrices all  $A, B \in \mathbb{F}^{n,n}$  and arbitrary scalars  $a, b \in \mathbb{F}$ , holds,

$$\begin{aligned} \text{tr}(aA + bB) &= \sum_{i=1}^n (aA + bB)_{ii} \\ &= \sum_{i=1}^n (aA_{ii} + bB_{ii}) \\ &= a \sum_{i=1}^n A_{ii} + b \sum_{i=1}^n B_{ii} \\ &= a \text{tr}(A) + b \text{tr}(B). \end{aligned}$$

This shows that the trace is a linear function. This establishes the Theorem.  $\square$

## 2.2.4. Exercises.

- 2.2.1.-** Construct an example of a  $3 \times 3$  matrix satisfying:
- Is symmetric and skew-symmetric.
  - Is Hermitian and symmetric.
  - Is Hermitian but not symmetric.

- 2.2.2.-** Find the numbers  $x, y, z$  solution of the equation

$$2 \begin{bmatrix} x+2 & y+3 \\ 3 & 0 \end{bmatrix} = \begin{bmatrix} 3 & 6 \\ y & z \end{bmatrix}^T$$

- 2.2.3.-** Given any square matrix  $A$  show that  $A+A^T$  is a symmetric matrix, while  $A-A^T$  is a skew-symmetric matrix.

- 2.2.4.-** Prove that there is only one way to express a matrix  $A \in \mathbb{F}^{n,n}$  as a sum of a symmetric matrix and a skew-symmetric matrix.

- 2.2.5.-** Prove the following statements:

- If  $A = [A_{ij}]$  is a skew-symmetric matrix, then holds  $A_{ii} = 0$ .
- If  $A = [A_{ij}]$  is a skew-Hermitian matrix, then the non-zero coefficients  $A_{ii}$  are purely imaginary.
- If  $A$  is a real and symmetric matrix, then  $B = iA$  is skew-Hermitian.

- 2.2.6.-** Prove that for all  $A, B \in \mathbb{F}^{m,n}$  and all  $a, b \in \mathbb{F}$  holds

$$(aA + bB)^* = \bar{a}A^* + \bar{b}B^*.$$

- 2.2.7.-** Prove that the transpose function  $T : \mathbb{F}^{m,n} \rightarrow \mathbb{F}^{n,m}$  and trace function  $\text{tr} : \mathbb{F}^{n,n} \rightarrow \mathbb{F}$  are linear functions.

## 2.3. MATRIX MULTIPLICATION

The operation of matrix multiplication originates in the composition of functions. We call this operation a multiplication since it reduces to the multiplication of real numbers in the case of  $1 \times 1$  real matrices. Unlike the multiplication of real numbers, the product of general matrices is not commutative, that is,  $AB \neq BA$  in the general case. This property reflects the fact that the composition of two functions is a non-commutative operation. In this Subsection we first introduce the multiplication of two matrices using the matrix vector product. We then introduce the formula for the components of the product of two matrices. Finally we show that the composition of two matrices is their matrix product.

**2.3.1. Algebraic definition.** Matrix multiplication is defined using the matrix-vector product. The product of two matrices is the matrix-vector product of one matrix with each column of the other matrix.

**Definition 2.3.1.** The *matrix multiplication* of the  $m \times n$  matrix  $A$  with the  $n \times \ell$  matrix  $B = [B_{:1}, \dots, B_{:\ell}]$ , denoted by  $AB$ , is the  $m \times \ell$  matrix given by  $AB = [AB_{:1}, \dots, AB_{:\ell}]$ .

The product is not defined for two arbitrary matrices, since the size of the matrices is important: The numbers of columns in the first matrix must match the numbers of rows in the second matrix,

$$\begin{array}{ccccc} A & \text{times} & B & \text{defines} & AB \\ m \times n & & n \times \ell & & m \times \ell \end{array}$$

We assign a name to matrices satisfying this property.

**Definition 2.3.2.** Matrices  $A$  and  $B$  are called *conformable* in the order  $A, B$ , iff the product  $AB$  is well defined.

**EXAMPLE 2.3.1:** Compute the product of the matrices  $A$  and  $B$  below, which are conformable in both orders  $AB$  and  $BA$ , where

$$A = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}, \quad B = \begin{bmatrix} 3 & 0 \\ 2 & -1 \end{bmatrix}.$$

**SOLUTION:** Following the definition above we compute the product in the order  $AB$ , namely,

$$AB = [AB_{:1}, AB_{:2}] = \left[ A \begin{bmatrix} 3 \\ 2 \end{bmatrix}, A \begin{bmatrix} 0 \\ -1 \end{bmatrix} \right] = \left[ \begin{bmatrix} 6-2 \\ -3+4 \end{bmatrix}, \begin{bmatrix} 0+1 \\ 0-2 \end{bmatrix} \right] \Rightarrow AB = \begin{bmatrix} 4 & 1 \\ 1 & -2 \end{bmatrix}.$$

Using the same definition we can compute the product in the opposite order, that is,

$$BA = [BA_{:1}, BA_{:2}] = \left[ B \begin{bmatrix} 2 \\ -1 \end{bmatrix}, B \begin{bmatrix} -1 \\ 2 \end{bmatrix} \right] = \left[ \begin{bmatrix} 6-0 \\ 4+1 \end{bmatrix}, \begin{bmatrix} -3+0 \\ -2-2 \end{bmatrix} \right] \Rightarrow BA = \begin{bmatrix} 6 & -3 \\ 5 & -4 \end{bmatrix}.$$

This is an example where we have that  $AB \neq BA$ . ◁

The following result gives a formula to compute the components of the product matrix in terms of the components of the individual matrices.

**Theorem 2.3.3.** Consider the  $m \times n$  matrix  $A = [A_{ij}]$  and the  $n \times \ell$  matrix  $B = [B_{jk}]$ , where the indices take values as follows:  $i = 1, \dots, m$ ,  $j = 1, \dots, n$  and  $k = 1, \dots, \ell$ . The components of the product matrix  $AB$  are given by

$$(AB)_{ik} = \sum_{j=1}^n A_{ij} B_{jk}. \quad (2.4)$$

We recall that the symbol  $\sum_{j=1}^n$  in Eq. (2.4) means to add up all the terms having the index  $j$  starting from  $j = 1$  until  $j = n$ , that is,

$$\sum_{j=1}^n A_{ij}B_{jk} = A_{i1}B_{1k} + A_{i2}B_{2k} + \cdots + A_{in}B_{nk}.$$

**Proof of Theorem 2.3.3:** The column  $k$  of the product  $\mathbf{AB}$  is given by the column vector  $\mathbf{AB}_{:k}$ . This is a vector with  $m$  components,

$$\mathbf{AB}_{:k} = \begin{bmatrix} \sum_{j=1}^n A_{1j}B_{jk} \\ \vdots \\ \sum_{j=1}^n A_{mj}B_{jk} \end{bmatrix},$$

therefore the  $i$ -th component of this vector is given by

$$(\mathbf{AB})_{ik} = \sum_{j=1}^n A_{ij}B_{jk}.$$

This establishes the Theorem. □

**EXAMPLE 2.3.2:** We now use Eq. (2.4) to find the product of matrices  $\mathbf{A}$  and  $\mathbf{B}$  in Example 2.3.1. The component  $(\mathbf{AB})_{11} = 4$  is obtained from the first row in matrix  $\mathbf{A}$  and the first column in matrix  $\mathbf{B}$  as follows:

$$\begin{bmatrix} \mathbf{2} & \mathbf{-1} \\ \mathbf{-1} & \mathbf{2} \end{bmatrix} \begin{bmatrix} \mathbf{3} & \mathbf{0} \\ \mathbf{2} & \mathbf{-1} \end{bmatrix} = \begin{bmatrix} \mathbf{4} & \mathbf{1} \\ \mathbf{1} & \mathbf{-2} \end{bmatrix}, \quad (2)(3) + (-1)(2) = 4;$$

The component  $(\mathbf{AB})_{12} = -1$  is obtained as follows:

$$\begin{bmatrix} \mathbf{2} & \mathbf{-1} \\ \mathbf{-1} & \mathbf{2} \end{bmatrix} \begin{bmatrix} \mathbf{3} & \mathbf{0} \\ \mathbf{2} & \mathbf{-1} \end{bmatrix} = \begin{bmatrix} \mathbf{4} & \mathbf{1} \\ \mathbf{1} & \mathbf{-2} \end{bmatrix}, \quad (2)(0) + (-1)(1) = -1;$$

The component  $(\mathbf{AB})_{21} = 1$  is obtained as follows:

$$\begin{bmatrix} \mathbf{2} & \mathbf{-1} \\ \mathbf{-1} & \mathbf{2} \end{bmatrix} \begin{bmatrix} \mathbf{3} & \mathbf{0} \\ \mathbf{2} & \mathbf{-1} \end{bmatrix} = \begin{bmatrix} \mathbf{4} & \mathbf{1} \\ \mathbf{1} & \mathbf{-2} \end{bmatrix}, \quad (-1)(3) + (2)(2) = 1;$$

And finally the component  $(\mathbf{AB})_{22} = -2$  is obtained as follows:

$$\begin{bmatrix} \mathbf{2} & \mathbf{-1} \\ \mathbf{-1} & \mathbf{2} \end{bmatrix} \begin{bmatrix} \mathbf{3} & \mathbf{0} \\ \mathbf{2} & \mathbf{-1} \end{bmatrix} = \begin{bmatrix} \mathbf{4} & \mathbf{1} \\ \mathbf{1} & \mathbf{-2} \end{bmatrix}, \quad (-1)(0) + (2)(-1) = -2. \quad \triangleleft$$

We have seen in Example 2.3.1 that the matrix product is not commutative, since in that example  $\mathbf{AB} \neq \mathbf{BA}$ . In that example the matrices were conformable in both orders  $\mathbf{AB}$  and  $\mathbf{BA}$ , although their products do not match. It can also be possible that two matrices  $\mathbf{A}$  and  $\mathbf{B}$  are conformable in the order  $\mathbf{AB}$  but they are not conformable in the opposite order. That is, the matrix product is possible in one order but not in the other order.

**EXAMPLE 2.3.3:** Consider the matrices

$$\mathbf{A} = \begin{bmatrix} 4 & 3 \\ 2 & 1 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$$

These matrices are conformable in the order  $\mathbf{AB}$  but not in the order  $\mathbf{BA}$ . In the first case we obtain

$$\mathbf{AB} = \begin{bmatrix} 4 & 3 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \Rightarrow \mathbf{AB} = \begin{bmatrix} 16 & 23 & 30 \\ 6 & 9 & 12 \end{bmatrix}.$$

The product  $\mathbf{BA}$  is not possible. □



**EXAMPLE 2.3.4:** Column vectors and row vectors are particular cases of matrices, they are  $n \times 1$  and  $1 \times n$  matrices, respectively. We denote row vectors as transpose of a column vector. Using this notation and the matrix product, compute both products  $\mathbf{v}^T \mathbf{u}$  and  $\mathbf{u} \mathbf{v}^T$ , where the vectors are given by

$$\mathbf{u} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}, \quad \mathbf{v} = \begin{bmatrix} 5 \\ 1 \end{bmatrix}.$$

**SOLUTION:** In the first case we multiply the matrices  $1 \times 2$  and  $2 \times 1$ , so the result is a  $1 \times 1$  matrix, a real number, given by,

$$\mathbf{v}^T \mathbf{u} = [5 \quad 1] \begin{bmatrix} 2 \\ 3 \end{bmatrix} \Rightarrow \mathbf{v}^T \mathbf{u} = 13.$$

In the second case we multiply the matrices  $2 \times 1$  and  $1 \times 2$ , so the result is a  $2 \times 2$  matrix,

$$\mathbf{u} \mathbf{v}^T = \begin{bmatrix} 2 \\ 3 \end{bmatrix} [5 \quad 1] \Rightarrow \mathbf{u} \mathbf{v}^T = \begin{bmatrix} 10 & 2 \\ 15 & 3 \end{bmatrix}.$$

◁

It is well-known that the product of two numbers is zero, then one of them must be zero. This property is not true in the case of matrix multiplication, as it can be seen below.

**EXAMPLE 2.3.5:** Compute the product  $\mathbf{AB}$  where  $\mathbf{A} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$  and  $\mathbf{B} = \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix}$ .

**SOLUTION:** It is simple to check that

$$\mathbf{AB} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \Rightarrow \mathbf{AB} = \mathbf{0}.$$

The product is the zero matrix, although  $\mathbf{A} \neq \mathbf{0}$  and  $\mathbf{B} \neq \mathbf{0}$ .

◁

**2.3.2. Matrix composition.** The product of two matrices originates in the composition of the linear functions defined by the matrices. This can be seen in the following result.

**Theorem 2.3.4.** *Given the  $m \times n$  matrix  $\mathbf{A} : \mathbb{F}^n \rightarrow \mathbb{F}^m$  and the  $n \times \ell$  matrix  $\mathbf{B} : \mathbb{F}^\ell \rightarrow \mathbb{F}^n$ , their composition is a function*

$$\mathbf{A} \circ \mathbf{B} : \mathbb{F}^\ell \xrightarrow{\mathbf{B}} \mathbb{F}^n \xrightarrow{\mathbf{A}} \mathbb{F}^m,$$

which is an  $m \times \ell$  matrix given by the matrix product of  $\mathbf{A}$  and  $\mathbf{B}$ , that is,  $\mathbf{A} \circ \mathbf{B} = \mathbf{AB}$ .

**Proof of Theorem 2.3.4:** The composition of the function  $\mathbf{A}$  and  $\mathbf{B}$  is defined for all  $\mathbf{x} \in \mathbb{F}^\ell$  as follows

$$(\mathbf{A} \circ \mathbf{B})\mathbf{x} = \mathbf{A}(\mathbf{B}\mathbf{x}).$$

Introduce the usual notation  $\mathbf{B} = [\mathbf{B}_{:1}, \dots, \mathbf{B}_{:\ell}]$ . Then, the composition  $\mathbf{A} \circ \mathbf{B}$  can be re-expressed as follows,

$$(\mathbf{A} \circ \mathbf{B})\mathbf{x} = \mathbf{A} \left( [\mathbf{B}_{:1}, \dots, \mathbf{B}_{:\ell}] \begin{bmatrix} x_1 \\ \vdots \\ x_\ell \end{bmatrix} \right) = \mathbf{A}(\mathbf{B}_{:1}x_1 + \dots + \mathbf{B}_{:\ell}x_\ell).$$

Since the matrix-vector product is a linear operation, we get

$$(\mathbf{A} \circ \mathbf{B})\mathbf{x} = \mathbf{A}(\mathbf{B}_{:1}x_1) + \dots + \mathbf{A}(\mathbf{B}_{:\ell}x_\ell) = (\mathbf{AB}_{:1})x_1 + \dots + (\mathbf{AB}_{:\ell})x_\ell$$

where the second equation comes from the definition of matrix multiplication. Then, using again the definition of matrix-vector product,

$$(A \circ B)\mathbf{x} = [AB_{:1}, \dots, AB_{:\ell}] \begin{bmatrix} x_1 \\ \vdots \\ x_\ell \end{bmatrix} \Rightarrow (A \circ B)\mathbf{x} = (AB)\mathbf{x}.$$

This establishes the Theorem.  $\square$

**EXAMPLE 2.3.6:** Find the matrix  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  that produces a rotation by an angle  $\theta_1$  counterclockwise and then another rotation by an angle  $\theta_2$  counterclockwise.

**SOLUTION:** Let us denote by  $R_\theta$  the matrix that performs a rotation on the plane by an angle  $\theta$  counterclockwise, that is,

$$R_\theta = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}.$$

The matrix  $T$  is the the composition  $T = R_{\theta_2} \circ R_{\theta_1}$ . Since the matrix of a composition of these two rotations can be obtained computing the matrix product of them, then we obtain

$$T = R_{\theta_2} R_{\theta_1} = \begin{bmatrix} \cos(\theta_2) & -\sin(\theta_2) \\ \sin(\theta_2) & \cos(\theta_2) \end{bmatrix} \begin{bmatrix} \cos(\theta_1) & -\sin(\theta_1) \\ \sin(\theta_1) & \cos(\theta_1) \end{bmatrix}.$$

The product above is given by

$$T = \begin{bmatrix} \cos(\theta_2) \cos(\theta_1) - \sin(\theta_2) \sin(\theta_1) & -[\sin(\theta_1) \cos(\theta_2) + \sin(\theta_2) \cos(\theta_1)] \\ \sin(\theta_1) \cos(\theta_2) + \sin(\theta_2) \cos(\theta_1) & -\sin(\theta_2) \sin(\theta_1) + \cos(\theta_2) \cos(\theta_1) \end{bmatrix}.$$

The formulas

$$\begin{aligned} \cos(\theta_1 + \theta_2) &= \cos(\theta_2) \cos(\theta_1) - \sin(\theta_2) \sin(\theta_1) \\ \sin(\theta_1 + \theta_2) &= \sin(\theta_1) \cos(\theta_2) + \sin(\theta_2) \cos(\theta_1), \end{aligned}$$

imply that

$$T = \begin{bmatrix} \cos(\theta_1 + \theta_2) & -\sin(\theta_1 + \theta_2) \\ \sin(\theta_1 + \theta_2) & \cos(\theta_1 + \theta_2) \end{bmatrix}.$$

Notice that we have obtained a result that it is intuitively clear: Two consecutive rotations on the plane is equivalent to a single rotation by an angle that is the sum of the individual rotations, that is,

$$R_{\theta_2} R_{\theta_1} = R_{\theta_2 + \theta_1}.$$

In particular, notice that the matrix multiplication when restricted to the set of all rotation matrices on the plane is a commutative operation, that is,

$$R_{\theta_2} R_{\theta_1} = R_{\theta_1} R_{\theta_2}, \quad \theta_1, \theta_2 \in \mathbb{R}.$$

$\triangleleft$

In Examples 2.3.7 and 2.3.8 below we show that the order of the functions in a composition change the resulting function. This is essentially the reason behind the non-commutativity of the matrix multiplication.

**EXAMPLE 2.3.7:** Find the matrix  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  that first performs a rotation by an angle  $\pi/2$  counterclockwise and then performs a reflection along the  $x_1 = x_2$  line on the plane.

**SOLUTION:** The matrix  $T$  is the composition the rotation  $R_{\pi/2}$  with the reflection function  $A$ , given by

$$R_{\pi/2} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}, \quad A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}.$$

The matrix  $T$  is given by

$$T = AR_{\pi/2} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \Rightarrow T = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}.$$

Therefore, the function  $T$  is a reflection along the horizontal line  $x_2 = 0$ .  $\triangleleft$

**EXAMPLE 2.3.8:** Find the matrix  $S : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  that first performs a reflection along the  $x_1 = x_2$  line on the plane and then performs a rotation by an angle  $\pi/2$  counterclockwise.

**SOLUTION:** The matrix  $S$  is the composition the reflection function  $A$  and then the rotation  $R_{\pi/2}$ , given in the previous Example 2.3.7. The matrix  $S$  is the given by

$$S = R_{\pi/2}A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \Rightarrow S = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}.$$

Therefore, the function  $S$  is a reflection along the vertical line  $x_1 = 0$ .  $\triangleleft$

**2.3.3. Main properties.** We summarize the main properties of the matrix multiplication.

**Theorem 2.3.5.** *The following properties hold for all  $m \times n$  matrix  $A$ ,  $n \times m$  matrix  $B$ ,  $n \times \ell$  matrices  $C$ ,  $D$ , and  $\ell \times k$  matrix  $E$ :*

- (a)  $AB \neq BA$  in the general case, so the product is non-commutative;
- (b)  $A(C + D) = AC + AD$ , and  $(C + D)E = CE + DE$ ;
- (c)  $A(CE) = (AC)E$ , associativity;
- (d)  $I_m A = A I_n = A$ ;
- (e)  $(AC)^T = C^T A^T$ , and  $(AC)^* = C^* A^*$ .
- (f)  $\text{tr}(AB) = \text{tr}(BA)$ .

**Proof of Theorem 2.3.5:**

**Part (a):** When  $m \neq n$  the matrix  $AB$  is  $m \times m$  while  $BA$  is  $n \times n$ , so they cannot be equal. When  $m = n$ , the matrices in Example 2.3.4 and 2.3.5 show that this product is not commutative.

**Part (b):** This property can be shown as follows:

$$\begin{aligned} A(C + D) &= A([\text{C}_{:1}, \dots, \text{C}_{:\ell}] + [\text{D}_{:1}, \dots, \text{D}_{:\ell}]) \\ &= A([\text{C}_{:1} + \text{D}_{:1}], \dots, [\text{C}_{:\ell} + \text{D}_{:\ell}]) \\ &= [A(\text{C}_{:1} + \text{D}_{:1}), \dots, A(\text{C}_{:\ell} + \text{D}_{:\ell})] \\ &= [(AC)_{:1} + (AD)_{:1}], \dots, [(AC)_{:\ell} + (AD)_{:\ell}] \\ &= [AC_{:1}, \dots, AC_{:\ell}] + [AD_{:1}, \dots, AD_{:\ell}] \\ &= AC + AD. \end{aligned}$$

The other equation is proven in a similar way.

**Part (c):** This property is proven using the component expression for the matrix product:

$$[A(CE)]_{ij} = \sum_{k=1}^n A_{ik}(CE)_{kj} = \sum_{k=1}^n A_{ik} \left( \sum_{l=1}^{\ell} C_{kl} E_{lj} \right);$$

however, the order of the sums can be interchanged,

$$\sum_{k=1}^n A_{ik} \left( \sum_{l=1}^{\ell} C_{kl} E_{lj} \right) = \sum_{l=1}^{\ell} \left( \sum_{k=1}^n A_{ik} C_{kl} \right) E_{lj};$$

So, from this last expression is not difficult to show that:

$$\sum_{l=1}^{\ell} \left( \sum_{k=1}^n A_{ik} C_{kl} \right) E_{lj} = \sum_{l=1}^{\ell} (AC)_{il} E_{lj} = [(AC)E]_{ij};$$

We have just proved that

$$[A(CE)]_{ij} = [(AC)E]_{ij}.$$

**Part (d):** We can use components again, recalling that the components of  $I_m$  are given by  $(I_m)_{ij} = 0$  if  $i \neq j$  and is  $(I_m)_{ii} = 1$ . Therefore,

$$(I_m A)_{ij} = \sum_{k=1}^m (I_m)_{ik} A_{kj} = (I_m)_{ii} A_{ij} = A_{ij}.$$

Analogously,

$$(A I_n)_{ij} = \sum_{k=1}^n A_{ik} (I_n)_{kj} = A_{ij} (I_n)_{jj} = A_{ij}.$$

**Part (e):** Use components once more:

$$[(AC)^T]_{ij} = (AC)_{ji} = \sum_{k=1}^n A_{jk} C_{ki} = \sum_{k=1}^n (A^T)_{kj} (C^T)_{ik} = \sum_{k=1}^n (C^T)_{ik} (A^T)_{kj} = [C^T A^T]_{ij}.$$

The second equation follows from the proof above and the property of the complex conjugate:

$$\overline{(AC)} = \overline{A} \overline{C}.$$

Indeed

$$(AC)^* = \overline{(AC)^T} = \overline{C^T A^T} = \overline{C^T} \overline{A^T} = C^* A^*.$$

**Part (f):** Recall that the trace of a matrix  $A$  is given by

$$\text{tr}(A) = A_{11} + \cdots + A_{nn} = \sum_{i=1}^n A_{ii}.$$

Then it is simple to see that

$$\text{tr}(AB) = \sum_{i=1}^m \left( \sum_{j=1}^n A_{ij} B_{ji} \right) = \sum_{i=1}^m \left( \sum_{j=1}^n B_{ji} A_{ij} \right) = \sum_{j=1}^n \left( \sum_{i=1}^m B_{ji} A_{ij} \right) = \text{tr}(BA).$$

This establishes the Theorem.  $\square$

We use the notation  $A^2 = AA$ ,  $A^3 = A^2A$ , and  $A^n = A^{n-1}A$ . Notice that the matrix product is not commutative, so the formula  $(a+b)^2 = a^2 + 2ab + b^2$  does not hold for matrices. Instead we have:

$$(A+B)^2 = (A+B)(A+B) = A^2 + AB + BA + B^2.$$

**2.3.4. Block multiplication.** The multiplication of two large matrices can be simplified in the case that each matrix can be subdivided in appropriate blocks. If these matrix blocks are conformable the multiplication of the original matrices reduces to the multiplication of the smaller matrix blocks. The next result presents a simple case.

**Theorem 2.3.6.** *If  $A$  is an  $m \times n$  matrix and  $B$  is an  $n \times \ell$  matrix having the following block decomposition,*

$$A = \begin{bmatrix} A_{11} & \vdots & A_{12} \\ \dots & \dots & \dots \\ A_{21} & \vdots & A_{22} \end{bmatrix}, \quad A = \begin{bmatrix} m_1 \times n_1 & \vdots & m_1 \times n_2 \\ \dots & \dots & \dots \\ m_2 \times n_1 & \vdots & m_2 \times n_2 \end{bmatrix},$$

$$B = \begin{bmatrix} B_{11} & \vdots & B_{12} \\ \dots & \dots & \dots \\ B_{21} & \vdots & B_{22} \end{bmatrix}, \quad B = \begin{bmatrix} n_1 \times \ell_1 & \vdots & n_1 \times \ell_2 \\ \dots & \dots & \dots \\ n_2 \times \ell_1 & \vdots & n_2 \times \ell_2 \end{bmatrix},$$

where  $m_1 + m_2 = m$ ,  $n_1 + n_2 = n$  and  $\ell_1 + \ell_2 = \ell$ , then the product  $AB$  has the form

$$AB = \begin{bmatrix} A_{11}B_{11} + A_{12}B_{21} & \vdots & A_{11}B_{12} + A_{12}B_{22} \\ \dots & \dots & \dots \\ A_{21}B_{11} + A_{22}B_{21} & \vdots & A_{21}B_{12} + A_{22}B_{22} \end{bmatrix}, \quad AB = \begin{bmatrix} m_1 \times \ell_1 & \vdots & m_1 \times \ell_2 \\ \dots & \dots & \dots \\ m_2 \times \ell_1 & \vdots & m_2 \times \ell_2 \end{bmatrix}.$$

The proof is a straightforward computation, so we omit it. This type of block decomposition is useful when several blocks are repeated inside the matrices  $A$  and  $B$ . This situation appears in the following example.

**EXAMPLE 2.3.9:** Use block multiplication to find the matrix  $AB$ , where

$$A = \begin{bmatrix} 1 & 2 & 1 & 0 \\ 3 & 4 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 0 & 1 & 2 \\ 0 & 1 & 3 & 4 \\ 0 & 0 & 1 & 2 \\ 0 & 0 & 3 & 4 \end{bmatrix}.$$

**SOLUTION:** These matrices have the following block structure:

$$A = \begin{bmatrix} 1 & 2 & \vdots & 1 & 0 \\ 3 & 4 & \vdots & 0 & 1 \\ \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & \vdots & 0 & 0 \\ 0 & 1 & \vdots & 0 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 0 & \vdots & 1 & 2 \\ 0 & 1 & \vdots & 3 & 4 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \vdots & 1 & 2 \\ 0 & 0 & \vdots & 3 & 4 \end{bmatrix},$$

so, introduce the matrices

$$I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad C = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}, \quad 0 = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix},$$

then, the original matrices have the block form

$$A = \begin{bmatrix} C & \vdots & I \\ \dots & \dots & \dots \\ I & \vdots & 0 \end{bmatrix}, \quad B = \begin{bmatrix} I & \vdots & C \\ \dots & \dots & \dots \\ 0 & \vdots & C \end{bmatrix}.$$

Then, the matrix  $AB$  has the form

$$AB = \begin{bmatrix} C & \vdots & C^2 + C \\ \dots & \dots & \dots \\ I & \vdots & C \end{bmatrix}.$$

So, the only calculation we need to do is the matrix  $C^2 + C$ , which is given by

$$C^2 + C = \begin{bmatrix} 8 & 12 \\ 18 & 26 \end{bmatrix}.$$

If we put all the information above together, we get

$$AB = \begin{bmatrix} 1 & 2 & \vdots & 8 & 12 \\ 3 & 4 & \vdots & 18 & 26 \\ \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & \vdots & 1 & 2 \\ 0 & 1 & \vdots & 3 & 4 \end{bmatrix} \Leftrightarrow AB = \begin{bmatrix} 1 & 2 & 8 & 12 \\ 3 & 4 & 18 & 26 \\ 1 & 0 & 1 & 2 \\ 0 & 1 & 3 & 4 \end{bmatrix}.$$

◁

**EXAMPLE 2.3.10:** Given arbitrary matrices  $A$ , that is  $n \times k$ , and  $B$ , that is  $k \times n$ , show that the  $(k+n) \times (k+n)$  matrix  $C$  below satisfies  $C^2 = I_{n+k}$ , where

$$C = \begin{bmatrix} I_k - BA & \vdots & B \\ \dots & \dots & \dots \\ 2A - ABA & \vdots & AB - I_n \end{bmatrix}$$

**SOLUTION:** Notice that  $AB$  is an  $n \times n$  matrix, while  $BA$  is an  $k \times k$  matrix, so the definition of  $C$  implies that

$$C = \begin{bmatrix} k \times k & \vdots & k \times n \\ \dots & \dots & \dots \\ n \times k & \vdots & n \times n \end{bmatrix}.$$

Using block multiplication we obtain:

$$C^2 = \begin{bmatrix} (I_k - BA) & \vdots & B \\ \dots & \dots & \dots \\ (2A - ABA) & \vdots & (AB - I_n) \end{bmatrix} \begin{bmatrix} (I_k - BA) & \vdots & B \\ \dots & \dots & \dots \\ (2A - ABA) & \vdots & (AB - I_n) \end{bmatrix} =$$

$$\begin{bmatrix} (I_k - BA)(I_k - BA) + B(2A - ABA) & \vdots & (I_k - BA)B + B(AB - I_n) \\ \dots & \dots & \dots \\ (2A - ABA)(I_k - BA) + (AB - I_n)(2A - ABA) & \vdots & (2A - ABA)B + (AB - I_n)(AB - I_n) \end{bmatrix}$$

Now, notice that the block (1, 1) above is:

$$(I_k - BA)(I_k - BA) + B(2A - ABA) = I_k - BA - BA + BABA + 2BA - BABA = I_k.$$

The block (1, 2) is:

$$(I_k - BA)B + B(AB - I_n) = B - BAB + BAB - B = 0.$$

The block (2, 1) is:

$$(2A - ABA)(I_k - BA) + (AB - I_n)(2A - ABA) =$$

$$2A - 2ABA - ABA + ABABA + 2ABA - ABABA - 2A + ABA = 0.$$

Finally, the block (2, 2) is:

$$(2A - ABA)B + (AB - I_n)(AB - I_n) = 2AB - ABAB + ABAB - AB - AB + I_n = I_n.$$

Therefore, we have shown that:  $C^2 = \begin{bmatrix} I_k & \vdots & 0 \\ \dots & \dots & \dots \\ 0 & \vdots & I_n \end{bmatrix} = I_{k+n}$ . ◁

**2.3.5. Matrix commutators.** Matrix multiplication is in general not commutative. Given two  $n \times n$  matrices  $A$  and  $B$ , their product  $AB$  is in general different from  $BA$ . We call the difference between these two matrices the commutator of  $A$  and  $B$ .

**Definition 2.3.7.** The *commutator* of two  $n \times n$  matrices  $A$  and  $B$ , is the matrix

$$[A, B] = AB - BA.$$

Furthermore, we say that the square matrices  $A, B$  *commute* iff holds  $[A, B] = 0$ .

Therefore, the symbol  $[A, B]$  denotes an  $n \times n$  matrix, the commutator of  $A$  and  $B$ . The commutator of two operators defined on an inner product space is an important concept in quantum mechanics. When we say an operator, we should think in something like a matrix, and when we say an inner product space, we should think in something like  $\mathbb{F}^n$ . An observation of a physical property, like position or momentum, of a quantum mechanical system, like an electron, is described by an operator acting on a vector, the latter describing the state of the quantum mechanical system. By observing nature, one discovers that two properties can be simultaneously measured without limit in the measurement precision iff the associated operators commute. If the operators do not commute, then their commutator determines the maximum precision of their simultaneous observation. This is the physical content of the *Heisenberg uncertainty principle*.

**EXAMPLE 2.3.11:** We have seen in Example 2.3.6 that two rotation matrices  $R(\theta_1)$  and  $R(\theta_2)$  commute for all  $\theta_1, \theta_2 \in \mathbb{R}$ , that is,

$$[R(\theta_1), R(\theta_2)] = R(\theta_1)R(\theta_2) - R(\theta_2)R(\theta_1) \Rightarrow [R(\theta_1), R(\theta_2)] = 0. \quad \triangleleft$$

**Theorem 2.3.8.** For all matrices  $A, B, C \in \mathbb{F}^{n,n}$  and scalars  $a, b, c \in \mathbb{F}$  holds:

- (a)  $[A, B] = -[B, A]$ , (antisymmetry);
- (b)  $[aA, bB] = ab[A, B]$ , (linearity);
- (c)  $[A, B + C] = [A, B] + [A, C]$ , (linearity on the right entry);
- (d)  $[A + B, C] = [A, C] + [B, C]$ , (linearity on the left entry);
- (e)  $[A, BC] = [A, B]C + B[A, C]$ , (right derivation property);
- (f)  $[AB, C] = [A, C]B + A[B, C]$ , (left derivation property);
- (g)  $[[A, B], C] + [[C, A], B] + [[B, C], A] = 0$ , (Jacobi property).

**Proof of Theorem 2.3.8:** All properties are simple to show.

Part (a):

$$[A, B] = AB - BA = -(BA - AB) = -[B, A].$$

Part (b):

$$[aA, bB] = aAbB - bBaA = ab(AB - BA) = ab[A, B].$$

Part (c):

$$[A, B + C] = A(B + C) - (B + C)A = AB + AC - BA - CA = [A, B] + [A, C].$$

Part (d): The proof is similar to the previous calculation.

Part (e):

$$[A, BC] = A(BC) - (BC)A = ABC + (BAC - BAC) - BCA = [A, B]C + B[A, C].$$

Part (f): The proof is similar to the previous calculation.

**Part (g):** We write down each term and we verify that they all add up to zero:

$$[[A, B], C] = (AB - BA)C - C(AB - BA);$$

$$[[C, A], B] = (CA - AC)B - B(CA - AC);$$

$$[[B, C], A] = (BC - CB)A - A(BC - CB).$$

These three equations add up to the zero matrix. This establishes the Theorem.  $\square$

We finally highlight that these properties implies that  $[A, A^m] = 0$  holds for all  $m \in \mathbb{N}$ .

**Further reading.** Almost every book in linear algebra explains matrix multiplication. See Section 3.5 in Meyer's book [3], and Section 3.6 for block multiplication. Also Strang's book [4]. The definition of commutators can be found in Section 2.3.1 in Hassani's book [1]. Commutators play an important role in the Spectral Theorem, which we study in Chapter 9.



## 2.3.6. Exercises.

2.3.1.- Given the matrices

$$A = \begin{bmatrix} 1 & -2 & 3 \\ 0 & -5 & 4 \\ 4 & -3 & 8 \end{bmatrix}, B = \begin{bmatrix} 1 & 2 \\ 0 & 4 \\ 3 & 7 \end{bmatrix},$$

and the vector  $C = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$ , compute the

following products, if possible:

- (a)  $AB$ ,  $BA$ ,  $CB$  and  $C^T B$ .  
 (b)  $A^2$ ,  $B^2$ ,  $C^T C$  and  $CC^T$ .

2.3.2.- Consider the matrices

$$A = \begin{bmatrix} 2 & 1 \\ 3 & 1 \end{bmatrix}, B = \begin{bmatrix} 1 & 1 \\ 3 & 0 \end{bmatrix}, C = \begin{bmatrix} 1 & 4 \\ 2 & 3 \end{bmatrix}.$$

- (a) Compute  $[A, B]$ .  
 (b) Find the product  $ABC$ .

2.3.3.- Find  $A^2$  and  $A^3$  for the matrix

$$A = \begin{bmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}.$$

2.3.4.- Given a real number  $a$  find the matrix  $A^n$ , where  $n$  is any positive integer and

$$A = \begin{bmatrix} 1 & a \\ 0 & 1 \end{bmatrix}.$$

2.3.5.- Given any square matrices  $A$ ,  $B$ , prove that

$$(A+B)^2 = A^2 + 2AB + B^2 \Leftrightarrow [A, B] = 0.$$

2.3.6.- Given  $a = 1/3$ , divide the matrix

$$A = \begin{bmatrix} 1 & 0 & 0 & a & a & a \\ 0 & 1 & 0 & a & a & a \\ 0 & 0 & 1 & a & a & a \\ 0 & 0 & 0 & a & a & a \\ 0 & 0 & 0 & a & a & a \\ 0 & 0 & 0 & a & a & a \end{bmatrix}$$

into appropriate blocks, and using block multiplication find the matrix  $A^{300}$ .2.3.7.- Prove that for all matrices  $A \in \mathbb{F}^{n,m}$  and  $B \in \mathbb{F}^{m,n}$  holds

$$\text{tr}(AB) = \text{tr}(BA).$$

2.3.8.- Let  $A$  be an  $m \times n$  matrix. Show that  $\text{tr}(A^T A) = 0$  iff  $A = 0$ .2.3.9.- Prove that: If  $A, B \in \mathbb{F}^{n,n}$  are symmetric matrices and commute, then their product  $AB$  is also a symmetric matrix.2.3.10.- Let  $A$  be an arbitrary  $n \times n$  matrix. Use the trace function to show that there exists no  $n \times n$  matrix  $X$  solution of the matrix equation

$$[A, X] = I_n.$$

## 2.4. INVERSE MATRIX

In this Section we introduce the concept of the inverse of a square matrix. Not every square matrix is invertible. In the case of  $2 \times 2$  matrices we present a condition on the matrix coefficients that is equivalent to the invertibility of the matrix, and we also present a formula for the inverse matrix. Later on in Section 3.2 we generalize this formula for  $n \times n$  matrices. The inverse of a matrix is useful to compute solutions to systems of linear equations.

**2.4.1. Main definition.** We start recalling that the matrix  $I_n \in \mathbb{F}^{n,n}$  is called the *identity matrix* iff holds that  $I_n \mathbf{x} = \mathbf{x}$  for all  $\mathbf{x} \in \mathbb{F}^n$ . Choosing the vector  $\mathbf{x} = \mathbf{e}_i$ , consisting of a 1 in the row  $i$  and zero every where else, it is simple to see that the components of the identity matrix are given by

$$I_n = [I_{ij}] \quad \text{with} \quad \begin{cases} I_{ii} = 1, \\ I_{ij} = 0, \quad i \neq j. \end{cases}$$

The cases  $n = 2$  and  $n = 3$  are shown below,

$$I_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad I_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

We also use the column vector notation for the identity matrix,

$$I_n = [\mathbf{e}_1, \dots, \mathbf{e}_n],$$

that is, we denote the identity column vectors  $I_{\cdot i} = \mathbf{e}_i$  for  $i = 1, \dots, n$ . For example, in the case of  $I_3$  we have that

$$I_3 = [\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3] = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \Rightarrow \mathbf{e}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{e}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{e}_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.$$

We are now ready to introduce the notion of the inverse matrix.

**Definition 2.4.1.** A matrix  $A \in \mathbb{F}^{n,n}$  is called *invertible* iff there exists a matrix, denoted as  $A^{-1}$ , such that  $(A^{-1})A = I_n$  and  $A(A^{-1}) = I_n$ .

Since the matrix product is non-commutative, the products  $A(A^{-1}) = I_n$  and  $(A^{-1})A = I_n$  must be specified in the definition above. Notice that we do not need to assume that the inverse matrix belongs to  $\mathbb{F}^{n,n}$ , since both products  $A(A^{-1})$  and  $(A^{-1})A$  are well-defined, we conclude that the inverse matrix must be  $n \times n$ .

**EXAMPLE 2.4.1:** Verify that the matrix and its inverse are given by

$$A = \begin{bmatrix} 2 & 2 \\ 1 & 3 \end{bmatrix}, \quad A^{-1} = \frac{1}{4} \begin{bmatrix} 3 & -2 \\ -1 & 2 \end{bmatrix}.$$

**SOLUTION:** We have to compute the products,

$$A(A^{-1}) = \begin{bmatrix} 2 & 2 \\ 1 & 3 \end{bmatrix} \frac{1}{4} \begin{bmatrix} 3 & -2 \\ -1 & 2 \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 4 & 0 \\ 0 & 4 \end{bmatrix} \Rightarrow A(A^{-1}) = I_2.$$

It is simple to check that the equation  $(A^{-1})A = I_2$  also holds. ◁

**EXAMPLE 2.4.2:** The only real numbers that are equal to its own inverses are  $a = 1$  and  $a = -1$ . This is not true in the case of matrices. Verify that the matrix  $A$  below is its own inverse, that is,

$$A = \begin{bmatrix} 1 & 1 \\ 0 & -1 \end{bmatrix} = A^{-1}.$$

**SOLUTION:** We have to compute the products,

$$A(A^{-1}) = \begin{bmatrix} 1 & 1 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 0 & -1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \Rightarrow A(A^{-1}) = I_2.$$

It is simple to check that the equation  $(A^{-1})A = I_2$  also holds.  $\triangleleft$

Not every square matrix is invertible. The following Example show a  $2 \times 2$  matrix with no inverse.

**EXAMPLE 2.4.3:** Show that matrix  $A = \begin{bmatrix} 1 & 2 \\ 3 & 6 \end{bmatrix}$  has no inverse.

**SOLUTION:** Suppose there exists the inverse matrix

$$A^{-1} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

Then, the following equation holds,

$$A(A^{-1}) = I_2 \Leftrightarrow \begin{bmatrix} 1 & 2 \\ 3 & 6 \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

The last equation implies

$$\begin{aligned} a + 2c &= 1, & b + 2d &= 0, \\ 3(a + 2c) &= 0, & 3(b + 2d) &= 1. \end{aligned}$$

However, both systems are inconsistent, so the inverse matrix  $A^{-1}$  does not exist.  $\triangleleft$

In the case of  $2 \times 2$  matrices there is a simple way to find out whether a matrix has inverse or not. If the  $2 \times 2$  matrix is invertible, then there is a simple formula for the inverse matrix. This is summarized in the following result.

**Theorem 2.4.2.** Given a  $2 \times 2$  matrix  $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ , introduce the number  $\Delta = ad - bc$ . The matrix  $A$  is invertible iff  $\Delta \neq 0$ . Furthermore, if  $A$  is invertible, its inverse is given by

$$A^{-1} = \frac{1}{\Delta} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}. \quad (2.5)$$

The number  $\Delta$  is called the *determinant of  $A$* , since it is the number that determines whether  $A$  is invertible or not, and soon we will see that it is the number that determines whether a system of linear equations has a unique solution or not. Also later on we will study generalizations to  $n \times n$  matrices of the Theorem above. That will require a generalization to  $n \times n$  matrices of the determinant  $\Delta$  of a matrix.

**EXAMPLE 2.4.4:** Compute the inverse of matrix  $A = \begin{bmatrix} 2 & 2 \\ 1 & 3 \end{bmatrix}$ , given in Example 2.4.1.

**SOLUTION:** Following Theorem 2.4.2 we first compute  $\Delta = 6 - 4 = 2$ . Since  $\Delta \neq 0$ , then  $A^{-1}$  exists and it is given by

$$A^{-1} = \frac{1}{2} \begin{bmatrix} 3 & -2 \\ -1 & 2 \end{bmatrix}.$$

$\triangleleft$

**EXAMPLE 2.4.5:** Theorem 2.4.2 says that the matrix in Example 2.4.3 is not invertible, since

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 6 \end{bmatrix} \Rightarrow \Delta = 6 - (3)(2) \Rightarrow \Delta = 0.$$

$\triangleleft$

**Proof of Theorem 2.4.2:** If the matrix  $A^{-1}$  exists, from the definition of inverse matrix it follows that  $A^{-1}$  must be  $2 \times 2$ . Suppose that the inverse of matrix  $A$  is given by

$$A^{-1} = \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \end{bmatrix}.$$

We first show that  $A^{-1}$  exists iff  $\Delta \neq 0$ .

The  $2 \times 2$  matrix  $A^{-1}$  exists iff the equation  $A(A^{-1}) = I_2$ , which is equivalent to the systems

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \Leftrightarrow \begin{cases} ax_1 + bx_2 = 1, & ay_1 + by_2 = 0, \\ cx_1 + dx_2 = 0, & cy_1 + dy_2 = 1. \end{cases} \quad (2.6)$$

Consider now the particular case  $a = 0$  and  $c = 0$ , which imply that  $\Delta = 0$ . In this case the equations above reduce to:

$$\begin{cases} bx_2 = 1, & by_2 = 0, \\ dx_2 = 0, & dy_2 = 1. \end{cases}$$

These systems have no solution, since from the first equation on the left  $b \neq 0$ , and from the first equation on the right we obtain that  $y_2 = 0$ . But this contradicts the second equation on the right, since  $dy_2$  is zero and so can never be equal to one. We then conclude there is no inverse matrix in this case.

Assume now that at least one of the coefficients  $a$  or  $c$  is non-zero, and let us return to Eqs. (2.6). Both systems can be solved using the following augmented matrix

$$\left[ \begin{array}{cc|cc} a & b & 1 & 0 \\ c & d & 0 & 1 \end{array} \right]$$

Now, perform the following Gauss operations:

$$\left[ \begin{array}{cc|cc} ac & bc & c & 0 \\ ac & ad & 0 & a \end{array} \right] \rightarrow \left[ \begin{array}{cc|cc} ac & bc & c & 0 \\ 0 & ad - bc & -c & a \end{array} \right] = \left[ \begin{array}{cc|cc} ac & bc & c & 0 \\ 0 & \Delta & -c & a \end{array} \right]$$

At least one of the source coefficients in the second row above is non-zero. Therefore, the system above is consistent iff  $\Delta \neq 0$ . This establishes the first part of the Theorem.

In order to prove the furthermore part, one can continue with the calculation above, and find the formula for the inverse matrix. This is a long calculation, since one has to study three different cases: the case  $a = 0$  and  $c \neq 0$ , the case  $a \neq 0$  and  $c = 0$ , and the case where both  $a$  and  $c$  are non-zero. It is faster to check that the expression in the Theorem is indeed the inverse of  $A$ . Since  $\Delta \neq 0$ , the matrix in Eq. (2.5) is well-defined. Then, the straightforward calculation

$$A(A^{-1}) = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \frac{1}{\Delta} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} = \frac{1}{\Delta} \begin{bmatrix} \Delta & -ab + ba \\ cd - dc & \Delta \end{bmatrix} = I_2.$$

It is not difficult to see that the second condition  $(A^{-1})A = I_2$  is also satisfied. This establishes the Theorem.  $\square$

There are many different ways to characterize  $n \times n$  invertible matrices. One possibility is to relate the existence of the inverse matrix to the solution of appropriate systems of linear equations. The following result summarizes a few of these characterizations.

**Theorem 2.4.3.** *Given a matrix  $A \in \mathbb{F}^{n,n}$ , the following statements are equivalent:*

- (a) *The matrix  $A^{-1}$  exists;*
- (b)  *$\text{rank}(A) = n$ ;*
- (c)  *$E_A = I_n$ ;*
- (d) *The homogeneous equation  $Ax = 0$  has a unique solution  $x = 0$ ;*
- (e) *The non-homogeneous equation  $Ax = b$  has a unique solution for every source  $b \in \mathbb{F}^n$ .*

**Proof of Theorem 2.4.3:** It is clear that properties (b)-(d) are all equivalent. Here we only show that (a) is equivalent to (b).

We start with (a)  $\Rightarrow$  (b): Since  $A^{-1}$  exists, the homogeneous equation  $Ax = 0$  has a unique solution  $x = 0$ . (Proof: Assume there are two solutions  $x_1$  and  $x_2$ , then  $A(x_2 - x_1) = 0$ , and so  $(x_2 - x_1) = A^{-1}A(x_2 - x_1) = A^{-1}0 = 0$ , and so  $x_2 = x_1$ .) But this implies there are no free variables in the solutions of  $Ax = 0$ , that is,  $E_A = I_n$ , that is,  $\text{rank}(A) = n$ .

We finish with (b)  $\Rightarrow$  (a): Since  $\text{rank}(A) = n$ , the non-homogeneous equation  $Ax = b$  has a unique solution  $x$  for every source vector  $b \in \mathbb{R}^n$ . In particular, there exist unique vectors  $x_1, \dots, x_n$  solutions of

$$Ax_1 = e_1, \quad \dots \quad Ax_n = e_n.$$

where  $e_i$  is the  $i$ -th column of the identity matrix  $I_n$ . This is equivalent to say that the matrix  $X = [x_1, \dots, x_n]$  satisfies the equation

$$AX = I_n.$$

This matrix  $X = A^{-1}$  since  $X$  also satisfies the equation  $XA = I_n$ . (Proof: Consider the identities

$$A - A = 0 \quad \Leftrightarrow \quad AXA - A = 0 \quad \Leftrightarrow \quad A(XA - I_n) = 0;$$

The last equation is equivalent to the  $n$  systems of equations  $Ay_i = 0$ , where  $y_i$  is the  $i$ -th column of the matrix  $(XA - I_n)$ ; Since  $\text{rank}(A) = n$ , each of these systems has a unique solution  $y_i = 0$ , that is,  $XA = I_n$ .) This establishes the Theorem.  $\square$

It is simple to see from Theorem 2.4.3 that an invertible matrix has a unique inverse.

**Corollary 2.4.4.** *An invertible matrix has a unique inverse.*

**Proof of Corollary 2.4.4:** Suppose that matrices  $X$  and  $Y$  are two inverses of an  $n \times n$  matrix  $A$ . Then,  $AX = I_n$  and  $AY = I_n$ , hence  $A(X - Y) = 0$ . The latter are  $n$  systems of linear equations, one for each column vector in  $(X - Y)$ . From Theorem 2.4.3 we know that  $\text{rank}(A) = n$ , so the only solution to these equations is the trivial solution, so each column vector vanishes, therefore  $X = Y$ .  $\square$

**2.4.2. Properties of invertible matrices.** They are summarized below.

**Theorem 2.4.5.** *If  $A$  and  $B$  are  $n \times n$  invertible matrices, then holds:*

- (a)  $(A^{-1})^{-1} = A$ ;
- (b)  $(AB)^{-1} = B^{-1}A^{-1}$ ;
- (c)  $(A^T)^{-1} = (A^{-1})^T$ .

**Proof of Theorem 2.4.5:** Since an invertible matrix is unique, we only have to verify these equations in (a)-(c).

**Part (a):** The inverse of  $A^{-1}$  is a matrix  $(A^{-1})^{-1}$  satisfying the equations

$$\left[(A^{-1})^{-1}\right](A^{-1}) = I_n, \quad (A^{-1})\left[(A^{-1})^{-1}\right] = I_n.$$

But matrix  $A$  satisfies precisely these equations, and recalling that the inverse of a matrix is unique, then  $A = (A^{-1})^{-1}$ .

**Part (b):** The proof is similar to the previous one. We verify that matrix  $(B^{-1})((A^{-1})$  satisfies the equations that  $(AB)^{-1}$  must satisfy. Then, they must be the same, since the inverse matrix is unique. Notice that,

$$\begin{aligned} (B^{-1}A^{-1})(AB) &= B^{-1}(A^{-1}A)B, & (AB)(B^{-1}A^{-1}) &= A(BB^{-1})A^{-1}, \\ &= (B^{-1})B, & &= A(A^{-1}), \\ &= I_n; & &= I_n. \end{aligned}$$

We then conclude that  $(AB)^{-1} = B^{-1}A^{-1}$ .

**Part (c):** Recall that  $(AB)^T = B^T A^T$ , therefore,

$$(A^{-1})A = I_n \Leftrightarrow [(A^{-1})A]^T = I_n^T \Leftrightarrow A^T(A^{-1})^T = I_n,$$

$$A(A^{-1}) = I_n \Leftrightarrow [A(A^{-1})]^T = I_n^T \Leftrightarrow (A^{-1})^T A^T = I_n.$$

Therefore,  $(A^{-1})^T = (A^T)^{-1}$ . This establishes the Theorem.  $\square$

The properties presented above are useful to solve equations involving matrices.

**EXAMPLE 2.4.6:** Find a matrix  $C$  solution of the equation  $(BC)^T - A = 0$ , where

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ -1 & -1 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & -1 \\ 1 & 2 \end{bmatrix}.$$

**SOLUTION:** The matrix  $B$  is invertible, since  $\Delta(B) = 3$ , and its inverse is given by

$$B^{-1} = \frac{1}{3} \begin{bmatrix} 2 & 1 \\ -1 & 1 \end{bmatrix}.$$

Therefore, matrix  $C$  is given by

$$(BC)^T = A \Leftrightarrow BC = A^T \Leftrightarrow C = B^{-1}A^T,$$

that is,

$$C = \frac{1}{3} \begin{bmatrix} 2 & 1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 3 & -2 \\ 2 & 4 & -1 \end{bmatrix} \Rightarrow C = \frac{1}{3} \begin{bmatrix} 4 & 10 & -5 \\ 1 & 1 & 1 \end{bmatrix}. \quad \triangleleft$$

**EXAMPLE 2.4.7:** The  $(n+k) \times (n+k)$  matrix  $C$  introduced in Example 2.3.10 satisfies the equation  $C^2 = I_{n+k}$ . Therefore, matrix  $C$  is its own inverse, that is,  $C^{-1} = C$ .  $\triangleleft$

**2.4.3. Computing the inverse matrix.** We now show how to use Gauss operations to find the inverse matrix in the case that such inverse exists. The main idea needed to compute the inverse of an  $n \times n$  matrix is summarized in the following Theorem. We emphasize that this is not a new result. We are just highlighting the main part of the proof in Theorem 2.4.3, (b)  $\Rightarrow$  (a).

**Theorem 2.4.6.** *Let  $A$  be an  $n \times n$  matrix. If the  $n$  systems of linear equations*

$$Ax_1 = e_1, \quad \dots \quad Ax_n = e_n, \quad (2.7)$$

*are all consistent, then the matrix  $A$  is invertible and its inverse is given by*

$$A^{-1} = [x_1, \dots, x_n].$$

*If at least one system in Eq. (2.7) is inconsistent, then matrix  $A$  is not invertible.*

**Proof of Theorem 2.4.6:** We first show that the consistency of all systems in Eq. (2.7) implies that matrix  $A$  is invertible. Indeed, if all systems in Eq. (2.7) are consistent, then the system  $Ax = b$  is also consistent for all  $b \in \mathbb{R}^n$  (Proof: The solution for the source  $b = b_1 e_1 + \dots + b_n e_n$  is simply  $x = b_1 x_1 + \dots + b_n x_n$ .) Therefore,  $\text{rank}(A) \geq n$ . Since  $A$  is an  $n \times n$  matrix, we conclude that  $\text{rank}(A) = n$ . Then, Theorem 2.4.3 implies that matrix  $A$  is invertible.

We now introduce the matrix  $X = [x_1, \dots, x_n]$  and we show that  $A^{-1} = X$ . From the definition of  $X$  we see that  $AX = I_n$ . Since  $\text{rank}(A) = n$ , the same argument given in the proof of Theorem 2.4.3 shows that  $XA = I_n$ . Since the inverse of a matrix is unique, we conclude that  $A^{-1} = X$ . This establishes the Theorem.  $\square$

This Theorem provides a method to find the inverse of a matrix.

**Corollary 2.4.7.** *The inverse of an invertible  $n \times n$  matrix  $A$  is computed with the Gauss operations such that  $[A|I_n] \longrightarrow [I_n|A^{-1}]$ .*

**Proof of Corollary 2.4.7:** One solves the  $n$  linear systems in Eq. (2.7). Since all these systems share the same coefficient matrix, matrix  $A$ , one can solve them all at the same time, introducing the augmented matrix

$$[A|e_1, \dots, e_n] = [A|I_n].$$

In the proof of Theorem 2.4.6 we show that  $\text{rank}(A) = n$ . Therefore, its reduced echelon form is  $E_A = I_n$ , so the Gauss method implies that

$$[A|I_n] = [A|e_1, \dots, e_n] \rightarrow [E_A|x_1, \dots, x_n] = [I_n|A^{-1}] \Leftrightarrow [A|I_n] \longrightarrow [I_n|A^{-1}].$$

This establishes the Corollary.  $\square$

**EXAMPLE 2.4.8:** Find the inverse of the matrix  $A = \begin{bmatrix} 2 & 4 \\ 1 & 3 \end{bmatrix}$ .

**SOLUTION:** We have to solve the two systems of equations  $Ax_1 = e_1$  and  $Ax_2 = e_2$ , that is,

$$\begin{bmatrix} 2 & 4 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 2 & 4 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}.$$

We solve both systems at the same time using the Gauss method on the augmented matrix

$$\left[ \begin{array}{cc|cc} 2 & 4 & 1 & 0 \\ 1 & 3 & 0 & 1 \end{array} \right] \rightarrow \left[ \begin{array}{cc|cc} 1 & 2 & 1/2 & 0 \\ 1 & 3 & 0 & 1 \end{array} \right] \rightarrow \left[ \begin{array}{cc|cc} 1 & 2 & 1/2 & 0 \\ 0 & 1 & -1/2 & 1 \end{array} \right] \rightarrow \left[ \begin{array}{cc|cc} 1 & 0 & 3/2 & -2 \\ 0 & 1 & -1/2 & 1 \end{array} \right]$$

therefore,  $A^{-1} = \frac{1}{2} \begin{bmatrix} 3 & -4 \\ -1 & 2 \end{bmatrix}$ .  $\triangleleft$

**EXAMPLE 2.4.9:** Find the inverse of the matrix  $A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 5 & 7 \\ 3 & 7 & 9 \end{bmatrix}$ .

**SOLUTION:** We compute the inverse matrix as follows:

$$\begin{aligned} & \left[ \begin{array}{ccc|ccc} 1 & 2 & 3 & 1 & 0 & 0 \\ 2 & 5 & 7 & 0 & 1 & 0 \\ 3 & 7 & 9 & 0 & 0 & 1 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|ccc} 1 & 2 & 3 & 1 & 0 & 0 \\ 0 & 1 & 1 & -2 & 1 & 0 \\ 0 & 1 & 0 & -3 & 0 & 1 \end{array} \right] \rightarrow \\ & \left[ \begin{array}{ccc|ccc} 1 & 0 & 1 & 5 & -2 & 0 \\ 0 & 1 & 1 & -2 & 1 & 0 \\ 0 & 0 & -1 & -1 & -1 & 1 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|ccc} 1 & 0 & 0 & 4 & -3 & 1 \\ 0 & 1 & 0 & -3 & 0 & 1 \\ 0 & 0 & 1 & 1 & 1 & -1 \end{array} \right] \Rightarrow A^{-1} = \begin{bmatrix} 4 & -3 & 1 \\ -3 & 0 & 1 \\ 1 & 1 & -1 \end{bmatrix}. \end{aligned}$$

**EXAMPLE 2.4.10:** Is matrix  $A = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}$  invertible?

**SOLUTION:** No, since  $\left[ \begin{array}{cc|cc} 1 & 2 & 1 & 0 \\ 2 & 4 & 0 & 1 \end{array} \right] \rightarrow \left[ \begin{array}{cc|cc} 1 & 2 & 1 & 0 \\ 0 & 0 & -2 & 1 \end{array} \right]$ , are inconsistent.  $\triangleleft$

**Further reading.** Almost any book in linear algebra introduces the inverse matrix. See Sections 3.7 and 3.9 in Meyer's book [3].

## 2.4.4. Exercises.

2.4.1.- When possible, find the inverse of the following matrices. (Check your answers using matrix multiplication.)

(a)

$$A = \begin{bmatrix} 1 & 2 \\ 1 & 3 \end{bmatrix};$$

(b)

$$A = \begin{bmatrix} -1 & -2 & 1 \\ 3 & 2 & -6 \\ 1 & 1 & -2 \end{bmatrix};$$

(c)

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}.$$

2.4.2.- Find the values of the constant  $k$  such that the matrix  $A$  below is not invertible,

$$A = \begin{bmatrix} 1 & 1 & -1 \\ 2 & 3 & k \\ 1 & k & 3 \end{bmatrix}.$$

2.4.3.- Consider the matrix

$$A = \begin{bmatrix} 0 & -1 & 0 \\ 2 & 0 & 1 \\ 0 & 1 & -1 \end{bmatrix},$$

(a) Find the inverse of matrix  $A$ .

(b) Use the previous part to find a solution to the linear system

$$Ax = b, \quad b = \begin{bmatrix} 1 \\ -1 \\ 3 \end{bmatrix}.$$

2.4.4.- Show that for every invertible matrix  $A$  holds that  $[A, A^{-1}] = 0$ .

2.4.5.- Find a matrix  $X$  such that the equation  $X = AX + B$  holds for

$$A = \begin{bmatrix} 0 & -1 & 0 \\ 0 & 0 & -1 \\ 0 & 0 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 3 & 3 \end{bmatrix}.$$

2.4.6.- If  $A$  is invertible and symmetric, then show that  $A^{-1}$  is also symmetric.

2.4.7.- Prove that: If the square matrix  $A$  satisfies  $A^2 = 0$ , then the matrix  $(I - A)$  is invertible.

2.4.8.- Prove that: If the square matrix  $A$  satisfies  $A^3 = 0$ , then the matrix  $(I - A)$  is invertible.

2.4.9.- Let  $A$  be a square matrix. Prove the following statements:

(a) If  $A$  contains a zero column then  $A$  is not invertible;(b) If one column is multiple of another column in  $A$ , then matrix  $A$  is not invertible.(c) Use the trace function to prove the following statement: If  $A$  is an  $m \times n$  matrix and  $B$  is an  $n \times m$  matrix such that  $AB = I_m$  and  $BA = I_n$ , then  $m = n$ .

2.4.10.- Consider the invertible matrices  $A \in \mathbb{F}^{r,r}$ ,  $B \in \mathbb{F}^{s,s}$  and the matrix  $C \in \mathbb{F}^{r,s}$ . Prove that the inverse of

$$\begin{bmatrix} A & \vdots & C \\ \dots & \dots & \dots \\ 0 & \vdots & B \end{bmatrix}$$

is given by

$$\begin{bmatrix} A^{-1} & \vdots & -A^{-1}CB^{-1} \\ \dots & \dots & \dots \\ 0 & \vdots & B^{-1} \end{bmatrix}.$$



## 2.5. NULL AND RANGE SPACES

**2.5.1. Definition of the spaces.** A matrix  $A \in \mathbb{F}^{m,n}$  defines two functions,  $A : \mathbb{F}^n \rightarrow \mathbb{F}^m$  and  $A^T : \mathbb{F}^m \rightarrow \mathbb{F}^n$ , which in turn, determine two sets in  $\mathbb{F}^n$  and two sets in  $\mathbb{F}^m$ . In this Section we define these sets, we call them null and range spaces, and study the main relations among them. We start introducing the two spaces associated with the function  $A$ .

**Definition 2.5.1.** Consider the matrix  $A \in \mathbb{F}^{m,n}$  defining the linear function  $A : \mathbb{F}^n \rightarrow \mathbb{F}^m$ . The **null space** of the function  $A$  is the set  $N(A) \subset \mathbb{F}^n$  given by

$$N(A) = \{x \in \mathbb{F}^n : Ax = 0\}.$$

The **range space** of the function  $A$  is set  $R(A) \subset \mathbb{F}^m$  given by

$$R(A) = \{y \in \mathbb{F}^m : y = Ax, \text{ for all } x \in \mathbb{F}^n\}.$$

In Fig. 25 we show a picture, usual in set theory, sketching the null and range spaces associated with a matrix  $A$ . One can see in these pictures that for a function  $A : \mathbb{F}^n \rightarrow \mathbb{F}^m$ , the null space is a subset in  $\mathbb{F}^n$ , while the range space is a subset in  $\mathbb{F}^m$ . Recall that the homogeneous equation  $Ax = 0$  always has the trivial solution  $x = 0$ . This property implies both that  $0 \in \mathbb{F}^n$  also belongs to  $N(A)$  and that  $0 \in \mathbb{F}^m$  also belongs to  $R(A)$ .

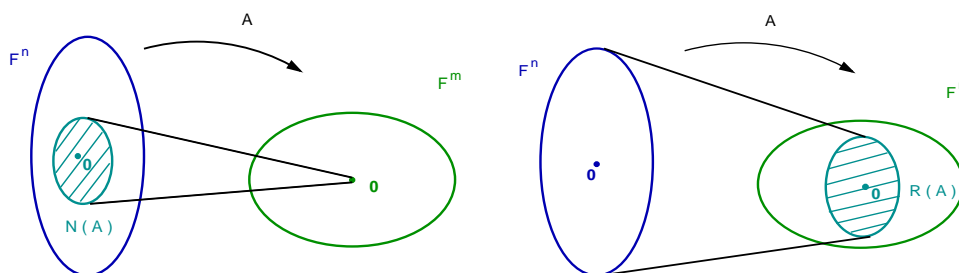


FIGURE 25. Sketch of the null space,  $N(A)$ , and the range space,  $R(A)$ , of a linear function determined by the  $m \times n$  matrix  $A$ .

The **four spaces** associated with a matrix  $A$  are  $N(A)$  and  $R(A)$  together with  $N(A^T)$  and  $R(A^T)$ . Notice that these null and range spaces are subsets of different spaces. More precisely, a matrix  $A \in \mathbb{F}^{m,n}$  defines the functions and subsets,

$$A : \mathbb{F}^n \rightarrow \mathbb{F}^m \quad \text{and} \quad A^T : \mathbb{F}^m \rightarrow \mathbb{F}^n,$$

$$N(A), R(A^T) \subset \mathbb{F}^n, \quad \text{while} \quad R(A), N(A^T) \subset \mathbb{F}^m.$$

**EXAMPLE 2.5.1:** Find the  $N(A)$  and  $R(A)$  for the function  $A : \mathbb{R}^3 \rightarrow \mathbb{R}^2$  given by

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 1 \end{bmatrix}.$$

**SOLUTION:** We first find  $N(A)$ , the null space of  $A$ . This is the set of elements  $x \in \mathbb{R}^3$  that are solutions of the equation  $Ax = 0$ . We use the Gauss method to find all such solutions,

$$\begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 0 & 0 & -5 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 0 \\ 0 & 0 & 1 \end{bmatrix} \Rightarrow \begin{cases} x_1 = -2x_2 \\ x_3 = 0, \\ x_2 : \text{free variable.} \end{cases}$$

Therefore, the elements in  $N(\mathbf{A})$  are given by

$$N(\mathbf{A}) \ni \mathbf{x} = \begin{bmatrix} -2x_2 \\ x_2 \\ 0 \end{bmatrix} = \begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix} x_2 \Rightarrow N(\mathbf{A}) = \text{Span}\left(\left\{ \begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix} \right\}\right) \subset \mathbb{R}^3.$$

We now find  $R(\mathbf{A})$ , the range space of  $\mathbf{A}$ . This is the set of elements  $\mathbf{y} \in \mathbb{R}^2$  that can be expressed as  $\mathbf{y} = \mathbf{A}\mathbf{x}$  for any  $\mathbf{x} \in \mathbb{R}^3$ . In our case this means

$$\mathbf{y} = \mathbf{A}\mathbf{x} = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} x_1 + \begin{bmatrix} 2 \\ 4 \end{bmatrix} x_2 + \begin{bmatrix} 3 \\ 1 \end{bmatrix} x_3.$$

What this last equation says, is that the elements of  $R(\mathbf{A})$  can be expressed as a linear combination of the column vectors of matrix  $\mathbf{A}$ . Therefore, the set  $R(\mathbf{A})$  is indeed the set of all possible linear combinations of the column vectors of matrix  $\mathbf{A}$ . We then conclude that

$$R(\mathbf{A}) = \text{Span}\left(\left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 2 \\ 4 \end{bmatrix}, \begin{bmatrix} 3 \\ 1 \end{bmatrix} \right\}\right).$$

Notice that

$$2 \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 2 \\ 4 \end{bmatrix} \Rightarrow \text{Span}\left(\left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 2 \\ 4 \end{bmatrix}, \begin{bmatrix} 3 \\ 1 \end{bmatrix} \right\}\right) = \text{Span}\left(\left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 3 \\ 1 \end{bmatrix} \right\}\right).$$

Therefore, we conclude that the smallest set whose span is  $R(\mathbf{A})$  contains two elements, the first and third columns of  $\mathbf{A}$ , that is,

$$R(\mathbf{A}) = \text{Span}\left(\left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 3 \\ 1 \end{bmatrix} \right\}\right) = \mathbb{R}^2.$$

◁

In the Example 2.5.1 above we have seen that both sets  $N(\mathbf{A})$  and  $R(\mathbf{A})$  can be expressed as spans of appropriate vectors. It is not surprising to see that the same property holds for  $N(\mathbf{A}^T)$  and  $R(\mathbf{A}^T)$ , as we observe in the following Example.

**EXAMPLE 2.5.2:** Find the  $N(\mathbf{A}^T)$  and  $R(\mathbf{A}^T)$  for the function  $\mathbf{A}^T : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ , where  $\mathbf{A}$  is the matrix in Example 2.5.1, that is,

$$\mathbf{A}^T = \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ 3 & 1 \end{bmatrix}.$$

**SOLUTION:** We start finding the  $N(\mathbf{A}^T)$ , that is, all vectors  $\mathbf{y} \in \mathbb{R}^2$  solutions of the homogeneous equation  $\mathbf{A}^T\mathbf{y} = \mathbf{0}$ . Using the Gauss method we obtain,

$$\begin{bmatrix} 1 & 2 \\ 2 & 4 \\ 3 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 \\ 0 & 0 \\ 0 & -5 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \Rightarrow \mathbf{y} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow N(\mathbf{A}^T) = \left\{ \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right\} \subset \mathbb{R}^2.$$

We now find the  $R(\mathbf{A}^T)$ . This is the set of  $\mathbf{x} \in \mathbb{R}^3$  that can be expressed as  $\mathbf{x} = \mathbf{A}^T\mathbf{y}$  for any  $\mathbf{y} \in \mathbb{R}^2$ , that is,

$$\mathbf{x} = \mathbf{A}^T\mathbf{y} = \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} y_1 + \begin{bmatrix} 2 \\ 4 \\ 1 \end{bmatrix} y_2.$$

As in the previous example, this last equation says that the elements of  $R(\mathbf{A}^T)$  can be expressed as a linear combination of the column vectors of matrix  $\mathbf{A}^T$ . Therefore, the set

$R(A^T)$  is indeed the set of all possible linear combinations of the column vectors of matrix  $A^T$ . We then conclude that

$$R(A^T) = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 2 \\ 4 \\ 1 \end{bmatrix}\right\}\right) \subset \mathbb{R}^3.$$

In Fig. 26 we have sketched the sets  $N(A)$  and  $R(A^T)$ , which are subsets of  $\mathbb{R}^3$ .  $\triangleleft$

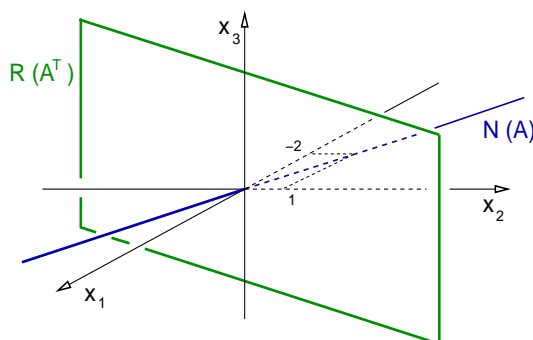


FIGURE 26. The horizontal blue line is the  $N(A)$ , and the vertical green plane is the  $R(A^T)$ , for the matrix  $A$  is given in Examples 2.5.1 and 2.5.2. Notice that the blue line is perpendicular to the green plane. We will see that this is always true for all  $m \times n$  matrices.

**EXAMPLE 2.5.3:** Find the  $N(A)$  and  $R(A)$  for the matrix

$$A = \begin{bmatrix} 1 & 3 & -1 \\ 2 & 6 & -2 \\ 3 & 9 & -3 \end{bmatrix}.$$

**SOLUTION:** Any element  $x \in N(A)$  must be solution of the homogeneous equation  $Ax = 0$ . The Gauss method implies,

$$\begin{bmatrix} 1 & 3 & -1 \\ 2 & 6 & -2 \\ 3 & 9 & -3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 3 & -1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \Rightarrow \begin{cases} x_1 = -3x_2 + x_3, \\ x_2, x_3 \text{ free variables.} \end{cases}$$

Therefore, every element  $x \in N(A)$  has the form

$$x = \begin{bmatrix} -3x_2 + x_3 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} -3 \\ 1 \\ 0 \end{bmatrix} x_2 + \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} x_3 \Rightarrow N(A) = \text{Span}\left(\left\{\begin{bmatrix} -3 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}\right\}\right).$$

The  $R(A)$  is the set of all  $y \in \mathbb{R}^3$  such that  $y = Ax$  for some  $x \in \mathbb{R}^3$ . Therefore,

$$y = \begin{bmatrix} 1 & 3 & -1 \\ 2 & 6 & -2 \\ 3 & 9 & -3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = x_1 \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + x_2 \begin{bmatrix} 3 \\ 6 \\ 9 \end{bmatrix} + x_3 \begin{bmatrix} -1 \\ -2 \\ -3 \end{bmatrix},$$

and we conclude that

$$R(A) = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 3 \\ 6 \\ 9 \end{bmatrix}, \begin{bmatrix} -1 \\ -2 \\ -3 \end{bmatrix}\right\}\right).$$

However, the expression above can be simplified noticing that

$$\begin{bmatrix} 3 \\ 6 \\ 9 \end{bmatrix} = 3 \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \quad \begin{bmatrix} -1 \\ -2 \\ -3 \end{bmatrix} = - \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}.$$

So the smallest sets whose span is  $R(\mathbf{A})$  contain only one vector, and a possible choice is the following:

$$R(\mathbf{A}) = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}\right\}\right).$$

◁

**2.5.2. Main properties.** The property of the null and range sets we have found for the matrix in Examples 2.5.1 and 2.5.2 also holds for every matrix. This will be an important property later on, so we give it a name.

**Definition 2.5.2.** A subset of  $U \subset \mathbb{F}^n$  is called **closed under linear combination** iff for all elements  $\mathbf{x}, \mathbf{y} \in U$  and all scalars  $a, b \in \mathbb{F}$  holds that  $(a\mathbf{x} + b\mathbf{y}) \in U$ .

In words, a set  $U \subset \mathbb{F}^n$  is closed under linear combination iff every linear combination of elements in  $U$  stays in  $U$ . In Chapter 4 we generalize the linear combination structure of  $\mathbb{F}^n$  into a structure we call a vector space; we will see in that Chapter that sets in a vector space which are closed under linear combinations are smaller vector spaces inside the original vector space, and will be called subspaces. We now state as a general result the property we found both in Examples 2.5.1 and 2.5.2.

**Theorem 2.5.3.** The sets  $N(\mathbf{A}) \subset \mathbb{F}^n$  and  $R(\mathbf{A}) \subset \mathbb{F}^m$  of a matrix  $\mathbf{A} \in \mathbb{F}^{m,n}$  are closed under linear combinations in  $\mathbb{F}^n$  and  $\mathbb{F}^m$ , respectively. Furthermore, denoting the matrix  $\mathbf{A} = [\mathbf{A}_{:1}, \dots, \mathbf{A}_{:n}]$ , we conclude that  $R(\mathbf{A}) = \text{Span}(\{\mathbf{A}_{:1}, \dots, \mathbf{A}_{:n}\})$ .

It is common in the literature to introduce the **column space** of an  $m \times n$  matrix  $\mathbf{A} = [\mathbf{A}_{:1}, \dots, \mathbf{A}_{:n}]$ , denoted as  $\text{Col}(\mathbf{A})$ , as the set of all linear combinations of the column vectors of  $\mathbf{A}$ , that is,  $\text{Col}(\mathbf{A}) = \text{Span}(\{\mathbf{A}_{:1}, \dots, \mathbf{A}_{:n}\})$ . The Proposition above then says that

$$R(\mathbf{A}) = \text{Col}(\mathbf{A}).$$

**Proof of Theorem 2.5.3:** The sets  $N(\mathbf{A})$  and  $R(\mathbf{A})$  are closed under linear combinations because the matrix-vector product is a linear operation. Consider two arbitrary elements  $\mathbf{x}_1, \mathbf{x}_2 \in N(\mathbf{A})$ , that is,  $\mathbf{A}\mathbf{x}_1 = \mathbf{0}$  and  $\mathbf{A}\mathbf{x}_2 = \mathbf{0}$ . Then, for any  $a, b \in \mathbb{F}$  holds

$$\mathbf{A}(a\mathbf{x}_1 + b\mathbf{x}_2) = a\mathbf{A}\mathbf{x}_1 + b\mathbf{A}\mathbf{x}_2 = \mathbf{0} \quad \Rightarrow \quad (a\mathbf{x}_1 + b\mathbf{x}_2) \in N(\mathbf{A}).$$

Therefore,  $N(\mathbf{A}) \subset \mathbb{F}^n$  is closed under linear combinations. Analogously, consider two arbitrary elements  $\mathbf{y}_1, \mathbf{y}_2 \in R(\mathbf{A})$ , that is, there exist  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{F}^n$  such that  $\mathbf{y}_1 = \mathbf{A}\mathbf{x}_1$  and  $\mathbf{y}_2 = \mathbf{A}\mathbf{x}_2$ . Then, for any  $a, b \in \mathbb{F}$  holds

$$(a\mathbf{y}_1 + b\mathbf{y}_2) = a\mathbf{A}\mathbf{x}_1 + b\mathbf{A}\mathbf{x}_2 = \mathbf{A}(a\mathbf{x}_1 + b\mathbf{x}_2) \quad \Rightarrow \quad (a\mathbf{y}_1 + b\mathbf{y}_2) \in R(\mathbf{A}).$$

Therefore,  $R(\mathbf{A}) \subset \mathbb{F}^m$  is closed under linear combinations. The furthermore part is proved as follows. Denote  $\mathbf{A} = [\mathbf{A}_{:1}, \dots, \mathbf{A}_{:n}]$ , then any element  $\mathbf{y} \in R(\mathbf{A})$  can be expressed as  $\mathbf{y} = \mathbf{A}\mathbf{x}$  for some  $\mathbf{x} = [x_i] \in \mathbb{F}^n$ , that is,

$$\mathbf{y} = \mathbf{A}\mathbf{x} = [\mathbf{A}_{:1}, \dots, \mathbf{A}_{:n}] \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \mathbf{A}_{:1}x_1 + \dots + \mathbf{A}_{:n}x_n \in \text{Span}(\{\mathbf{A}_{:1}, \dots, \mathbf{A}_{:n}\}).$$

This implies that  $R(A) \subset \text{Col}(A)$ . The opposite inclusion,  $\text{Col}(A) \subset R(A)$  is trivial, since any element of the form  $y = A_{:1}x_1 + \cdots + A_{:n}x_n \in \text{Col}(A)$  also belongs to  $R(A)$ , since  $y = Ax$ , where  $x = [x_i]$ . This establishes the Theorem.  $\square$

The null and range spaces of a square matrix characterize whether the matrix is invertible or not. The following result is a simple rewriting of Theorem 2.4.3.

**Theorem 2.5.4.** *Given a matrix  $A \in \mathbb{F}^{n,n}$ , the following statements are equivalent:*

- (a) *The matrix  $A^{-1}$  exists;*
- (b)  *$N(A) = \{0\}$ ;*
- (c)  *$R(A) = \mathbb{F}^n$ .*

**Proof of Theorem 2.5.4:** This result follows from Theorem 2.4.3. More precisely, part (b) follows from Theorem 2.4.3 part (c). While part (c) follows from Theorem 2.4.3 part (e).  $\square$

**2.5.3. Gauss operations.** In this section we study the relations between the four spaces associated with the matrices  $A$  and  $B$  in the case that matrix  $B$  can be obtained from matrix  $A$  by performing Gauss operations. We use the following notation:  $A \xrightarrow{\text{row}} B$  to indicate that  $A$  can be transformed into matrix  $B$  by performing Gauss operations on the rows of  $A$ .

**Theorem 2.5.5.** *If matrices  $A, B \in \mathbb{F}^{m,n}$ , then*

- (a)  $A \xrightarrow{\text{row}} B \Leftrightarrow N(A) = N(B)$ ;
- (b)  $A \xrightarrow{\text{row}} B \Leftrightarrow R(A^T) = R(B^T)$ .

It is simple to see that this result is true. Gauss operations do not change the solutions of linear systems, so the linear system  $Ax = 0$  has exactly the same solutions  $x$  as the linear system  $Bx = 0$ , that is,  $N(A) = N(B)$ . The second property must be also true, since Gauss operations on the rows of  $A$  are equivalent to Gauss operations on the columns of  $A^T$ . Now, it is simple to see that each of the Gauss operations on the columns of  $A^T$  do not change the  $\text{Col}(A^T)$ , hence  $R(A^T) = R(B^T)$ .

One could say that the paragraph above is enough for a proof of the Theorem. However, we would like to have a more detailed presentation of the ideas above. One way is to use matrix multiplication to express the Gauss operation property that they do not change the solutions to linear systems. One can prove the following: If the  $m \times n$  matrices  $A$  and  $B$  are related by Gauss operations on their rows, then there exists an  $m \times m$  invertible matrix  $G$  such that  $GA = B$ . The proof of this property is simple. Each one of the Gauss operations is associated with an invertible matrix,  $E$ , called an elementary Gauss matrix. Every elementary Gauss matrix is invertible, since every Gauss operation can always be reversed. The result of several Gauss operations on a matrix  $A$  is the product of the appropriate elementary Gauss matrices in the same order as the Gauss operations are performed. If the matrix  $B$  is obtained from matrix  $A$  by doing Gauss operations given by matrices  $E_i$ , for  $i = 1, \dots, k$ , in that order, we can express the result of the Gauss method as follows:

$$E_k \cdots E_1 A = B, \quad G = E_k \cdots E_1 \Rightarrow GA = B.$$

Since each elementary Gauss matrix is invertible, then the matrix  $G$  is also invertible. The following example shows all  $3 \times 3$  elementary Gauss matrices.

**EXAMPLE 2.5.4:** Find all the elementary Gauss matrices which operate on  $3 \times n$  matrices.

**SOLUTION:** In the case of  $3 \times n$  matrices, the elementary Gauss matrices are  $3 \times 3$ . We present these matrices for each one of the three main types of Gauss operations. Consider the matrices  $E_i$  for  $i = 1, 2, 3$  are given by

$$E_1 = \begin{bmatrix} k & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad E_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & k & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad E_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & k \end{bmatrix};$$

then, the product  $E_i A$  represents the Gauss operation of multiplying by  $k$  the first, second and third row of  $A$ , respectively. Consider the matrices  $E_i$  for  $i = 4, 5, 6$  given by

$$E_4 = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad E_5 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}, \quad E_6 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix};$$

then, the product  $E_i A$  for  $i = 4, 5, 6$  represents the Gauss operation of interchanging the first and second, the first and third, and the second and third rows of  $A$ , respectively. Finally, consider the matrices  $E_i$  for  $j = 7, \dots, 12$  given by

$$E_7 = \begin{bmatrix} 1 & 0 & 0 \\ k & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad E_8 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ k & 0 & 1 \end{bmatrix}, \quad E_9 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & k & 1 \end{bmatrix},$$

$$E_{10} = \begin{bmatrix} 1 & k & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad E_{11} = \begin{bmatrix} 1 & 0 & k \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad E_{12} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & k \\ 0 & 0 & 1 \end{bmatrix};$$

then, the product  $E_i A$  for  $i = 7, \dots, 12$  represents the Gauss operation of multiplying by  $k$  one row of  $A$  and add the result to another row of  $A$ .  $\triangleleft$

**Proof of Theorem 2.5.5:** Recall the comment below Theorem 2.5.5: If the  $m \times n$  matrices  $A$  and  $B$  are related by Gauss operations on their rows, then there exists an  $m \times m$  invertible matrix  $G$  such that  $GA = B$ . This observation is the key to show that  $N(A) = N(B)$ , since given any element  $x \in N(A)$

$$Ax = 0 \quad \Leftrightarrow \quad GAx = 0,$$

where the equivalence follows from  $G$  being invertible. Then it is simple to see that

$$0 = GAx = Bx \quad \Leftrightarrow \quad x \in N(B).$$

Therefore, we have shown that  $N(A) = N(B)$ .

We now show the converse statement. If  $N(A) = N(B)$  this means that their reduced echelon forms are the same, that is,  $E_A = E_B$ . This means that there exist Gauss operations on the rows of  $A$  that transform it into matrix  $B$ .

We now show that  $R(A^T) = R(B^T)$ . Given any element  $x \in R(A^T)$  we know that there exists an element  $y \in \mathbb{F}^m$  such that

$$x = A^T y = A^T G^T (G^T)^{-1} y = (GA)^T \tilde{y} = B^T \tilde{y}, \quad \tilde{y} = (G^T)^{-1} y.$$

We have shown that given any  $x \in R(A^T)$ , then  $x \in R(B^T)$ , that is,  $R(A^T) \subset R(B^T)$ . The opposite implication is proven in the same way: Given any  $x \in R(B^T)$  there exists  $\tilde{y} \in \mathbb{F}^m$  such that

$$x = B^T \tilde{y} = B^T (G^T)^{-1} G^T \tilde{y} = (G^{-1} B)^T y = A^T y, \quad y = G^T \tilde{y}.$$

We have shown that given any  $x \in R(B^T)$ , then  $x \in R(A^T)$ , that is,  $R(B^T) \subset R(A^T)$ . Therefore,  $R(A^T) = R(B^T)$ .

We now show that the converse statement. Assume that  $R(A^T) = R(B^T)$ . This means that every row in matrix  $A$  is a linear combination of the rows in matrix  $B$ . This also means that there exists Gauss operations on the rows of  $A$  such that transform  $A$  into  $B$ . This establishes the Theorem.  $\square$

**EXAMPLE 2.5.5:** Verify Theorem 2.5.5 for matrix  $A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 1 \end{bmatrix}$  and  $E_A = \begin{bmatrix} 1 & 2 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ .

**SOLUTION:** Since  $E_A$  is the reduced echelon form of matrix  $A$ , then  $R(A^T) = R((E_A)^T)$ . Consider the matrix  $A$  and its reduced echelon form  $E_A$  given by

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 1 \end{bmatrix}, \quad E_A = \begin{bmatrix} 1 & 2 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

Their respective transposed matrices are given by

$$A^T = \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ 3 & 1 \end{bmatrix}, \quad (E_A)^T = \begin{bmatrix} 1 & 0 \\ 2 & 0 \\ 0 & 1 \end{bmatrix}.$$

The result of Theorem 2.5.5 is that

$$R(A^T) = R((E_A)^T) \Leftrightarrow \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 2 \\ 4 \\ 1 \end{bmatrix}\right\}\right) = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}\right\}\right).$$

We can verify that this result above is correct, since the column vectors in  $A^T$  are linear combinations of the column vectors in  $(E_A)^T$ , as the following equations show,

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix} + 3 \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \quad \begin{bmatrix} 2 \\ 4 \\ 1 \end{bmatrix} = 2 \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.$$

◁

The Example 2.5.5 above helps understand one important meaning of Gauss operations. *Gauss operations on the matrix  $A$  rows are linear combination of the matrix  $A^T$  columns.* This is the reason why the interchange change  $A \leftrightarrow B$  by doing Gauss operations on rows leaves the the range spaces of their transposes invariant, that is,  $R(A^T) = R(B^T)$ .

**EXAMPLE 2.5.6:** Given the matrices  $A = \begin{bmatrix} 1 & 1 & 5 \\ 2 & 0 & 6 \\ 1 & 2 & 7 \end{bmatrix}$ ,  $B = \begin{bmatrix} 1 & -4 & 4 \\ 4 & -8 & 6 \\ 0 & -4 & 5 \end{bmatrix}$ , verify whether

$$R(A^T) = R(B^T)? \quad R(A) = R(B)? \quad N(A) = N(B)? \quad N(A^T) = N(B^T)?$$

**SOLUTION:** We base our answer in Theorem 2.5.5 and an extra observation. First, let  $E_A$  and  $E_B$  be the reduced echelon forms of  $A$  and  $B$ , respectively, and let  $E_{A^T}$  and  $E_{B^T}$  be the reduced echelon forms of  $A^T$  and  $B^T$  respectively. The extra observation is the following:  $E_A = E_B$  iff  $A \xrightarrow{\text{row}} B$ . This observation and Theorem 2.5.5 imply that  $E_A = E_B$  is equivalent to  $R(A^T) = R(B^T)$  and it is also equivalent to  $N(A) = N(B)$ . We then find  $E_A$  and  $E_B$ ,

$$A = \begin{bmatrix} 1 & 1 & 5 \\ 2 & 0 & 6 \\ 1 & 2 & 7 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 1 & 5 \\ 0 & -2 & -4 \\ 0 & 1 & 2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 0 \end{bmatrix} = E_A,$$

$$B = \begin{bmatrix} 1 & -4 & 4 \\ 4 & -8 & 6 \\ 0 & -4 & 5 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -4 & 4 \\ 0 & 8 & -10 \\ 0 & -4 & 5 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -4 & 4 \\ 0 & 4 & -5 \\ 0 & 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & -5/4 \\ 0 & 0 & 0 \end{bmatrix} = E_B.$$

Since  $E_A \neq E_B$ , we conclude that  $R(A^T) \neq R(B^T)$ . This result also says that  $N(A) \neq N(B)$ . So for the first and third questions, the answer is no.

A similar argument also says that  $E_{A^T} = E_{B^T}$  iff  $A^T \xrightarrow{\text{row}} B^T$ . This observation and Theorem 2.5.5 imply that  $E_{A^T} = E_{B^T}$  is equivalent to  $R(A) = R(B)$  and it is also equivalent

to  $N(A^T) = N(B^T)$ . We then find  $E_{A^T}$  and  $E_{B^T}$ ,

$$A^T = \begin{bmatrix} 1 & 2 & 1 \\ 1 & 0 & 2 \\ 5 & 6 & 7 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 1 \\ 0 & -2 & 1 \\ 0 & -4 & 2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 2 \\ 0 & 1 & -1/2 \\ 0 & 0 & 0 \end{bmatrix} = E_{A^T},$$

$$B^T = \begin{bmatrix} 1 & 4 & 0 \\ -4 & -8 & -4 \\ 4 & 6 & 5 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 4 & 0 \\ 0 & 8 & -4 \\ 0 & -10 & 5 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 4 & 0 \\ 0 & 2 & -1 \\ 0 & 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 2 \\ 0 & 1 & -1/2 \\ 0 & 0 & 0 \end{bmatrix} = E_{B^T}.$$

Since  $E_{A^T} = E_{B^T}$ , we conclude that  $R(A) = R(B)$ . This result on the reduced echelon forms also says that  $N(A^T) = N(B^T)$ . So for the second and four questions, the answer is yes.  $\triangleleft$

We finish this Section with two results concerning the ranks of a matrix and its transpose. The first result says that transposition operation on a matrix does not change its rank.

**Theorem 2.5.6.** *Every matrix  $A \in \mathbb{F}^{m,n}$  satisfies that  $\text{rank}(A) = \text{rank}(A^T)$ .*

We delay the proof of the first result to Chapter 6, where we introduce the notion of inner product in a vector space. Using an inner product in  $\mathbb{F}^n$ , the proof of Theorem 2.5.6 below is simple. We now introduce a particular name for those matrices having the maximum possible rank.

**Definition 2.5.7.** *An  $m \times n$  matrix  $A$  has **full rank** iff  $\text{rank}(A) = \min(m, n)$ .*

Our last result of this Section concerns full rank matrices and relates the rank of a matrix  $A$  to the size of both  $N(A)$  and  $N(A^T)$ .

**Theorem 2.5.8.** *If matrix  $A \in \mathbb{F}^{m,n}$  has full rank, then hold:*

- (a) *If  $m = n$ , then  $\text{rank}(A) = \text{rank}(A^T) = n = m \Leftrightarrow \{0\} = N(A) = N(A^T) \subset \mathbb{F}^n$ ;*
- (b) *If  $m < n$ , then  $\text{rank}(A) = \text{rank}(A^T) = m < n \Leftrightarrow \begin{cases} \{0\} \subsetneq N(A) \subset \mathbb{F}^n, \\ \{0\} = N(A^T) \subset \mathbb{F}^m; \end{cases}$*
- (c) *If  $m > n$ , then  $\text{rank}(A) = \text{rank}(A^T) = n < m \Leftrightarrow \begin{cases} \{0\} = N(A) \subset \mathbb{F}^n, \\ \{0\} \subsetneq N(A^T) \subset \mathbb{F}^m. \end{cases}$*

**Proof of Theorem 2.5.8:** We start recalling that the rank of a matrix  $A$  is the number of pivot columns in its reduced echelon form  $E_A$ .

If an  $m \times n$  matrix  $A$  has  $\text{rank}(A) = n$ , this means two things: First,  $n \leq m$ ; and second, that every column in  $E_A$  has a pivot. The latter property implies that there is no free variables in the solution of the equation  $Ax = 0$ , and so  $x = 0$  is the unique solution. We conclude that  $N(A) = \{0\}$ . In order to study  $N(A^T)$  we need to consider the two possible cases  $n = m$  or  $n < m$ . If  $n = m$ , then the matrices  $A$  and  $A^T$  are square, and the same argument about free variables applies to solutions of  $A^T y = 0$ , so we conclude that  $N(A^T) = \{0\}$ . This proves (a). In the case that  $n < m$ , then there are free variables in the solution of the equation  $A^T y = 0$ , therefore  $\{0\} \subsetneq N(A^T)$ . This proves (c).

If an  $m \times n$  matrix  $A$  has  $\text{rank}(A) = m$ , recalling that  $\text{rank}(A) = \text{rank}(A^T)$ , this means two things: First,  $m \leq n$ ; and second, that every column in  $E_{A^T}$  has a pivot. The latter property shows that  $A^T$  is full rank, so the argument above shows that  $N(A^T) = \{0\}$ . Since the case  $n = m$  has already been studied, so we only need to consider the case  $m < n$ . In this case there are free variables in the solution to the equation  $Ax = 0$ , therefore  $\{0\} \subsetneq N(A)$ . This proves (b), and we have established the Theorem.  $\square$

**Further reading.** Section 4.2 in Meyer's book [3] follows closely the descriptions of the four spaces given here, although in more depth.



## 2.5.4. Exercises.

2.5.1.- Find  $R(A)$  and  $R(A^T)$  for the matrix  $A$  below and express them as the span of the smallest possible set of vectors, where

$$A = \begin{bmatrix} 1 & 2 & 2 & 3 \\ 2 & 4 & 1 & 3 \\ 3 & 6 & 1 & 4 \end{bmatrix}.$$

2.5.2.- Find the  $N(A)$ ,  $R(A)$ ,  $N(A^T)$  and  $R(A^T)$  for the matrix  $A$  below and express them as the span of the smallest possible set of vectors, where

$$A = \begin{bmatrix} 1 & 2 & 1 & 1 & 5 \\ -2 & -4 & 0 & 4 & -2 \\ 1 & 2 & 2 & 4 & 9 \end{bmatrix}.$$

2.5.3.- Let  $A$  be a  $3 \times 3$  matrix such that

$$R(A) = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix}\right\}\right),$$

$$N(A) = \text{Span}\left(\left\{\begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix}\right\}\right).$$

(a) Show that the linear system  $Ax = \mathbf{b}$  is consistent for

$$\mathbf{b} = \begin{bmatrix} 1 \\ 8 \\ 5 \end{bmatrix}.$$

(b) Show that the system  $Ax = \mathbf{b}$  above has infinitely many solutions.

2.5.4.- Prove:

- (a)  $Ax = \mathbf{b}$  is consistent iff  $\mathbf{b} \in R(A)$ .
- (b) The consistent system  $Ax = \mathbf{b}$  has a unique solution iff  $N(A) = \{\mathbf{0}\}$ .

2.5.5.- Prove: A matrix  $A \in \mathbb{F}^{n,n}$  is invertible iff  $R(A) = \mathbb{F}^n$ .

2.5.6.- Let  $A \in \mathbb{F}^{n,n}$  be an invertible matrix. Find the spaces  $N(A)$ ,  $R(A)$ ,  $N(A^T)$ , and  $R(A^T)$ .

2.5.7.- Consider the matrices

$$A = \begin{bmatrix} 1 & 5 & 3 \\ 2 & 1 & -3 \\ 1 & 3 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 0 & -2 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix}.$$

Answer the questions below. If the answer is *yes*, give a proof; if the answer is *no*, give a counter-example.

- (a) Is  $R(A) = R(B)$ ?
- (b) Is  $R(A^T) = R(B^T)$ ?
- (c) Is  $N(A) = N(B)$ ?
- (d) Is  $N(A^T) = N(B^T)$ ?

## 2.6. LU-FACTORIZATION

A factorization of a number means to decompose that number as a product of appropriate factors. For example, the prime factorization of an integer number means to decompose that integer as a product of prime numbers, like  $70 = (2)(5)(7)$ . In a similar way, a factorization of a matrix means to decompose that matrix, using matrix multiplication, as a product of appropriate factors. In this Section we introduce a particular type of factorization, called LU-factorization, which stands for *lower triangular-upper triangular* factorization. A given matrix  $A$  is expressed as a product  $A = LU$ , where  $L$  is lower triangular and  $U$  is upper triangular. This type of factorization can be useful to solve systems of linear equations having  $A$  as the coefficient matrix. The LU-factorization of  $A$  reduces the number of algebraic operations needed to solve the linear system, saving computer time when solving these systems using a computer. However, not every matrix has this type of factorization. We provide sufficient conditions on  $A$  that guarantee its LU-factorization.

**2.6.1. Main definitions.** We start with few basic definitions.

**Definition 2.6.1.** An  $m \times n$  matrix is called **upper triangular** iff every element below the diagonal vanishes, and **lower triangular** iff every element above the diagonal vanishes.

**EXAMPLE 2.6.1:** The matrices  $U_1$  and  $U_2$  below are upper triangular, while  $L_1$  and  $L_2$  are lower triangular,

$$U_1 = \begin{bmatrix} 2 & 3 & 4 \\ 0 & 5 & 6 \\ 0 & 0 & 1 \end{bmatrix}, \quad U_2 = \begin{bmatrix} 2 & 3 & 4 & 3 \\ 0 & 5 & 6 & 2 \\ 0 & 0 & 1 & 5 \end{bmatrix}, \quad L_1 = \begin{bmatrix} 2 & 0 & 0 \\ 3 & 4 & 0 \\ 5 & 6 & 7 \end{bmatrix}, \quad L_2 = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 3 & 4 & 0 & 0 \\ 5 & 6 & 7 & 0 \end{bmatrix}.$$

&lt;

**Definition 2.6.2.** An  $m \times n$  matrix  $A$  has an **LU-factorization** iff there exists a square  $m \times m$  lower triangular matrix  $L$  and an  $m \times n$  upper triangular matrix  $U$  such that  $A = LU$ .

In the particular case that  $A$  is a square matrix both  $L$  and  $U$  are square matrices.

**EXAMPLE 2.6.2:** Verify that matrix  $L = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix}$  and matrix  $U = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$  are the

LU-factorization of matrix  $A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 2 \\ 1 & 2 & 3 \end{bmatrix}$ .

**SOLUTION:** We need to verify that  $A = LU$ . This is indeed the case, since

$$LU = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 2 \\ 1 & 2 & 3 \end{bmatrix} = A,$$

that is,  $A = LU$ .

&lt;

**EXAMPLE 2.6.3:** Verify that matrix  $L = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & -3 & 1 \end{bmatrix}$  and matrix  $U = \begin{bmatrix} 2 & 4 & -1 \\ 0 & 3 & 1 \\ 0 & 0 & 0 \end{bmatrix}$  are

the LU-factorization of matrix  $A = \begin{bmatrix} 2 & 4 & -1 \\ -4 & -5 & 3 \\ 2 & -5 & -4 \end{bmatrix}$ .

**SOLUTION:** We need to verify that  $A = LU$ . This is the case, since

$$LU = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & -3 & 1 \end{bmatrix} \begin{bmatrix} 2 & 4 & -1 \\ 0 & 3 & 1 \\ 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 4 & -1 \\ -4 & -5 & 3 \\ 2 & -5 & -4 \end{bmatrix} = A \Rightarrow A = LU,$$

◁

**EXAMPLE 2.6.4:** Verify that matrix  $L = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 3 & 2 & 1 \end{bmatrix}$  and matrix  $U = \begin{bmatrix} 1 & 4 \\ 0 & -3 \\ 0 & 0 \end{bmatrix}$  are the

LU-factorization of matrix  $A = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}$ .

**SOLUTION:** We need to verify that  $A = LU$ . This is the case, since

$$LU = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} 1 & 4 \\ 0 & -3 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix} = A \Rightarrow A = LU.$$

◁

**EXAMPLE 2.6.5:** Verify that matrix  $L = \begin{bmatrix} 1 & 0 \\ 2 & 1 \end{bmatrix}$  and matrix  $U = \begin{bmatrix} 1 & 3 & 5 \\ 0 & -2 & -4 \end{bmatrix}$  are the

LU-factorization of matrix  $A = \begin{bmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \end{bmatrix}$ .

**SOLUTION:** We need to verify that  $A = LU$ . This is the case, since

$$LU = \begin{bmatrix} 1 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 & 3 & 5 \\ 0 & -2 & -4 \end{bmatrix} = \begin{bmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \end{bmatrix} = A \Rightarrow A = LU.$$

◁

**2.6.2. A sufficient condition.** We now provide sufficient conditions on a matrix that imply that such matrix has an LU-factorization.

**Theorem 2.6.3.** *If the  $m \times n$  matrix  $A$  can be transformed into an echelon form without row exchanges and without row rescaling, then  $A$  has an LU-factorization,  $A = LU$ . The matrix  $U$  is the  $m \times n$  echelon form mentioned above. Matrix  $L = [L_{ij}]$  is an  $m \times m$  lower triangular matrix satisfying two conditions: First, the coefficients  $L_{ii} = 1$ ; second, the coefficients  $L_{ij}$  below the diagonal are equal to the multiple of row  $j$  that is subtracted from row  $i$  in order to annihilate the  $(i, j)$  position during the Gauss elimination method.*

This result is usually found in the literature for the case of square matrices,  $m = n$ . A reason for this is that more often than not only square matrices appear in applications, like finding solutions of a  $n \times n$  linear system  $Ax = b$ , where matrix  $A$  is not only square but also invertible. We have decided in these notes to present the most general version of the LU-factorization in the statement above without proof, and we sketch a proof in the case of  $2 \times 2$  and  $3 \times 3$  matrices only. Before this sketch of a proof we show few examples in order to have a better understanding of the statement in Theorem 2.6.3.

**EXAMPLE 2.6.6:** Find the LU-factorization of matrix  $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ .

**SOLUTION:** Theorem 2.6.3 says that  $U$  is an echelon form of  $A$  obtained without row exchanges and without row rescaling, while  $L$  has the form  $L = \begin{bmatrix} 1 & 0 \\ L_{21} & 1 \end{bmatrix}$ , that is, we need to

find only the coefficient  $L_{21}$ . In this simple example we can summarize the whole procedure in the following diagram:

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \xrightarrow{\text{row}(2) - 3\text{row}(1)} \begin{bmatrix} 1 & 2 \\ 0 & -2 \end{bmatrix} = U, \quad L = \begin{bmatrix} 1 & 0 \\ 3 & 1 \end{bmatrix}.$$

This is what we have done: An echelon form for  $A$  is obtained with only one Gauss operation: Subtract from the second row the first row multiplied by 3. The resulting echelon form of  $A$  is matrix  $U$  already. And  $L_{21} = 3$ , since we multiplied by 3 the first row and we subtracted it from the second row. So we have both  $U$  and  $L$ , and the result is:

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 0 & -2 \end{bmatrix} = LU.$$

◁

**EXAMPLE 2.6.7:** Find the LU-factorization of the matrix  $A = \begin{bmatrix} 2 & 4 & -1 \\ -4 & -5 & 3 \\ 2 & -5 & -4 \end{bmatrix}$ .

**SOLUTION:** We recall that we have to find the coefficients below the diagonal in matrix

$$L = \begin{bmatrix} 1 & 0 & 0 \\ L_{21} & 1 & 0 \\ L_{31} & L_{32} & 1 \end{bmatrix}.$$

From the first Gauss operation we obtain  $L_{21}$  as follows:

$$\begin{bmatrix} 2 & 4 & -1 \\ -4 & -5 & 3 \\ 2 & -5 & -4 \end{bmatrix} \xrightarrow{\text{row}(2) - (-2)\text{row}(1)} \begin{bmatrix} 2 & 4 & -1 \\ 0 & 3 & 1 \\ 2 & -5 & -4 \end{bmatrix} \Rightarrow L = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ L_{31} & L_{32} & 1 \end{bmatrix}.$$

From the second Gauss operation we obtain  $L_{31}$  as follows:

$$\begin{bmatrix} 2 & 4 & -1 \\ 0 & 3 & 1 \\ 2 & -5 & -4 \end{bmatrix} \xrightarrow{\text{row}(3) - 1\text{row}(1)} \begin{bmatrix} 2 & 4 & -1 \\ 0 & 3 & 1 \\ 0 & -9 & -3 \end{bmatrix} \Rightarrow L = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & L_{32} & 1 \end{bmatrix}.$$

From the third Gauss operation we obtain  $L_{32}$  as follows:

$$\begin{bmatrix} 2 & 4 & -1 \\ 0 & 3 & 1 \\ 0 & -9 & -3 \end{bmatrix} \xrightarrow{\text{row}(3) - (-3)\text{row}(2)} \begin{bmatrix} 2 & 4 & -1 \\ 0 & 3 & 1 \\ 0 & 0 & 0 \end{bmatrix} \Rightarrow L = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & -3 & 1 \end{bmatrix}.$$

We then conclude,

$$L = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & -3 & 1 \end{bmatrix}, \quad U = \begin{bmatrix} 2 & 4 & -1 \\ 0 & 3 & 1 \\ 0 & 0 & 0 \end{bmatrix},$$

and we have the decomposition

$$A = \begin{bmatrix} 2 & 4 & -1 \\ -4 & -5 & 3 \\ 2 & -5 & -4 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & -3 & 1 \end{bmatrix} \begin{bmatrix} 2 & 4 & -1 \\ 0 & 3 & 1 \\ 0 & 0 & 0 \end{bmatrix} = LU.$$

◁

In the following Example we show a matrix  $A$  that does not have an LU-factorization. The reason is that in the procedure to find  $U$  appears a diagonal coefficient that vanishes. Since row interchanges are prohibited, then there is no LU-factorization in this case.

**EXAMPLE 2.6.8:** Show that matrix  $A = \begin{bmatrix} 2 & 4 & -1 \\ -4 & -8 & 3 \\ 2 & -5 & -4 \end{bmatrix}$  has no LU-factorization.

**SOLUTION:** From the first Gauss operation we obtain  $L_{21}$  as follows:

$$\begin{bmatrix} 2 & 4 & -1 \\ -4 & -8 & 3 \\ 2 & -5 & -4 \end{bmatrix} \xrightarrow{\text{row}(2) - (-2)\text{row}(1)} \begin{bmatrix} 2 & 4 & -1 \\ 0 & 0 & 1 \\ 2 & -5 & -4 \end{bmatrix} \Rightarrow L = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ L_{31} & L_{32} & 1 \end{bmatrix}.$$

From the second Gauss operation we obtain  $L_{31}$  as follows:

$$\begin{bmatrix} 2 & 4 & -1 \\ 0 & 0 & 1 \\ 2 & -5 & -4 \end{bmatrix} \xrightarrow{\text{row}(3) - 1\text{row}(1)} \begin{bmatrix} 2 & 4 & -1 \\ 0 & 0 & 1 \\ 0 & -9 & -3 \end{bmatrix} \Rightarrow L = \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & L_{32} & 1 \end{bmatrix}.$$

However, we cannot continue the Gauss method to find an echelon form for  $A$  without interchanging the second and third rows. Therefore, **matrix  $A$  has no LU-factorization.**  $\triangleleft$

In the Example 2.6.7 we also had a diagonal coefficient in matrix  $U$  that vanishes, the coefficient  $U_{33} = 0$ . However, in this case  $A$  did have an LU-factorization, since this vanishing coefficient was in the last row of matrix  $U$ , and no further Gauss operations were needed.

**Proof of Theorem 2.6.3:** We only give a proof in the case that matrix  $A$  is  $2 \times 2$  or  $3 \times 3$ . This would give an idea how to construct a proof in the general case. This generalization does not involve new ideas, only a more sophisticated notation.

Assume that matrix  $A$  is  $2 \times 2$ , that is

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}.$$

Matrix  $A$  is assumed to satisfy the following property:  $A$  can be transformed into echelon form by the only Gauss operation of multiplying a row and add that result to another row. If the coefficient  $A_{21} = 0$ , then matrix  $A$  is upper triangular already and it has the trivial LU-factorization  $A = I_2 A$ . If the coefficient  $A_{21} \neq 0$ , then from Example 2.6.8 we know that the assumption on  $A$  implies that the coefficient  $A_{11} \neq 0$ . Then, we can perform the Gauss operation  $EA = U$ , that is

$$EA = \begin{bmatrix} 1 & 0 \\ -\frac{A_{21}}{A_{11}} & 1 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} \\ 0 & \frac{A_{22}A_{11} - A_{21}A_{12}}{A_{11}} \end{bmatrix} = U.$$

Matrix  $U$  is the upper triangular, and denoting

$$U_{22} = \frac{A_{22}A_{11} - A_{21}A_{12}}{A_{11}}, \quad L_{21} = \frac{A_{21}}{A_{11}},$$

we obtain that

$$U = \begin{bmatrix} A_{11} & A_{12} \\ 0 & U_{22} \end{bmatrix}, \quad E = \begin{bmatrix} 1 & 0 \\ -L_{21} & 1 \end{bmatrix} \Rightarrow E^{-1} = \begin{bmatrix} 1 & 0 \\ L_{21} & 1 \end{bmatrix}.$$

We conclude that matrix  $A$  has an LU-factorization

$$A = \begin{bmatrix} 1 & 0 \\ L_{21} & 1 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ 0 & U_{22} \end{bmatrix}.$$

Assume now that matrix  $A$  is  $3 \times 3$ , that is

$$A = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix}.$$

Once again, matrix  $A$  is assumed to satisfy the following property:  $A$  can be transformed into echelon form by the only Gauss operation of multiplying a row and add that result to another row. If any of the coefficients  $A_{21}$  or  $A_{13}$  is non-zero, then the assumption on  $A$  implies that the coefficient  $A_{11} \neq 0$ . Then, we can perform the Gauss operation  $E_2 E_1 A = B$ , that is

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -L_{31} & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ -L_{21} & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & B_{22} & B_{23} \\ 0 & B_{32} & B_{33} \end{bmatrix},$$

where

$$E_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -L_{31} & 0 & 1 \end{bmatrix}, \quad E_1 = \begin{bmatrix} 1 & 0 & 0 \\ -L_{21} & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

we used the notation

$$L_{21} = \frac{A_{21}}{A_{11}}, \quad L_{31} = \frac{A_{31}}{A_{11}},$$

and the  $B_{ij}$  coefficients can be easily computed. Now, if the coefficient  $B_{32} \neq 0$ , then the assumption on  $A$  implies that the coefficient  $B_{22} \neq 0$ . We then assume that  $B_{22} \neq 0$  and we proceed one more step. We can perform the Gauss operation  $E_3 B = U$ , that is

$$E_3 B = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -L_{31} & 1 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & B_{22} & B_{23} \\ 0 & B_{32} & B_{33} \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & B_{22} & B_{23} \\ 0 & 0 & U_{33} \end{bmatrix} = U,$$

where we used the notation  $L_{31} = \frac{B_{32}}{B_{22}}$  and the coefficient  $U_{33}$  can be easily computed. This product can be expressed as follows,

$$E_3 E_2 E_1 A = U \quad \Rightarrow \quad A = E_1^{-1} E_2^{-1} E_3^{-1} U.$$

It is not difficult to see that

$$E_1^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ L_{21} & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad E_2^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ L_{31} & 0 & 1 \end{bmatrix}, \quad E_3^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & L_{32} & 1 \end{bmatrix},$$

which implies that

$$E_1^{-1} E_2^{-1} E_3^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ L_{21} & 1 & 0 \\ L_{31} & L_{32} & 1 \end{bmatrix} = L.$$

We have then shown that matrix  $A$  has the LU-factorization  $A = LU$ , where

$$L = \begin{bmatrix} 1 & 0 & 0 \\ L_{21} & 1 & 0 \\ L_{31} & L_{32} & 1 \end{bmatrix}, \quad U = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ 0 & B_{22} & B_{23} \\ 0 & 0 & U_{33} \end{bmatrix}.$$

It is not difficult to see that in the case where the coefficient  $B_{32} = 0$ , then the expression  $A = (E_1^{-1} E_2^{-1}) B$  is already the LU-factorization of matrix  $A$ . This establishes the Theorem for  $2 \times 2$  and  $3 \times 3$  matrices.  $\square$

**2.6.3. Solving linear systems.** Suppose we want to solve a linear system with coefficient matrix  $A$ , and suppose that this matrix has an LU-factorization. This extra information is useful to save computational time when solving the linear system. We state how this can be done in the following result.

**Theorem 2.6.4.** Fix an  $m \times n$  matrix  $A$  having an LU-factorization  $A = LU$  and a vector  $\mathbf{b} \in \mathbb{F}^m$ . The vector  $\mathbf{x} \in \mathbb{F}^n$  is solution of the linear system  $A\mathbf{x} = \mathbf{b}$  iff holds

$$U\mathbf{x} = \mathbf{y}, \quad \text{where } L\mathbf{y} = \mathbf{b}. \quad (2.8)$$

First solve the second equation for  $\mathbf{y}$ , and then use this vector  $\mathbf{y}$  to solve for vector  $\mathbf{x}$ . This is faster than solving directly for  $\mathbf{x}$  in the original equation  $A\mathbf{x} = \mathbf{b}$ . The reason is that forward substitution can be used to solve for  $\mathbf{y}$ , since  $L$  is lower triangular. Then, backward substitution can be used to solve for  $\mathbf{x}$ , since  $U$  is upper triangular.

**Proof of Theorem 2.6.4:** Replace on the second equation in (2.8) the vector  $\mathbf{y}$  defined by the first equation in (2.8). Hence, the vector  $\mathbf{x}$  is solution of the system in (2.8) iff holds

$$L(U\mathbf{x}) = \mathbf{b}.$$

Since  $A = LU$ , the equation above is equivalent to  $A\mathbf{x} = \mathbf{b}$ . This establishes the Theorem.  $\square$

**EXAMPLE 2.6.9:** Use the LU-factorization of matrix  $A$  below, to find the solution  $\mathbf{x}$  to the system  $A\mathbf{x} = \mathbf{b}$ , where

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 2 \\ 1 & 2 & 3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}.$$

**SOLUTION:** In Example 2.6.2 we have shown that  $A = LU$  with

$$L = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix}, \quad U = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}.$$

We first find  $\mathbf{y}$  solution of the system  $L\mathbf{y} = \mathbf{b}$ , then, having  $\mathbf{y}$ , we find  $\mathbf{x}$  solution of  $U\mathbf{x} = \mathbf{y}$ . For the first system we have

$$\left[ \begin{array}{ccc|c} 1 & 0 & 0 & 1 \\ 1 & 1 & 0 & 2 \\ 1 & 1 & 1 & 3 \end{array} \right] \Rightarrow \begin{cases} y_1 = 1, \\ y_1 + y_2 = 2, \\ y_1 + y_2 + y_3 = 3, \end{cases} \Rightarrow \mathbf{y} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}.$$

This system is solved for  $\mathbf{y}$  using forward substitution. For the second system we have

$$\left[ \begin{array}{ccc|c} 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{array} \right] \Rightarrow \begin{cases} x_1 + x_2 + x_3 = 1, \\ x_2 + x_3 = 1, \\ x_3 = 1, \end{cases} \Rightarrow \mathbf{x} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.$$

This system is solved for  $\mathbf{x}$  using back substitution.  $\triangleleft$

**Further reading.** There exists a vast literature on matrix factorization. See Section 3.10 in Meyer's book [3] for few types of generalizations of the LU-factorizations, for example, admitting row interchanges in the process to obtain the factorization, or the LDU-factorization.

**2.6.4. Exercises.****2.6.1.-** Find the LU-factorization of matrix

$$A = \begin{bmatrix} 5 & 2 \\ -15 & -3 \end{bmatrix}.$$

**2.6.2.-** Find the LU-factorization of matrix

$$A = \begin{bmatrix} 2 & 1 & 3 \\ 4 & 6 & 7 \end{bmatrix}.$$

**2.6.3.-** Find the LU-factorization of matrix

$$A = \begin{bmatrix} 2 & 1 & 2 \\ 4 & 5 & 5 \\ 6 & -3 & 5 \end{bmatrix}.$$

**2.6.4.-** Determine if the matrix below has an LU-factorization,

$$A = \begin{bmatrix} 1 & 2 & 4 & 17 \\ 3 & 6 & -12 & 3 \\ 2 & 3 & -3 & 2 \\ 0 & 2 & -2 & 6 \end{bmatrix}.$$

**2.6.5.-** Find the LU-factorization of a tridiagonal matrix

$$T = \begin{bmatrix} 2 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 1 \end{bmatrix}.$$

**2.6.6.-** Use the LU-factorization to find the solutions to the system  $Ax = b$ , where

$$A = \begin{bmatrix} 2 & 2 & 2 \\ 4 & 7 & 7 \\ 6 & 18 & 22 \end{bmatrix}, \quad b = \begin{bmatrix} 12 \\ 24 \\ 12 \end{bmatrix}.$$

**2.6.7.-** Find the values of the number  $c$  such that matrix  $A$  below has no LU-factorization,

$$A = \begin{bmatrix} c & 2 & 0 \\ 1 & c & 1 \\ 0 & 1 & c \end{bmatrix}.$$



## CHAPTER 3. DETERMINANTS

## 3.1. DEFINITIONS AND PROPERTIES

A determinant is a scalar associated to a square matrix which can be used to determine whether the matrix is invertible or not, and this property is the origin of its name. The determinant has a clear geometrical meaning in the case of  $2 \times 2$  and  $3 \times 3$  matrices. In the former case the absolute value of the determinant is the area of a parallelogram formed with the matrix column vectors; in the latter case the absolute value of the determinant is the volume of the parallelepiped formed with the matrix column vectors. The determinant for  $n \times n$  matrices is introduced as a suitable generalization of these properties. In this Section we present the determinant for  $2 \times 2$  and  $3 \times 3$  matrices and we study their main properties. We then present the definition of determinant for  $n \times n$  matrices and we mention without proof its main properties.

**3.1.1. Determinant of  $2 \times 2$  matrices.** We start introducing the determinant as a scalar-valued function on  $2 \times 2$  matrices.

**Definition 3.1.1.** The *determinant* of a  $2 \times 2$  matrix  $A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$  is the value of the function  $\det : \mathbb{F}^{2,2} \rightarrow \mathbb{F}$  given by  $\det(A) = A_{11}A_{22} - A_{12}A_{21}$ .

Depending on the context we will use for the determinant any of the following notations,

$$\det(A) = |A| = \Delta = \begin{vmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{vmatrix}.$$

For example,  $\begin{vmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{vmatrix} = A_{11}A_{22} - A_{12}A_{21}$ .

**EXAMPLE 3.1.1:** The value of the determinant can be any real number, as the following three cases show:

$$\begin{vmatrix} 1 & 2 \\ 3 & 4 \end{vmatrix} = 4 - 6 = -2, \quad \begin{vmatrix} 2 & 1 \\ 3 & 4 \end{vmatrix} = 8 - 3 = 5, \quad \begin{vmatrix} 1 & 2 \\ 2 & 4 \end{vmatrix} = 4 - 4 = 0.$$

◁

We have seen in Theorem 2.4.2 in Section 2.4 that a  $2 \times 2$  matrix is invertible iff its determinant is non-zero. We now see that there is an interesting geometrical interpretation of this property.

**Theorem 3.1.2.** Given a matrix  $A = [A_1, A_2] \in \mathbb{R}^{2,2}$ , the absolute value of its determinant,  $|\det(A)|$ , is the area of the parallelogram formed by the vectors  $A_1, A_2$ .

**Proof of Theorem 3.1.2:** Denote the matrix  $A$  as follows

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \Rightarrow A_1 = \begin{bmatrix} a \\ c \end{bmatrix}, \quad A_2 = \begin{bmatrix} b \\ d \end{bmatrix}.$$

The case where all the coefficients in matrix  $A$  are positive is shown in Fig. 27. We can see in this Figure that the area of the parallelogram formed by the vectors  $A_1$  and  $A_2$  is related to the area of the rectangle with sides  $b$  and  $c$ . More precisely, the area of the parallelogram is equal to the area of the rectangle minus the area of the two triangles marked with a “−” sign in Fig. 27 plus the area of the triangle marked with “+” sign. Denoting the area of the parallelogram by  $A_p$ , we obtain the following equation,

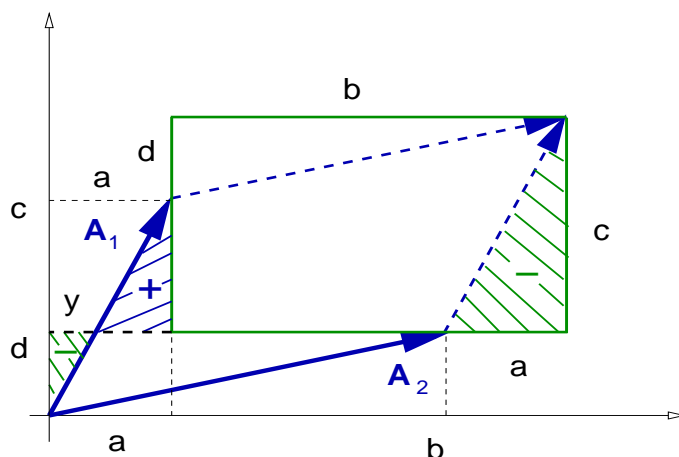


FIGURE 27. The geometrical meaning of the determinant of a  $2 \times 2$  matrix is that its absolute value is the area of the parallelogram formed by the matrix column vectors.

$$\begin{aligned}
 A_p &= bc - \frac{ac}{2} - \frac{yd}{2} + \frac{(a-y)(c-d)}{2} \\
 &= bc - \frac{ac}{2} - \frac{yd}{2} + \frac{ac}{2} - \frac{ad}{2} - \frac{yc}{2} + \frac{yd}{2} \\
 &= bc - \frac{ad}{2} - \frac{yc}{2}.
 \end{aligned}$$

Similar triangles implies that

$$\frac{y}{d} = \frac{a}{c} \Rightarrow yc = ad.$$

Introducing this relation in the equation above it we obtain

$$A_p = bc - \frac{ad}{2} - \frac{ad}{2} \Rightarrow A_p = ad - bc = |\det(A)|.$$

We consider only this case in our proof. The remaining cases can be studied in a similar way. This establishes the Theorem.  $\square$

**EXAMPLE 3.1.2:** Find the area of the parallelogram formed by  $\mathbf{a} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$  and  $\mathbf{b} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$ .

**SOLUTION:** If we consider the matrix  $A = [\mathbf{a}, \mathbf{b}] = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$ , then the area  $A$  of the parallelogram formed by the vectors  $\mathbf{a}$ ,  $\mathbf{b}$  is  $A = |\det(A)|$ , that is,

$$\det(A) = \begin{vmatrix} 1 & 2 \\ 2 & 1 \end{vmatrix} = 1 - 4 = -3 \Rightarrow A = 3.$$

$\triangleleft$

We have seen that the determinant is a scalar-valued function on the space of  $2 \times 2$  matrices, whose absolute value is the area of the parallelogram formed by the matrix column vectors. This relation of the determinant with areas on the plane is at the origin of the following properties.

**Theorem 3.1.3.** For every vectors  $\mathbf{a}_1, \mathbf{a}_2, \mathbf{b} \in \mathbb{F}^2$  and scalar  $k \in \mathbb{F}$ , the determinant function  $\det : \mathbb{F}^{2,2} \rightarrow \mathbb{F}$  satisfies:

- (a)  $\det([\mathbf{a}_1, \mathbf{a}_2]) = -\det([\mathbf{a}_2, \mathbf{a}_1])$ ;
- (b)  $\det([k\mathbf{a}_1, \mathbf{a}_2]) = \det([\mathbf{a}_1, k\mathbf{a}_2]) = k \det([\mathbf{a}_1, \mathbf{a}_2])$ ;
- (c)  $\det([\mathbf{a}_1, k\mathbf{a}_1]) = 0$ ;
- (d)  $\det([\mathbf{a}_1 + \mathbf{b}, \mathbf{a}_2]) = \det([\mathbf{a}_1, \mathbf{a}_2]) + \det([\mathbf{b}, \mathbf{a}_2])$ .

**Proof of Theorem 3.1.3:** Introduce the notation

$$\mathbf{a}_1 = \begin{bmatrix} a \\ c \end{bmatrix}, \quad \mathbf{a}_2 = \begin{bmatrix} b \\ d \end{bmatrix}, \quad \Rightarrow \quad [\mathbf{a}_1, \mathbf{a}_2] = \begin{bmatrix} a & b \\ c & d \end{bmatrix}.$$

Part (a):

$$\left. \begin{aligned} \det([\mathbf{a}_1, \mathbf{a}_2]) &= \begin{vmatrix} a & b \\ c & d \end{vmatrix} = ad - bc \\ \det([\mathbf{a}_2, \mathbf{a}_1]) &= \begin{vmatrix} b & a \\ d & c \end{vmatrix} = bc - ad, \end{aligned} \right\} \Rightarrow \det([\mathbf{a}_1, \mathbf{a}_2]) = -\det([\mathbf{a}_2, \mathbf{a}_1]).$$

Part (b):

$$\begin{aligned} \det([k\mathbf{a}_1, \mathbf{a}_2]) &= \begin{vmatrix} ka & b \\ kc & d \end{vmatrix} = k(ad - bc) = k \begin{vmatrix} a & b \\ c & d \end{vmatrix} = k \det([\mathbf{a}_1, \mathbf{a}_2]), \\ \det([\mathbf{a}_1, k\mathbf{a}_2]) &= \begin{vmatrix} a & kb \\ c & kd \end{vmatrix} = k(ad - bc) = k \begin{vmatrix} a & b \\ c & d \end{vmatrix} = k \det([\mathbf{a}_1, \mathbf{a}_2]). \end{aligned}$$

Part (c):

$$\det([\mathbf{a}_1, k\mathbf{a}_1]) = \det([k\mathbf{a}_1, \mathbf{a}_1]) = -\det([\mathbf{a}_1, k\mathbf{a}_1]) \Rightarrow \det([\mathbf{a}_1, k\mathbf{a}_1]) = 0.$$

Part (d): Introduce the notation  $\mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$ . Then,

$$\begin{aligned} \det([\mathbf{a}_1 + \mathbf{b}, \mathbf{a}_2]) &= \begin{vmatrix} (a + b_1) & b \\ (c + b_2) & d \end{vmatrix} \\ &= (a + b_1)d - (c + b_2)b \\ &= (ad - cb) + (b_1d - b_2b) \\ &= \begin{vmatrix} a & b \\ c & d \end{vmatrix} + \begin{vmatrix} b_1 & b \\ b_2 & d \end{vmatrix} \\ &= \det([\mathbf{a}_1, \mathbf{a}_2]) + \det([\mathbf{b}, \mathbf{a}_2]). \end{aligned}$$

This establishes the Theorem. □

Moreover, the determinant function also satisfies the following further properties.

**Theorem 3.1.4.** Let  $\mathbf{A}, \mathbf{B} \in \mathbb{F}^{2,2}$ . Then it holds:

- (a) Matrix  $\mathbf{A}$  is invertible iff  $\det(\mathbf{A}) \neq 0$ ;
- (b)  $\det(\mathbf{A}) = \det(\mathbf{A}^T)$ ;
- (c)  $\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B})$ .
- (d) If matrix  $\mathbf{A}$  is invertible, then  $\det(\mathbf{A}^{-1}) = \frac{1}{\det(\mathbf{A})}$ .

**Proof of Theorem 3.1.4:**

Part (a): It has been proven in Theorem 2.4.2.

Part (b): Denote  $\mathbf{A} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ . Then holds

$$\det(\mathbf{A}) = \begin{vmatrix} a & b \\ c & d \end{vmatrix} = ad - bc = \begin{vmatrix} a & c \\ b & d \end{vmatrix} = \det(\mathbf{A}^T).$$

Part (c): Denoting  $\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$  and  $\mathbf{B} = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}$  we obtain,

$$\mathbf{AB} = \begin{bmatrix} (A_{11}B_{11} + A_{12}B_{21}) & (A_{11}B_{12} + A_{12}B_{22}) \\ (A_{21}B_{11} + A_{22}B_{21}) & (A_{21}B_{12} + A_{22}B_{22}) \end{bmatrix}.$$

Therefore,

$$\begin{aligned} \det(\mathbf{AB}) &= (A_{11}B_{11} + A_{12}B_{21})(A_{21}B_{12} + A_{22}B_{22}) - (A_{11}B_{12} + A_{12}B_{22})(A_{21}B_{11} + A_{22}B_{21}) \\ &= A_{11}B_{11}A_{21}B_{12} + A_{11}B_{11}A_{22}B_{22} + A_{12}B_{21}A_{21}B_{12} + A_{12}B_{21}A_{22}B_{22} \\ &\quad - A_{21}B_{11}A_{11}B_{12} - A_{21}B_{11}A_{12}B_{22} - A_{22}B_{21}A_{11}B_{12} - A_{22}B_{21}A_{12}B_{22} \\ &= A_{11}B_{11}A_{22}B_{22} + A_{12}B_{21}A_{21}B_{12} - A_{21}B_{11}A_{12}B_{22} - A_{22}B_{21}A_{11}B_{12} \\ &= (A_{11}A_{22} - A_{21}A_{12})B_{11}B_{22} - (A_{11}A_{22} - A_{12}A_{21})B_{12}B_{21} \\ &= (A_{11}A_{22} - A_{21}A_{12})(B_{11}B_{22} - B_{12}B_{21}) \\ &= \det(\mathbf{A}) \det(\mathbf{B}). \end{aligned}$$

Part (d): Since  $\mathbf{A}(\mathbf{A}^{-1}) = \mathbf{I}_2$ , we obtain

$$1 = \det(\mathbf{I}_2) = \det(\mathbf{A}(\mathbf{A}^{-1})) = \det(\mathbf{A}) \det(\mathbf{A}^{-1}) \Rightarrow \det(\mathbf{A}^{-1}) = \frac{1}{\det(\mathbf{A})}.$$

This establishes the Theorem.  $\square$

We finally mention that  $\det(\mathbf{A}) = \det(\mathbf{A}^T)$  implies that Theorem 3.1.3 can be generalized from properties involving columns of the matrix into properties involving rows of the matrix. Introduce the following notation that generalizes column vectors to include row vectors:

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = [\mathbf{a}_{:1}, \mathbf{a}_{:2}] = \begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \end{bmatrix},$$

where

$$\mathbf{a}_{:1} = \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix}, \quad \mathbf{a}_{:2} = \begin{bmatrix} A_{12} \\ A_{22} \end{bmatrix}, \quad \mathbf{a}_{1:} = [A_{11} \quad A_{12}], \quad \mathbf{a}_{2:} = [A_{21} \quad A_{22}].$$

The first two vectors above are the usual column vectors of matrix  $\mathbf{A}$ , and the last two are its row vectors. When working with both, column vectors and row vectors, we use the notation above; when working only with column vectors, we drop the colon and we denote them, for example, as  $\mathbf{A}_{:1} = \mathbf{A}_1$ . Using this notation, it is simple to verify the following result.

**Theorem 3.1.5.** For all column vectors  $\mathbf{a}_1^T, \mathbf{a}_2^T, \mathbf{b}_1^T \in \mathbb{F}^2$  and scalars  $k \in \mathbb{F}$ , the determinant function  $\det : \mathbb{F}^{2,2} \rightarrow \mathbb{F}$  satisfies:

- (a)  $\det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \end{bmatrix}\right) = -\det\left(\begin{bmatrix} \mathbf{a}_{2:} \\ \mathbf{a}_{1:} \end{bmatrix}\right);$
- (b)  $\det\left(\begin{bmatrix} k \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \end{bmatrix}\right) = \det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ k \mathbf{a}_{2:} \end{bmatrix}\right) = k \det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \end{bmatrix}\right);$
- (c)  $\det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ k \mathbf{a}_{1:} \end{bmatrix}\right) = 0;$
- (d)  $\det\left(\begin{bmatrix} \mathbf{a}_{1:} + \mathbf{b}_{1:} \\ \mathbf{a}_{2:} \end{bmatrix}\right) = \det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \end{bmatrix}\right) + \det\left(\begin{bmatrix} \mathbf{b}_{1:} \\ \mathbf{a}_{2:} \end{bmatrix}\right).$

The proof is left as an exercise.

**3.1.2. Determinant of  $3 \times 3$  matrices.** The determinant for  $3 \times 3$  matrices is defined recursively in terms of  $2 \times 2$  determinants.

**Definition 3.1.6.** The *determinant* on  $3 \times 3$  matrices is the function  $\det : \mathbb{F}^{3,3} \rightarrow \mathbb{F}$ ,

$$\det(\mathbf{A}) = \begin{vmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{vmatrix} = A_{11} \begin{vmatrix} A_{22} & A_{23} \\ A_{32} & A_{33} \end{vmatrix} - A_{12} \begin{vmatrix} A_{21} & A_{23} \\ A_{31} & A_{33} \end{vmatrix} + A_{13} \begin{vmatrix} A_{21} & A_{22} \\ A_{31} & A_{32} \end{vmatrix}.$$

We see that the determinant of a  $3 \times 3$  matrix is computed as a linear combination of the determinants of  $2 \times 2$  matrices,

$$\det(\mathbf{A}) = A_{11} \det(\overset{\circ}{\mathbf{A}}_{11}) - A_{12} \det(\overset{\circ}{\mathbf{A}}_{12}) + A_{13} \det(\overset{\circ}{\mathbf{A}}_{13}),$$

where we introduced the  $2 \times 2$  matrices

$$\overset{\circ}{\mathbf{A}}_{11} = \begin{bmatrix} A_{22} & A_{23} \\ A_{32} & A_{33} \end{bmatrix}, \quad \overset{\circ}{\mathbf{A}}_{12} = \begin{bmatrix} A_{21} & A_{23} \\ A_{31} & A_{33} \end{bmatrix}, \quad \overset{\circ}{\mathbf{A}}_{13} = \begin{bmatrix} A_{21} & A_{22} \\ A_{31} & A_{32} \end{bmatrix};$$

These matrices are three of the nine matrices called the minors of matrix  $\mathbf{A}$ .

**Definition 3.1.7.** The  $2 \times 2$  matrix  $\overset{\circ}{\mathbf{A}}_{ij}$  is called the  *$(i, j)$ -minor* of a  $3 \times 3$  matrix  $\mathbf{A}$  iff  $\overset{\circ}{\mathbf{A}}_{ij}$  is obtained from  $\mathbf{A}$  by removing the row  $i$  and the column  $j$ . The  *$(i, j)$ -cofactor* of matrix  $\mathbf{A}$  is the number  $C_{ij} = (-1)^{(i+j)} \det(\overset{\circ}{\mathbf{A}}_{ij})$ .

These definitions are common in the literature, and they allow to write down a compact expression for the determinant of a  $3 \times 3$  matrix. Indeed, the expression in Definition 3.1.6 now has the form

$$\det(\mathbf{A}) = A_{11}C_{11} + A_{12}C_{12} + A_{13}C_{13}.$$

Of course, all the complexity in the definition of determinant is now hidden in the definition of the cofactors. The latter expression is not simpler than the one in Definition 3.1.6, it only looks simpler.

**EXAMPLE 3.1.3:** Find the determinant of the matrix  $\mathbf{A} = \begin{bmatrix} 1 & 3 & -1 \\ 2 & 1 & 1 \\ 3 & 2 & 1 \end{bmatrix}$ .

**SOLUTION:** We use the definition above, that is,

$$\begin{aligned} \begin{vmatrix} 1 & 3 & -1 \\ 2 & 1 & 1 \\ 3 & 2 & 1 \end{vmatrix} &= 1 \begin{vmatrix} 1 & 1 \\ 2 & 1 \end{vmatrix} - 3 \begin{vmatrix} 2 & 1 \\ 3 & 1 \end{vmatrix} + (-1) \begin{vmatrix} 2 & 1 \\ 3 & 2 \end{vmatrix} \\ &= (1 - 2) - 3(2 - 3) - (4 - 3) \\ &= 1 + 3 - 1 \\ &= 1. \end{aligned}$$

We conclude that  $\det(\mathbf{A}) = 1$ . ◁

As in the  $2 \times 2$  case, the determinant of a real-valued  $3 \times 3$  matrix can be any real number, positive, negative of zero. The absolute value of the determinant has the geometrical meaning, see Fig. 28.

**Theorem 3.1.8.** The number  $|\det(\mathbf{A})|$ , the absolute value of the determinant of matrix  $\mathbf{A} = [\mathbf{A}_{:1}, \mathbf{A}_{:2}, \mathbf{A}_{:3}]$ , is the volume of the parallelepiped formed by the vectors  $\mathbf{A}_{:1}$ ,  $\mathbf{A}_{:2}$ ,  $\mathbf{A}_{:3}$ .

We do not give a prove of this statement. See the references at the end of the Section.

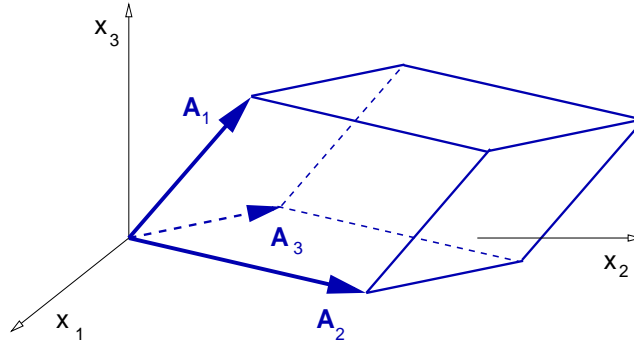


FIGURE 28. The geometrical meaning of the determinant of a  $3 \times 3$  matrix is that its absolute value is the volume of the parallelepiped formed by the matrix column vectors.

**EXAMPLE 3.1.4:** Show that the determinant of an upper triangular or a lower triangular matrix is the product of its diagonal elements.

**SOLUTION:** Denote by  $A$  an upper-triangular matrix. Then, the definition of determinant of a  $3 \times 3$  matrix implies

$$\det(A) = \begin{vmatrix} A_{11} & A_{12} & A_{13} \\ 0 & A_{22} & A_{23} \\ 0 & 0 & A_{33} \end{vmatrix} = A_{11} \begin{vmatrix} A_{22} & A_{23} \\ 0 & A_{33} \end{vmatrix} - A_{12} \begin{vmatrix} 0 & A_{23} \\ 0 & A_{33} \end{vmatrix} + A_{13} \begin{vmatrix} 0 & A_{22} \\ 0 & 0 \end{vmatrix}.$$

We obtain that  $\det(A) = A_{11}A_{22}A_{33}$ . Suppose now that  $A$  is lower-triangular. The definition of determinant of a  $3 \times 3$  matrix implies

$$\det(A) = \begin{vmatrix} A_{11} & 0 & 0 \\ A_{21} & A_{22} & 0 \\ A_{31} & A_{32} & A_{33} \end{vmatrix} = A_{11} \begin{vmatrix} A_{22} & 0 \\ A_{32} & A_{33} \end{vmatrix} - 0 \begin{vmatrix} A_{21} & 0 \\ A_{31} & A_{33} \end{vmatrix} + 0 \begin{vmatrix} A_{21} & A_{22} \\ A_{31} & A_{32} \end{vmatrix}.$$

We obtain that  $\det(A) = A_{11}A_{22}A_{33}$ . ◁

The determinant function on  $3 \times 3$  satisfies a generalization of the properties proven for  $2 \times 2$  matrices in Theorem 3.1.3.

**Theorem 3.1.9.** For all vectors  $\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \mathbf{b} \in \mathbb{F}^3$  and scalars  $k \in \mathbb{F}$ , the determinant function  $\det : \mathbb{F}^{3,3} \rightarrow \mathbb{F}$  satisfies:

- (a)  $\det([\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3]) = -\det([\mathbf{a}_2, \mathbf{a}_1, \mathbf{a}_3]) = -\det([\mathbf{a}_1, \mathbf{a}_3, \mathbf{a}_2]);$
- (b)  $\det([k\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3]) = k \det([\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3]);$
- (c)  $\det([\mathbf{a}_1, k\mathbf{a}_1, \mathbf{a}_3]) = 0;$
- (d)  $\det([\mathbf{a}_1 + \mathbf{b}, \mathbf{a}_2, \mathbf{a}_3]) = \det([\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3]) + \det([\mathbf{b}, \mathbf{a}_2, \mathbf{a}_3]).$

The proof of this Theorem is left as an exercise. The property (a) implies that all the remaining properties (b)-(d) also hold for all the column vectors. That is, from properties (a) and (b) one also shows

$$\det([k\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3]) = \det([\mathbf{a}_1, k\mathbf{a}_2, \mathbf{a}_3]) = \det([\mathbf{a}_1, \mathbf{a}_2, k\mathbf{a}_3]) = k \det([\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3]);$$

from properties (a) and (c) one also shows

$$\det([\mathbf{a}_1, \mathbf{a}_2, k\mathbf{a}_1]) = 0, \quad \det([\mathbf{a}_1, \mathbf{a}_2, k\mathbf{a}_2]) = 0;$$

and from properties (a) and (d) one also shows

$$\begin{aligned}\det([\mathbf{a}_1, (\mathbf{a}_2 + \mathbf{b}), \mathbf{a}_3]) &= \det([\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3]) + \det([\mathbf{a}_1, \mathbf{b}, \mathbf{a}_3]), \\ \det([\mathbf{a}_1, \mathbf{a}_2, (\mathbf{a}_3 + \mathbf{b})]) &= \det([\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3]) + \det([\mathbf{a}_1, \mathbf{a}_2, \mathbf{b}]).\end{aligned}$$

The following result is a generalization of Theorem 3.1.4.

**Theorem 3.1.10.** *Let  $\mathbf{A}, \mathbf{B} \in \mathbb{F}^{3,3}$ . Then it holds:*

- (a) *Matrix  $\mathbf{A}$  is invertible iff  $\det(\mathbf{A}) \neq 0$ ;*
- (b)  *$\det(\mathbf{A}) = \det(\mathbf{A}^T)$ ;*
- (c)  *$\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B})$ .*
- (d) *If matrix  $\mathbf{A}$  is invertible, then  $\det(\mathbf{A}^{-1}) = \frac{1}{\det(\mathbf{A})}$ .*

The proof of this Theorem is left as an exercise. Like in the case of matrices  $2 \times 2$ , Theorem 3.1.9 can be generalized from properties involving column vector to properties involving row vectors.

**Theorem 3.1.11.** *For all vectors  $\mathbf{a}_{1:}^T, \mathbf{a}_{2:}^T, \mathbf{a}_{3:}^T, \mathbf{b}_{1:}^T \in \mathbb{F}^3$  and scalars  $k \in \mathbb{F}$ , the determinant function  $\det : \mathbb{F}^{3,3} \rightarrow \mathbb{F}$  satisfies:*

- (a)  $\det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \\ \mathbf{a}_{3:} \end{bmatrix}\right) = -\det\left(\begin{bmatrix} \mathbf{a}_{2:} \\ \mathbf{a}_{1:} \\ \mathbf{a}_{3:} \end{bmatrix}\right); = -\det\left(\begin{bmatrix} \mathbf{a}_{3:} \\ \mathbf{a}_{2:} \\ \mathbf{a}_{1:} \end{bmatrix}\right); = -\det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{3:} \\ \mathbf{a}_{2:} \end{bmatrix}\right);$
- (b)  $\det\left(\begin{bmatrix} k \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \\ \mathbf{a}_{3:} \end{bmatrix}\right) = k \det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \\ \mathbf{a}_{3:} \end{bmatrix}\right);$
- (c)  $\det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ k \mathbf{a}_{1:} \\ \mathbf{a}_{3:} \end{bmatrix}\right) = 0;$
- (d)  $\det\left(\begin{bmatrix} (\mathbf{a}_{1:} + \mathbf{b}_{1:}) \\ \mathbf{a}_{2:} \\ \mathbf{a}_{3:} \end{bmatrix}\right) = \det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \\ \mathbf{a}_{3:} \end{bmatrix}\right) + \det\left(\begin{bmatrix} \mathbf{b}_{1:} \\ \mathbf{a}_{2:} \\ \mathbf{a}_{3:} \end{bmatrix}\right).$

The proof is left as an exercise. The property (a) implies that all the remaining properties (b)-(d) also hold for all the row vectors. That is, from properties (a) and (b) one also shows

$$\det\left(\begin{bmatrix} k \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \\ \mathbf{a}_{3:} \end{bmatrix}\right) = \det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ k \mathbf{a}_{2:} \\ \mathbf{a}_{3:} \end{bmatrix}\right) = \det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \\ k \mathbf{a}_{3:} \end{bmatrix}\right) = k \det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \\ \mathbf{a}_{3:} \end{bmatrix}\right);$$

from properties (a) and (c) one also shows

$$\det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \\ k \mathbf{a}_{1:} \end{bmatrix}\right) = 0; \quad \det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ k \mathbf{a}_{2:} \\ \mathbf{a}_{2:} \end{bmatrix}\right) = 0;$$

from properties (a) and (d) one also shows

$$\begin{aligned}\det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ (\mathbf{a}_{2:} + \mathbf{b}_{2:}) \\ \mathbf{a}_{3:} \end{bmatrix}\right) &= \det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \\ \mathbf{a}_{3:} \end{bmatrix}\right) + \det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{b}_{2:} \\ \mathbf{a}_{3:} \end{bmatrix}\right), \\ \det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \\ (\mathbf{a}_{3:} + \mathbf{b}_{3:}) \end{bmatrix}\right) &= \det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \\ \mathbf{a}_{3:} \end{bmatrix}\right) + \det\left(\begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \\ \mathbf{b}_{3:} \end{bmatrix}\right).\end{aligned}$$

**3.1.3. Determinant of  $n \times n$  matrices.** The notion of determinant can be generalized to  $n \times n$  matrices with  $n \geq 1$ . One way to find an appropriate generalization is to realize that the determinant of  $2 \times 2$  and  $3 \times 3$  matrices are related to the notion of areas and volumes, respectively. One then studies the properties of areas and volumes, like the ones given in Theorem 3.1.3 and 3.1.9. It can be shown that generalizations of these properties determine a unique function  $\det : \mathbb{F}^{n,n} \rightarrow \mathbb{F}$ . In this notes we only present the final result, the definition of determinant for  $n \times n$  matrices. The reader is referred to the literature for a constructive proof of the function determinant.

**Definition 3.1.12.** The *determinant* of a matrix  $A = [A_{ij}] \in \mathbb{F}^{n,n}$ , with  $n \geq 1$  and  $i, j = 1, \dots, n$ , is the value of the function  $\det : \mathbb{F}^{n,n} \rightarrow \mathbb{F}$  is given by

$$\det(A) = A_{11} C_{11} + \dots + A_{1n} C_{1n}, \quad (3.1)$$

where the number  $C_{ij}$ , called  *$(i, j)$ -cofactor* of matrix  $A$ , is the scalar given by

$$C_{ij} = (-1)^{(i+j)} \det(\overset{\circ}{A}_{ij}),$$

and where the matrices  $\overset{\circ}{A}_{ij} \in \mathbb{F}^{(n-1),(n-1)}$ , called the  *$(i, j)$ -minors* of  $A$ , are obtained from matrix  $A$  by eliminating the row  $i$  and the column  $j$ , which are highlighted below

$$\overset{\circ}{A}_{ij} = \begin{bmatrix} A_{11} & \cdots & \mathbf{A_{1j}} & \cdots & A_{1n} \\ \vdots & & \vdots & & \vdots \\ \mathbf{A_{i1}} & \cdots & \mathbf{A_{ij}} & \cdots & \mathbf{A_{in}} \\ \vdots & & \vdots & & \vdots \\ A_{n1} & \cdots & \mathbf{A_{nj}} & \cdots & A_{nn} \end{bmatrix}.$$

**EXAMPLE 3.1.5:** Use Eq. (3.1) to compute the determinant of a  $2 \times 2$  matrix.

**SOLUTION:** Consider the matrix  $A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$ . The four minor matrices are given by

$$\overset{\circ}{A}_{11} = A_{22}, \quad \overset{\circ}{A}_{12} = A_{21}, \quad \overset{\circ}{A}_{21} = A_{12}, \quad \overset{\circ}{A}_{22} = A_{11}.$$

This means, generalizing the notion of determinant to numbers by  $\det(a) = a$ , that the cofactors are

$$C_{11} = A_{22}, \quad C_{12} = -A_{21}, \quad C_{21} = -A_{12}, \quad C_{22} = A_{11}.$$

So, we obtain that

$$\det(A) = A_{11}C_{11} + A_{12}C_{12} \Leftrightarrow \det(A) = A_{11}A_{22} - A_{12}A_{21}.$$

Notice that we do not need to compute all four cofactors to compute the determinant of the matrix; just two of them are enough.  $\triangleleft$

**EXAMPLE 3.1.6:** Use Eq. (3.1) to compute the determinant of a  $3 \times 3$  matrix.

**SOLUTION:** Consider the matrix

$$A = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix}.$$



The nine minor matrices are given by

$$\begin{array}{lll} \overset{\circ}{A}_{11} = \begin{bmatrix} A_{22} & A_{23} \\ A_{32} & A_{33} \end{bmatrix} & \overset{\circ}{A}_{12} = \begin{bmatrix} A_{21} & A_{23} \\ A_{31} & A_{33} \end{bmatrix} & \overset{\circ}{A}_{13} = \begin{bmatrix} A_{21} & A_{22} \\ A_{31} & A_{32} \end{bmatrix} \\ \overset{\circ}{A}_{21} = \begin{bmatrix} A_{12} & A_{13} \\ A_{32} & A_{33} \end{bmatrix} & \overset{\circ}{A}_{22} = \begin{bmatrix} A_{11} & A_{13} \\ A_{31} & A_{33} \end{bmatrix} & \overset{\circ}{A}_{23} = \begin{bmatrix} A_{11} & A_{12} \\ A_{31} & A_{32} \end{bmatrix} \\ \overset{\circ}{A}_{31} = \begin{bmatrix} A_{12} & A_{13} \\ A_{22} & A_{23} \end{bmatrix} & \overset{\circ}{A}_{32} = \begin{bmatrix} A_{11} & A_{13} \\ A_{21} & A_{23} \end{bmatrix} & \overset{\circ}{A}_{33} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}. \end{array}$$

Then, Def. 3.1.12 agrees with Def. 3.1.6, since

$$\det(\mathbf{A}) = A_{11}C_{11} + A_{12}C_{12} + A_{13}C_{13}$$

is equivalent to

$$\det(\mathbf{A}) = A_{11} \begin{vmatrix} A_{22} & A_{23} \\ A_{32} & A_{33} \end{vmatrix} - A_{12} \begin{vmatrix} A_{21} & A_{23} \\ A_{31} & A_{33} \end{vmatrix} + A_{13} \begin{vmatrix} A_{21} & A_{22} \\ A_{31} & A_{32} \end{vmatrix}.$$

◁

We now state several results that summarize the main properties of the determinant of  $n \times n$  matrices. We state the results without proof. The first result says that the determinant of a matrix can be computed using expansions along any row or any column in the matrix.

**Theorem 3.1.13.** *The determinant of an  $n \times n$  matrix  $\mathbf{A} = [A_{ij}]$  can be computed expanding along any row or any column of matrix  $\mathbf{A}$ , that is,*

$$\begin{aligned} \det(\mathbf{A}) &= A_{i1}C_{i1} + \cdots + A_{in}C_{in}, & i &= 1, \dots, n, \\ &= A_{1j}C_{1j} + \cdots + A_{nj}C_{nj}, & j &= 1, \dots, n. \end{aligned}$$

**EXAMPLE 3.1.7:** Show all possible expansions of the determinant of a  $3 \times 3$  matrix  $\mathbf{A} = [A_{ij}]$ .

**SOLUTION:** The expansions along each of the three rows are the following:

$$\begin{aligned} \det(\mathbf{A}) &= A_{11}C_{11} + A_{12}C_{12} + A_{13}C_{13} \\ &= A_{21}C_{21} + A_{22}C_{22} + A_{23}C_{23} \\ &= A_{31}C_{31} + A_{32}C_{32} + A_{33}C_{33}; \end{aligned}$$

The expansions along each of the three columns are the following:

$$\begin{aligned} \det(\mathbf{A}) &= A_{11}C_{11} + A_{21}C_{21} + A_{31}C_{31} \\ &= A_{12}C_{12} + A_{22}C_{22} + A_{32}C_{32} \\ &= A_{13}C_{13} + A_{23}C_{23} + A_{33}C_{33}. \end{aligned}$$

◁

**EXAMPLE 3.1.8:** Compute the determinant of  $\mathbf{A} = \begin{bmatrix} 1 & 3 & -1 \\ 2 & 1 & 1 \\ 3 & 2 & 1 \end{bmatrix}$  with an expansion by the third column.

**SOLUTION:** The expansion by the third column is the following,

$$\begin{aligned} \begin{vmatrix} 1 & 3 & -1 \\ 2 & 1 & 1 \\ 3 & 2 & 1 \end{vmatrix} &= (-1) \begin{vmatrix} 2 & 1 \\ 3 & 2 \end{vmatrix} - (1) \begin{vmatrix} 1 & 3 \\ 3 & 2 \end{vmatrix} + (1) \begin{vmatrix} 1 & 3 \\ 2 & 1 \end{vmatrix} \\ &= -(4 - 3) - (2 - 9) + (1 - 6) \\ &= -1 + 7 - 5 \\ &= 1, \end{aligned}$$

That is,  $\det(\mathbf{A}) = 1$ . This result agrees with Example 3.1.3.  $\triangleleft$

The generalization of Theorem 3.1.9 to  $n \times n$  matrices is the following.

**Theorem 3.1.14.** For all vectors  $\mathbf{a}_i, \mathbf{a}_j, \mathbf{b} \in \mathbb{F}^n$ , with  $i, j = 1, \dots, n$ , and for all scalars  $k \in \mathbb{F}$ , the determinant function  $\det : \mathbb{F}^{n,n} \rightarrow \mathbb{F}$  satisfies:

- (a)  $\det([\mathbf{a}_1, \dots, \mathbf{a}_i, \dots, \mathbf{a}_j, \dots, \mathbf{a}_n]) = -\det([\mathbf{a}_1, \dots, \mathbf{a}_j, \dots, \mathbf{a}_i, \dots, \mathbf{a}_n])$ ;
- (b)  $\det([\mathbf{a}_1, \dots, k\mathbf{a}_i, \dots, \mathbf{a}_n]) = k \det([\mathbf{a}_1, \dots, \mathbf{a}_i, \dots, \mathbf{a}_n])$ ;
- (c)  $\det([\mathbf{a}_1, \dots, k\mathbf{a}_i, \dots, \mathbf{a}_i, \dots, \mathbf{a}_n]) = 0$ ;
- (d)  $\det([\mathbf{a}_1, \dots, (\mathbf{a}_i + \mathbf{b}), \dots, \mathbf{a}_n]) = \det([\mathbf{a}_1, \dots, \mathbf{a}_i, \dots, \mathbf{a}_n]) + \det([\mathbf{b}, \dots, \mathbf{b}, \dots, \mathbf{a}_n])$ .

The generalization of Theorem 3.1.10 to  $n \times n$  matrices is the following.

**Theorem 3.1.15.** Let  $\mathbf{A}, \mathbf{B} \in \mathbb{F}^{n,n}$ . Then it holds:

- (a) Matrix  $\mathbf{A}$  is invertible iff  $\det(\mathbf{A}) \neq 0$ ;
- (b)  $\det(\mathbf{A}) = \det(\mathbf{A}^T)$ ;
- (c)  $\det(\overline{\mathbf{A}}) = \overline{\det(\mathbf{A})}$ ;
- (d)  $\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B})$ .
- (e) If matrix  $\mathbf{A}$  is invertible, then  $\det(\mathbf{A}^{-1}) = \frac{1}{\det(\mathbf{A})}$ .

The proof of this Proposition is left as an exercise. Like in the case of matrices  $2 \times 2$ , Theorem 3.1.9 can be generalized from properties involving column vector to properties involving row vectors.

**Theorem 3.1.16.** For all vectors  $\mathbf{a}_i^T, \mathbf{a}_j^T, \mathbf{b}_i^T \in \mathbb{F}^n$  and scalars  $k \in \mathbb{F}$ , the determinant function  $\det : \mathbb{F}^{n,n} \rightarrow \mathbb{F}$  satisfies:

$$\begin{aligned} \begin{vmatrix} \mathbf{a}_1: \\ \vdots \\ \mathbf{a}_i: \\ \vdots \\ \mathbf{a}_j: \\ \vdots \\ \mathbf{a}_n: \end{vmatrix} &= - \begin{vmatrix} \mathbf{a}_1: \\ \vdots \\ \mathbf{a}_j: \\ \vdots \\ \mathbf{a}_i: \\ \vdots \\ \mathbf{a}_n: \end{vmatrix}; & \begin{vmatrix} \mathbf{a}_1: \\ \vdots \\ k\mathbf{a}_i: \\ \vdots \\ \mathbf{a}_n: \end{vmatrix} &= k \begin{vmatrix} \mathbf{a}_1: \\ \vdots \\ \mathbf{a}_i: \\ \vdots \\ \mathbf{a}_n: \end{vmatrix}; & \begin{vmatrix} \mathbf{a}_1: \\ \vdots \\ k\mathbf{a}_i: \\ \vdots \\ \mathbf{a}_i: \\ \vdots \\ \mathbf{a}_n: \end{vmatrix} &= 0; & \begin{vmatrix} \mathbf{a}_1: \\ \vdots \\ (\mathbf{a}_i + \mathbf{b}_i): \\ \vdots \\ \mathbf{a}_n: \end{vmatrix} &= \begin{vmatrix} \mathbf{a}_1: \\ \vdots \\ \mathbf{a}_i: \\ \vdots \\ \mathbf{a}_n: \end{vmatrix} + \begin{vmatrix} \mathbf{a}_1: \\ \vdots \\ \mathbf{b}_i: \\ \vdots \\ \mathbf{a}_n: \end{vmatrix}. \end{aligned}$$

## 3.1.4. Exercises.

3.1.1.- Find the determinant of matrices

$$A = \begin{bmatrix} 2 & 1 & 1 \\ 6 & 2 & 1 \\ -2 & 2 & 1 \end{bmatrix}, B = \begin{bmatrix} 2 & 1 & 1 \\ 6 & 0 & 1 \\ -2 & 0 & 1 \end{bmatrix}.$$

3.1.2.- Find the volume of the parallelepiped formed by the vectors

$$x_1 = \begin{bmatrix} 3 \\ 0 \\ -4 \end{bmatrix}, x_2 = \begin{bmatrix} 0 \\ 2 \\ 0 \end{bmatrix}, x_3 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}.$$

3.1.3.- Find the determinants of the upper and lower triangular matrices

$$U = \begin{bmatrix} 1 & 2 & 3 \\ 0 & 4 & 5 \\ 0 & 0 & 6 \end{bmatrix}, L = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 3 & 0 \\ 4 & 5 & 6 \end{bmatrix}.$$

3.1.4.- Find the determinant of the matrix

$$A = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}.$$

3.1.5.- Give an example to show that in general

$$\det(A + B) \neq \det(A) + \det(B).$$

3.1.6.- If  $A \in \mathbb{F}^{n,n}$  express

$$\det(2A), \quad \det(-A), \quad \det(A^2)$$

in terms of  $\det(A)$  and  $n$ .

3.1.7.- Given matrices  $A, P \in \mathbb{F}^{n,n}$ , with  $P$  invertible, let  $B = P^{-1}AP$ . Prove that  $\det(B) = \det(A)$ .

3.1.8.- Prove that for all matrix  $A \in \mathbb{F}^{n,n}$  holds that  $\det(A^*) = \overline{\det(A)}$ .

3.1.9.- Prove that for all  $A \in \mathbb{F}^{n,n}$  holds that  $\det(A^*A) \geq 0$ .

3.1.10.- Prove that for all  $A \in \mathbb{F}^{n,n}$  and all  $k \in \mathbb{F}$  holds that  $\det(kA) = k^n \det(A)$ .

3.1.11.- Prove that a skew-symmetric matrix  $A \in \mathbb{F}^{n,n}$ , with  $n$  odd, must satisfy that  $\det(A) = 0$ .

3.1.12.- Let  $A \in \mathbb{F}^{n,n}$  be a matrix satisfying that  $A^T A = I_n$ . Prove that

$$\det(A) = \pm 1.$$

## 3.2. APPLICATIONS

In the previous section we mentioned that the absolute value of the determinant is the area of a parallelogram formed by the column vectors of a  $2 \times 2$  matrix, while it is the volume of the parallelepiped formed by the column vectors of a  $3 \times 3$  matrix. There are several other applications of determinants. They determine whether a matrix  $A \in \mathbb{F}^{n,n}$  is invertible or not, since a matrix  $A$  is invertible iff  $\det(A) \neq 0$ . They determine whether a square system of linear equations has a unique solution for every source vector. The determinant is the key tool to find a formula for the inverse matrix, hence a formula for the solution components of a square linear system, called Cramer's rule. We discuss these to applications in more detail below.

**3.2.1. Inverse matrix formula.** The determinant of a matrix plays a crucial role in a formula for the inverse matrix. The existence of this formula does not change the fact that one can always compute the inverse matrix using Gauss operations. This formula for the inverse matrix is important though, since it explicitly shows how the inverse matrix depends on the coefficients of the original matrix. It thus shows how the inverse changes when the original matrix changes. The formula for the inverse matrix shows that the inverse matrix is a continuous function of the original matrix.

**Theorem 3.2.1.** *Given a matrix  $A = [A_{ij}] \in \mathbb{F}^{n,n}$ , let  $C = [C_{ij}]$  be its cofactor matrix, where the cofactors are given by  $C_{ij} = (-1)^{(i+j)} \det(\mathring{A}_{ij})$ , and matrix  $\mathring{A}_{ij} \in \mathbb{F}^{(n-1),(n-1)}$  is the  $(i, j)$ -minor of matrix  $A$ . If  $\det(A) \neq 0$ , then the inverse of matrix  $A$  is given by*

$$A^{-1} = \frac{1}{\det(A)} C^T, \quad \text{that is,} \quad (A^{-1})_{ij} = \frac{1}{\det(A)} C_{ji}.$$

**Proof of Theorem 3.2.1:** Since  $\det(A) \neq 0$  we know that matrix  $A$  is invertible. We only need to verify that the expression for  $A^{-1}$  given in Theorem 3.2.1 is correct. That is, we need to show that  $C^T A = \det(A) I_n = A C^T$ . We start with the product

$$(C^T A)_{ij} = (C_{1i} A_{1j} + \cdots + C_{ni} A_{nj}).$$

Notice that this component  $(C^T A)_{ij}$  is precisely the expansion along the column  $i$  of the determinant of the following matrix: The matrix constructed from matrix  $A$  by placing the column vector  $A_j$  in both columns  $i$  and  $j$ . In the particular case that  $i < j$ , this component  $(C^T A)_{ij}$  has the form

$$(C_{1i} A_{1j} + \cdots + C_{ni} A_{nj}) = \begin{vmatrix} A_{11} & \cdots & \mathbf{A}_{1j} & \cdots & A_{1j} & \cdots & A_{1n} \\ \vdots & & \vdots & & \vdots & & \vdots \\ A_{i1} & \cdots & \mathbf{A}_{ij} & \cdots & A_{ij} & \cdots & A_{in} \\ \vdots & & \vdots & & \vdots & & \vdots \\ A_{j1} & \cdots & \mathbf{A}_{jj} & \cdots & A_{jj} & \cdots & A_{jn} \\ \vdots & & \vdots & & \vdots & & \vdots \\ A_{n1} & \cdots & \mathbf{A}_{nj} & \cdots & A_{nj} & \cdots & A_{nn} \end{vmatrix}, \quad i < j,$$

where we have highlighted in blue the column  $i$  which is occupied by the column vector  $A_j$ . The determinant on the right hand side vanishes, since for  $i < j$  the matrix has the columns  $i$  and  $j$  repeated. A similar situation happens for  $i > j$ . So, we conclude that the non-diagonal elements of  $(C^T A)$  vanish. The diagonal elements are given by

$$(C^T A)_{ii} = (C_{1i} A_{1i} + \cdots + C_{ni} A_{ni}) = \det(A),$$

the determinant of  $A$  expanded along the  $i$ -th column. This shows that  $C^T A = \det(A) I_n$ . A similar analysis shows that  $A C^T = \det(A) I_n$ . This establishes the Theorem.  $\square$

**EXAMPLE 3.2.1:** Use the formula in Theorem 3.2.1 to find the inverse of matrix

$$A = \begin{bmatrix} 1 & 3 & -1 \\ 2 & 1 & 1 \\ 3 & 2 & 1 \end{bmatrix}.$$

**SOLUTION:** We already know from Example 3.1.3 that  $\det(A) = 1$ . We now need to compute all the cofactors:

$$\begin{aligned} C_{11} &= \begin{vmatrix} 2 & 1 \\ 3 & 1 \end{vmatrix} & C_{12} &= -\begin{vmatrix} 1 & 1 \\ 3 & 1 \end{vmatrix} & C_{13} &= \begin{vmatrix} 2 & 1 \\ 3 & 2 \end{vmatrix} \\ &= -1; & &= 1; & &= 1; \\ C_{21} &= -\begin{vmatrix} 3 & -1 \\ 2 & 1 \end{vmatrix} & C_{22} &= \begin{vmatrix} 1 & -1 \\ 3 & 1 \end{vmatrix} & C_{23} &= -\begin{vmatrix} 1 & 3 \\ 3 & 2 \end{vmatrix} \\ &= -5; & &= 4; & &= 7; \\ C_{31} &= \begin{vmatrix} 3 & -1 \\ 1 & 1 \end{vmatrix} & C_{32} &= -\begin{vmatrix} 1 & -1 \\ 2 & 1 \end{vmatrix} & C_{33} &= \begin{vmatrix} 1 & 3 \\ 2 & 1 \end{vmatrix} \\ &= 4; & &= -3; & &= -5. \end{aligned}$$

Therefore, the cofactor matrix is given by

$$C = \begin{bmatrix} -1 & 1 & 1 \\ -5 & 4 & 7 \\ 4 & -3 & -5 \end{bmatrix},$$

and the formula  $A^{-1} = C^T / \det(A)$  together with  $\det(A) = 1$  imply that

$$A^{-1} = \begin{bmatrix} -1 & -5 & 4 \\ 1 & 4 & -3 \\ 1 & 7 & -5 \end{bmatrix}.$$

◁

The formula in Theorem 3.2.1 is useful to compute individual components of the inverse of a matrix.

**EXAMPLE 3.2.2:** Find the coefficients  $(A^{-1})_{12}$  and  $(A^{-1})_{32}$  of the matrix

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & 2 \\ 1 & 1 & 3 \end{bmatrix}.$$

**SOLUTION:** We first need to compute

$$\det(A) = (1) \begin{vmatrix} -1 & 2 \\ 1 & 3 \end{vmatrix} - (1) \begin{vmatrix} 1 & 2 \\ 1 & 3 \end{vmatrix} + (1) \begin{vmatrix} 1 & -1 \\ 1 & 1 \end{vmatrix} = -5 - 1 + 2 \Rightarrow \det(A) = -4.$$

The formula in Theorem 3.2.1 implies that

$$(A^{-1})_{12} = \frac{1}{-4} C_{21}, \quad C_{21} = (-1) \begin{vmatrix} 1 & 1 \\ 1 & 3 \end{vmatrix} = -2, \quad \Rightarrow \quad (A^{-1})_{12} = \frac{1}{2}.$$

$$(A^{-1})_{32} = \frac{1}{-4} C_{23}, \quad C_{23} = (-1) \begin{vmatrix} 1 & 1 \\ 1 & 1 \end{vmatrix} = 0, \quad \Rightarrow \quad (A^{-1})_{32} = 0.$$

◁

**3.2.2. Cramer's rule.** We know that given an invertible matrix  $A \in \mathbb{F}^{n,n}$ , a system of linear equations  $Ax = b$  has a unique solution  $x$  for every source vector  $b \in \mathbb{F}^n$ . This solution can be written in terms of the inverse matrix  $A^{-1}$  as  $x = A^{-1}b$ . The formula for the inverse matrix given in Theorem 3.2.1 provides an explicit expression for the solution  $x$ . The result is known as Cramer's rule, and it is summarized below.

**Theorem 3.2.2.** *If the matrix  $A = [A_1, \dots, A_n] \in \mathbb{F}^{n,n}$  is invertible, then the system of linear equations  $Ax = b$  has a unique solution  $x = [x_i]$  for every  $b \in \mathbb{F}^n$  given by*

$$x_i = \frac{\det(A_i(b))}{\det(A)},$$

with matrix  $A_i(b) = [A_1, \dots, b, \dots, A_n]$  where vector  $b$  is placed in the  $i$ -th column.

**EXAMPLE 3.2.3:** Use Cramer's rule to find the solution  $x$  of the linear system  $Ax = b$ , where

$$A = \begin{bmatrix} 3 & -2 \\ -5 & 6 \end{bmatrix}, \quad b = \begin{bmatrix} 7 \\ -5 \end{bmatrix}.$$

**SOLUTION:** We first need to compute the determinant of  $A$ , that is,

$$\det(A) = 18 - 10 \Rightarrow \det(A) = 8.$$

Then, we need to find the matrices  $A_1(b)$  and  $A_2(b)$ , given by

$$A_1(b) = [b, A_2] = \begin{bmatrix} 7 & -2 \\ -5 & 6 \end{bmatrix}, \quad A_2(b) = [A_1, b] = \begin{bmatrix} 3 & 7 \\ -5 & -5 \end{bmatrix}.$$

We finally compute their determinants,

$$\det(A_1(b)) = 42 - 10 = 32, \quad \det(A_2(b)) = -15 + 35 = 20.$$

So the solution is,

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad x_1 = \frac{32}{8} = 4, \quad x_2 = \frac{20}{8} = \frac{5}{2}, \quad \Rightarrow \quad x = \begin{bmatrix} 4 \\ 5/2 \end{bmatrix}.$$

◁

**Proof of Theorem 3.2.2:** Since matrix  $A$  is invertible, the solution to the linear system is  $x = A^{-1}b$ . Using the formula for the inverse matrix given in Theorem 3.2.1, the solution  $x = [x_i]$  can be written as

$$x = \frac{1}{\det(A)} C^T b \Rightarrow x_i = \frac{1}{\det(A)} \sum_{j=1}^n C_{ji} b_j.$$

Notice that the sum in the last equation above is precisely the expansion along the column  $i$  of the determinant of the matrix  $A_j(b)$ , that is,

$$\det(A_i(b)) = \begin{vmatrix} A_{11} & \cdots & b_1 & \cdots & A_{1n} \\ \vdots & & \vdots & & \vdots \\ A_{n1} & \cdots & b_n & \cdots & A_{nn} \end{vmatrix} = b_1 C_{1i} + \cdots + b_n C_{ni} = \sum_{j=1}^n b_j C_{ji},$$

where the vector  $b$  replaces the vector  $A_i$  in the column  $i$  of matrix  $A$ . So, we conclude that

$$x_i = \frac{\det(A_i(b))}{\det(A)}.$$

This establishes the Theorem. □

**3.2.3. Determinants and Gauss operations.** Gauss elimination operations can be used to compute the determinant of a matrix. If a Gauss operation on matrix  $A$  produces matrix  $B$ , then  $\det(A)$  is related to  $\det(B)$  in a precise way. Here is the relation.

**Theorem 3.2.3.** *Let  $A, B \in \mathbb{F}^{m,n}$  be matrices related by a Gauss operation, that is,  $A \rightarrow B$  by a Gauss operation. Then, the following statements hold:*

(a) *If matrix  $B$  is the result of adding to a row in  $A$  a multiple of another row in  $A$ , then*

$$\det(A) = \det(B);$$

(b) *If matrix  $B$  is the result of interchanging two rows in  $A$ , then*

$$\det(A) = -\det(B);$$

(c) *If matrix  $B$  is the result of the multiplication of one row in  $A$  by a scalar  $\frac{1}{k} \in \mathbb{F}$ , then*

$$\det(A) = k \det(B).$$

**Proof of Theorem 3.2.3:** We use the row vector notation for matrices  $A$  and  $B$ ,

$$A = \begin{bmatrix} A_{1:} \\ \vdots \\ A_{n:} \end{bmatrix}, \quad B = \begin{bmatrix} B_{1:} \\ \vdots \\ B_{n:} \end{bmatrix}.$$

**Part (a):** Matrix  $B$  results from multiplying row  $j$  of  $A$  by  $k$  and adding that to row  $i$ ,

$$B = \begin{bmatrix} A_{1:} \\ \vdots \\ A_{i:} + kA_{j:} \\ \vdots \\ A_{j:} \\ \vdots \\ A_{n:} \end{bmatrix} \Rightarrow \det(B) = \begin{vmatrix} A_{1:} \\ \vdots \\ A_{i:} + kA_{j:} \\ \vdots \\ A_{j:} \\ \vdots \\ A_{n:} \end{vmatrix} = \begin{vmatrix} A_{1:} \\ \vdots \\ A_{i:} \\ \vdots \\ A_{j:} \\ \vdots \\ A_{n:} \end{vmatrix} + k \begin{vmatrix} A_{1:} \\ \vdots \\ \vdots \\ \vdots \\ A_{j:} \\ \vdots \\ A_{n:} \end{vmatrix} = \begin{vmatrix} A_{1:} \\ \vdots \\ A_{i:} \\ \vdots \\ A_{j:} \\ \vdots \\ A_{n:} \end{vmatrix} = \det(A).$$

**Part (b):** Matrix  $B$  is the result of the interchange of rows  $i$  and  $j$  in matrix  $A$ , that is,

$$A = \begin{bmatrix} A_{1:} \\ \vdots \\ A_{i:} \\ \vdots \\ A_{j:} \\ \vdots \\ A_{n:} \end{bmatrix} \rightarrow B = \begin{bmatrix} A_{1:} \\ \vdots \\ A_{j:} \\ \vdots \\ A_{i:} \\ \vdots \\ A_{n:} \end{bmatrix} \Rightarrow \det(B) = \begin{vmatrix} A_{1:} \\ \vdots \\ A_{j:} \\ \vdots \\ A_{i:} \\ \vdots \\ A_{n:} \end{vmatrix} = - \begin{vmatrix} A_{1:} \\ \vdots \\ A_{i:} \\ \vdots \\ A_{j:} \\ \vdots \\ A_{n:} \end{vmatrix} = -\det(A).$$

**Part (c):** Matrix  $B$  is the result of the multiplication of row  $i$  of  $A$  by  $\frac{1}{k} \in \mathbb{F}$ , that is,

$$A = \begin{bmatrix} A_{1:} \\ \vdots \\ A_{i:} \\ \vdots \\ A_{n:} \end{bmatrix} \rightarrow B = \begin{bmatrix} A_{1:} \\ \vdots \\ \frac{1}{k} A_{i:} \\ \vdots \\ A_{n:} \end{bmatrix} \Rightarrow \det(B) = \begin{vmatrix} A_{1:} \\ \vdots \\ \frac{1}{k} A_{i:} \\ \vdots \\ A_{n:} \end{vmatrix} = \frac{1}{k} \begin{vmatrix} A_{1:} \\ \vdots \\ A_{i:} \\ \vdots \\ A_{n:} \end{vmatrix} = \frac{1}{k} \det(A).$$

This establishes the Theorem.  $\square$

**EXAMPLE 3.2.4:** Use Gauss operations to transform matrix  $A$  below into upper-triangular form, and use that calculation to find  $\det(A)$ , where,

$$A = \begin{bmatrix} 1 & 3 & -1 \\ 2 & 1 & 1 \\ 3 & 2 & 1 \end{bmatrix}.$$

**SOLUTION:** We perform Gauss operations to transform  $A$  into upper-triangular form:

$$\det(A) = \begin{vmatrix} 1 & 3 & -1 \\ 2 & 1 & 1 \\ 3 & 2 & 1 \end{vmatrix} = \begin{vmatrix} 1 & 3 & -1 \\ 0 & -5 & 3 \\ 0 & -7 & 4 \end{vmatrix} = (-5) \begin{vmatrix} 1 & 3 & -1 \\ 0 & 1 & -3/5 \\ 0 & -7 & 4 \end{vmatrix} = (-5) \begin{vmatrix} 1 & 3 & -1 \\ 0 & 1 & -3/5 \\ 0 & 0 & -1/5 \end{vmatrix}.$$

We only need to reduce matrix  $A$  into upper-triangular form, not into reduced echelon form, since the determinant of an upper-triangular matrix is simple enough to find, just the product of its diagonal elements. In our case we find that

$$\det(A) = (-5) \begin{vmatrix} 1 & 3 & -1 \\ 0 & 1 & -3/5 \\ 0 & 0 & -1/5 \end{vmatrix} = (-5)(1)(1)\left(-\frac{1}{5}\right) \Rightarrow \det(A) = 1.$$

◁



## 3.2.4. Exercises.

**3.2.1.-** Use determinants to compute  $A^{-1}$ , where

$$A = \begin{bmatrix} 2 & 1 & 1 \\ 6 & 2 & 1 \\ -2 & 2 & 1 \end{bmatrix}.$$

**3.2.2.-** Find the coefficients  $(A^{-1})_{12}$  and  $(A^{-1})_{32}$ , where

$$A = \begin{bmatrix} 1 & 5 & 7 \\ 2 & 1 & 0 \\ 4 & 1 & 3 \end{bmatrix}.$$

**3.2.3.-** Use Gauss operations to reduce matrices  $A$  and  $B$  below to upper triangular form and evaluate their determinant, where

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 1 \\ 1 & 4 & 4 \end{bmatrix}, B = \begin{bmatrix} 1 & 3 & 5 \\ -1 & 4 & 2 \\ 3 & -2 & 4 \end{bmatrix}.$$

**3.2.4.-** Use Gauss operations to prove the formula

$$\begin{vmatrix} 1 & a & a^2 \\ 1 & b & b^2 \\ 1 & c & c^2 \end{vmatrix} = (b-a)(c-a)(c-b).$$

**3.2.5.-** Use the  $\det(A)$  to find the values of the constant  $k$  such that the system  $Ax = b$  has a unique solution for every source vector  $b$ , where

$$A = \begin{bmatrix} 1 & k & 0 \\ 0 & 1 & -1 \\ k & 0 & 1 \end{bmatrix}.$$

**3.2.6.-** Use Cramer's rule to find the solution to the linear system

$$ax_1 + bx_2 = 1$$

$$cx_1 + dx_2 = 0.$$

where  $a, b, c, d \in \mathbb{R}$ .

**3.2.7.-** Use Cramer's rule to find the solution to the linear system

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}.$$

## CHAPTER 4. VECTOR SPACES

## 4.1. SPACES AND SUBSPACES

A vector space is any set where the linear combination operation is defined on its elements. In a vector space, also called a linear space, the elements are not important. The actual elements that constitute the vector space are left unspecified, only the relation among them is determined. An example of a vector space is the set  $\mathbb{F}^n$  of  $n$ -vectors with the operation of linear combination studied in Chapter 1. Another example is the set  $\mathbb{F}^{m,n}$  of all  $m \times n$  matrices with the operation of linear combination studied in Chapter 2. We now define a vector space and comment its main properties. A subspace is introduced later on as a smaller vector space inside the original one. We end this Section with the concept of span of a set of vectors, which is a way to construct a subspace from any subset in a vector space.

**Definition 4.1.1.** A set  $V$  is a **vector space** over the scalar field  $\mathbb{F} \in \{\mathbb{R}, \mathbb{C}\}$  iff there are two operations defined on  $V$ , called **vector addition** and **scalar multiplication** with the following properties: For all  $\mathbf{u}, \mathbf{v}, \mathbf{w} \in V$  the vector addition satisfies

- (A1)  $\mathbf{u} + \mathbf{v} \in V$ , (closure of  $V$ );  
 (A2)  $\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$ , (commutativity);  
 (A3)  $(\mathbf{u} + \mathbf{v}) + \mathbf{w} = \mathbf{u} + (\mathbf{v} + \mathbf{w})$ , (associativity);  
 (A4)  $\exists \mathbf{0} \in V : \mathbf{0} + \mathbf{u} = \mathbf{u} \quad \forall \mathbf{u} \in V$ , (existence of a neutral element);  
 (A5)  $\forall \mathbf{u} \in V \quad \exists (-\mathbf{u}) \in V : \mathbf{u} + (-\mathbf{u}) = \mathbf{0}$ , (existence of an opposite element);

furthermore, for all  $a, b \in \mathbb{F}$  the scalar multiplication satisfies

- (M1)  $a\mathbf{u} \in V$ , (closure of  $V$ );  
 (M2)  $1\mathbf{u} = \mathbf{u}$ , (neutral element);  
 (M3)  $a(b\mathbf{u}) = (ab)\mathbf{u}$ , (associativity);  
 (M4)  $a(\mathbf{u} + \mathbf{v}) = a\mathbf{u} + a\mathbf{v}$ , (distributivity);  
 (M5)  $(a + b)\mathbf{u} = a\mathbf{u} + b\mathbf{u}$ , (distributivity).

The definition of a vector space does not specify the elements of the set  $V$ , it only determines the properties of the vector addition and scalar multiplication operations. We use the convention that elements in a vector space, called vectors, are represented in boldface. Nevertheless, we allow several exceptions, the first two of them are given in Examples 4.1.1 and 4.1.2. We now present several examples of vector spaces.

**EXAMPLE 4.1.1:** A vector space is the set  $\mathbb{F}^n$  of  $n$ -vectors  $\mathbf{u} = [u_i]$  with components  $u_i \in \mathbb{F}$  and the operations of column vector addition and scalar multiplication given by

$$[u_i] + [v_i] = [u_i + v_i], \quad a[u_i] = [au_i].$$

This space of column vectors was introduced in Chapter 1. Elements in these vector spaces are not represented in boldface, instead we keep the previous sanserif font,  $\mathbf{u} \in \mathbb{F}^n$ . The reason for this notation will be clear in Sect. 4.4.  $\triangleleft$

**EXAMPLE 4.1.2:** A vector space is the set  $\mathbb{F}^{m,n}$  of  $m \times n$  matrices  $\mathbf{A} = [A_{ij}]$  with matrix coefficients  $A_{ij} \in \mathbb{F}$  and the operations of addition and scalar multiplication given by

$$[A_{ij}] + [B_{ij}] = [A_{ij} + B_{ij}], \quad a[A_{ij}] = [aA_{ij}],$$

These operations were introduced in Chapter 2. As in the previous example, elements in these vector spaces are not represented in boldface, instead we keep the previous capital sanserif font,  $\mathbf{A} \in \mathbb{F}^{m,n}$ . The reason for this notation will be clear in Sect. 4.4.  $\triangleleft$

**EXAMPLE 4.1.3:** Let  $\mathbb{P}_n(U)$  be the set of all polynomials having degree  $n \geq 0$  and domain  $U \subset \mathbb{F}$ , that is,

$$\mathbb{P}_n(U) = \{p(x) = a_0 + a_1x + \cdots + a_nx^n, \text{ with } a_0, \dots, a_n \in \mathbb{F} \text{ and } x \in U \subset \mathbb{F}\}.$$

The set  $\mathbb{P}_n(U)$  together with the addition of polynomials  $(p+q)(x) = p(x) + q(x)$  and the scalar multiplication  $(ap)(x) = ap(x)$  is a vector space. In the case  $U = \mathbb{R}$  we use the notation  $\mathbb{P}_n = \mathbb{P}_n(\mathbb{R})$ .  $\triangleleft$

**EXAMPLE 4.1.4:** Let  $C^k([a, b], \mathbb{F})$  be the set of scalar valued functions with domain  $[a, b] \subset \mathbb{R}$  with the  $k$ -th derivative being a continuous function, that is,

$$C^k([a, b], \mathbb{F}) = \{f : [a, b] \rightarrow \mathbb{F} \text{ such that } f^{(k)} \text{ is continuous}\}.$$

The set  $C^k([a, b], \mathbb{F})$  together with the addition of functions  $(f+g)(x) = f(x) + g(x)$  and the scalar multiplication  $(af)(x) = af(x)$  is a vector space. The particular case  $C^k(\mathbb{R}, \mathbb{R})$  is denoted simply as  $C^k$ . The set of real valued continuous function is then  $C^0$ .  $\triangleleft$

**EXAMPLE 4.1.5:** Let  $\ell$  be the set of absolute convergent series, that is,

$$\ell = \left\{ a = \sum a_n : a_n \in \mathbb{F} \text{ and } \sum_{n=0}^{\infty} |a_n| \text{ exists} \right\}.$$

The set  $\ell$  with the addition of series  $a+b = \sum(a_n + b_n)$  and the scalar multiplication  $ca = \sum ca_n$  is a vector space.  $\triangleleft$

The properties (A1)-(M5) given in the definition of vector space are not redundant. As an example, these properties do not include the condition that the neutral element  $\mathbf{0}$  is unique, since it follows from the definition.

**Theorem 4.1.2.** *The element  $\mathbf{0}$  in a vector space is unique.*

**Proof Theorem 4.1.2:** Suppose that there exist two neutral elements  $\mathbf{0}_1$  and  $\mathbf{0}_2$  in the vector space  $V$ , that is,

$$\mathbf{0}_1 + \mathbf{u} = \mathbf{u} \quad \text{and} \quad \mathbf{0}_2 + \mathbf{u} = \mathbf{u} \quad \text{for all } \mathbf{u} \in V$$

Taking  $\mathbf{u} = \mathbf{0}_2$  in the first equation above, and  $\mathbf{u} = \mathbf{0}_1$  in the second equation above we obtain that

$$\mathbf{0}_1 + \mathbf{0}_2 = \mathbf{0}_2, \quad \mathbf{0}_2 + \mathbf{0}_1 = \mathbf{0}_1.$$

These equations above simply that the two neutral elements must be the same, since

$$\mathbf{0}_2 = \mathbf{0}_1 + \mathbf{0}_2 = \mathbf{0}_2 + \mathbf{0}_1 = \mathbf{0}_1;$$

where in the second equation we used that the addition operation is commutative. This establishes the Theorem.  $\square$

**Theorem 4.1.3.** *It holds that  $0\mathbf{u} = \mathbf{0}$  for all element  $\mathbf{u}$  in a vector space  $V$ .*

**Proof Theorem 4.1.3:** For every  $\mathbf{u} \in V$  holds

$$\mathbf{u} = 1\mathbf{u} = (1+0)\mathbf{u} = 1\mathbf{u} + 0\mathbf{u} = \mathbf{u} + 0\mathbf{u} = 0\mathbf{u} + \mathbf{u} \quad \Rightarrow \quad \mathbf{u} = 0\mathbf{u} + \mathbf{u}.$$

This last equation says that  $0\mathbf{u}$  is a neutral element,  $\mathbf{0}$ . Theorem 4.1.2 says that the neutral element is unique, so we conclude that, for all  $\mathbf{u} \in V$  holds that

$$0\mathbf{u} = \mathbf{0}.$$

This establishes the Theorem.  $\square$

Also notice that the property (A5) in the definition of vector space says that the opposite element exists, but it does not say whether it is unique. The opposite element is actually unique.

**Theorem 4.1.4.** *The opposite element  $-\mathbf{u}$  in a vector space is unique.*

**Proof Theorem 4.1.4:** Suppose there are two opposite elements  $-\mathbf{u}_1$  and  $-\mathbf{u}_2$  to the element  $\mathbf{u} \in V$ , that is,

$$\mathbf{u} + (-\mathbf{u}_1) = \mathbf{0}, \quad \mathbf{u} + (-\mathbf{u}_2) = \mathbf{0}.$$

Therefore,

$$\begin{aligned} (-\mathbf{u}_1) &= \mathbf{0} + (-\mathbf{u}_1) \\ &= \mathbf{u} + (-\mathbf{u}_2) + (-\mathbf{u}_1) \\ &= (-\mathbf{u}_2) + \mathbf{u} + (-\mathbf{u}_1) \\ &= (-\mathbf{u}_2) + \mathbf{0} \\ &= \mathbf{0} + (-\mathbf{u}_2) \\ &= (-\mathbf{u}_2) \Rightarrow (-\mathbf{u}_1) = (-\mathbf{u}_2). \end{aligned}$$

This establishes the Theorem. □

Finally, notice that the element  $(-\mathbf{u})$  opposite to  $\mathbf{u}$  is actually the element  $(-1)\mathbf{u}$ .

**Theorem 4.1.5.** *It holds that  $(-1)\mathbf{u} = (-\mathbf{u})$ .*

**Proof Theorem 4.1.5:**

$$\mathbf{0} = 0\mathbf{u} = (1 - 1)\mathbf{u} = 1\mathbf{u} + (-1)\mathbf{u} = \mathbf{u} + (-1)\mathbf{u}.$$

Hence  $(-1)\mathbf{u}$  is an opposite element of  $\mathbf{u}$ . Since Theorem 4.1.4 says that the opposite element is unique, we conclude that  $(-1)\mathbf{u} = (-\mathbf{u})$ . This establishes the Theorem. □

**4.1.1. Subspaces.** We now introduce the notion of a subspace of a vector space, which is essentially a smaller vector space inside the original vector space. Subspaces are important in physics, since physical processes frequently take place not inside the whole space but in a particular subspace. For instance, planetary motion does not take place in the whole space but it is confined to a plane.

**Definition 4.1.6.** *The subset  $W \subset V$  of a vector space  $V$  over  $\mathbb{F}$  is called a **subspace** of  $V$  iff for all  $\mathbf{u}, \mathbf{v} \in W$  and all  $a, b \in \mathbb{F}$  holds that  $a\mathbf{u} + b\mathbf{v} \in W$ .*

A subspace is a particular type of set in a vector space. Is a set where all possible linear combinations of two elements in the set results in another element in the same set. In other words, elements outside the set cannot be reached by linear combinations of elements within the set. For this reason a subspace is called a *closed set under linear combinations*. The following statement is usually helpful

**Theorem 4.1.7.** *If  $W \subset V$  is a subspace of a vector space  $V$ , then  $\mathbf{0} \in W$ .*

This statement says that  $\mathbf{0} \notin W$  implies that  $W$  is not a subspace. However, if actually  $\mathbf{0} \in W$ , this fact alone does not prove that  $W$  is a subspace. One must show that  $W$  is closed under linear combinations.

**Proof of Theorem 4.1.7:** Since  $W$  is closed under linear combinations, given any element  $\mathbf{u} \in W$ , the trivial linear combination  $0\mathbf{u} = \mathbf{0}$  must belong to  $W$ , hence  $\mathbf{0} \in W$ . This establishes the Theorem. □

**EXAMPLE 4.1.6:** Show that the set  $W \subset \mathbb{R}^3$  given by  $W = \{\mathbf{u} = [u_i] \in \mathbb{R}^3 : u_3 = 0\}$  is a subspace of  $\mathbb{R}^3$ :

**SOLUTION:** Given two arbitrary elements  $u, v \in W$  we must show that  $au + bv \in W$  for all  $a, b \in \mathbb{R}$ . Since  $u, v \in W$  we know that

$$u = \begin{bmatrix} u_1 \\ u_2 \\ 0 \end{bmatrix}, \quad v = \begin{bmatrix} v_1 \\ v_2 \\ 0 \end{bmatrix}.$$

Therefore

$$au + bv = \begin{bmatrix} au_1 + bv_1 \\ au_2 + bv_2 \\ 0 \end{bmatrix} \in W,$$

since the third component vanishes, which makes the linear combination an element in  $W$ . Hence,  $W$  is a subspace of  $\mathbb{R}^3$ . In Fig. 29 we see the plane  $u_3 = 0$ . It is a subspace, since not only  $0 \in W$ , but any linear combination of vectors on the plane stays on the plane.  $\triangleleft$

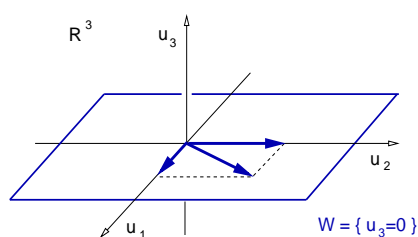


FIGURE 29. The horizontal plane  $u_3 = 0$  is a subspace of  $\mathbb{R}^3$ .

**EXAMPLE 4.1.7:** Show that the set  $W = \{u = [u_i] \in \mathbb{R}^2 : u_2 = 1\}$  is not a subspace of  $\mathbb{R}^2$ .

**SOLUTION:** The set  $W$  is not a subspace, since  $0 \notin W$ . This is enough to show that  $W$  is not a subspace. Another proof is that the addition of two vectors in the set is a vector outside the set, as can be seen by the following calculation,

$$u = \begin{bmatrix} u_1 \\ 1 \end{bmatrix} \in W, \quad v = \begin{bmatrix} v_1 \\ 1 \end{bmatrix} \in W \quad \Rightarrow \quad u + v = \begin{bmatrix} u_1 + v_1 \\ 2 \end{bmatrix} \notin W.$$

The second component in the addition above is 2, not 1, hence this addition does not belong to  $W$ . An example of this calculation is given in Fig. 30.  $\triangleleft$

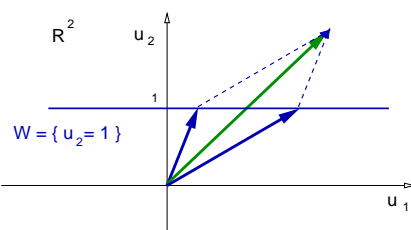


FIGURE 30. The horizontal line  $u_2 = 1$  is not a subspace of  $\mathbb{R}^2$ .

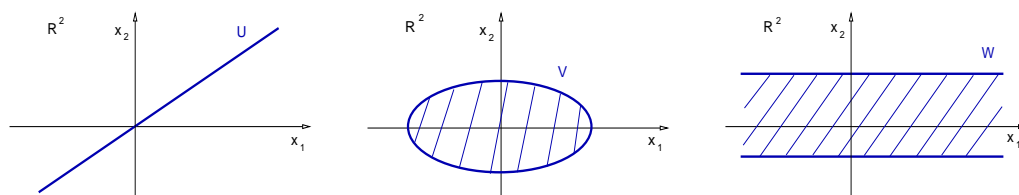


FIGURE 31. Three subsets,  $U$ ,  $V$ , and  $W$ , of  $\mathbb{R}^2$ . Only the set  $U$  is a subspace.

**EXAMPLE 4.1.8:** Determine which one of the sets given in Fig. 31 is a subspace of  $\mathbb{R}^2$ .

**SOLUTION:** The set  $U$  is a vector space, since any linear combination of vectors parallel to the line is again a vector parallel to the line. The sets  $V$  and  $W$  are not subspaces, since given a vector  $\mathbf{u}$  in these spaces, a the vector  $a\mathbf{u}$  does not belong to these sets for a number  $a \in \mathbb{R}$  big enough. This argument is sketched in Fig. 32.

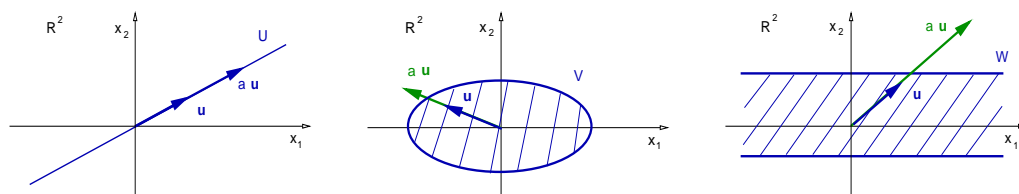


FIGURE 32. Three subsets,  $U$ ,  $V$ , and  $W$ , of  $\mathbb{R}^2$ . Only the set  $U$  is a subspace.

◀

**4.1.2. The span of finite sets.** If a set is not a subspace there is a way to increase it into a subspace. Define a new set including all possible linear combinations of elements in the old set.

**Definition 4.1.8.** The *span* of a finite set  $S = \{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  in a vector space  $V$  over  $\mathbb{F}$ , denoted as  $\text{Span}(S)$ , is the set given by

$$\text{Span}(S) = \{\mathbf{u} \in V : \mathbf{u} = c_1 \mathbf{u}_1 + \dots + c_n \mathbf{u}_n, \quad c_1, \dots, c_n \in \mathbb{F}\}.$$

The following result remarks that the span of a set is a subspace.

**Theorem 4.1.9.** Given a finite set  $S$  in a vector space  $V$ ,  $\text{Span}(S)$  is a subspace of  $V$ .

**Proof of Theorem 4.1.9:** Since  $\text{Span}(S)$  contains all possible linear combinations of the elements in  $S$ , then  $\text{Span}(S)$  is closed under linear combinations. This establishes the Theorem.  $\square$

**EXAMPLE 4.1.9:** The subspace  $\text{Span}(\{\mathbf{v}\})$ , that is, the set of all possible linear combinations of the vector  $\mathbf{v}$ , is formed by all vectors of the form  $c\mathbf{v}$ , and these vectors belong to a line containing  $\mathbf{v}$ . The subspace  $\text{Span}(\{\mathbf{v}, \mathbf{w}\})$ , that is, the set of all linear combinations of two vectors  $\mathbf{v}$ ,  $\mathbf{w}$ , is the plane containing both vectors  $\mathbf{v}$  and  $\mathbf{w}$ . See Fig. 33 for the case of the vector spaces  $\mathbb{R}^2$  and  $\mathbb{R}^3$ , respectively.  $\triangleleft$

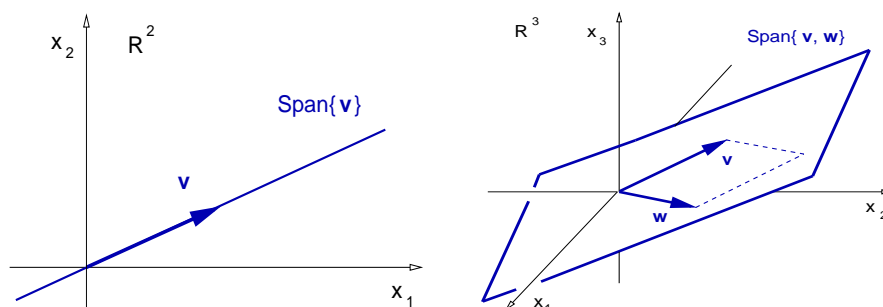


FIGURE 33. Examples of the span of a set of a single vector, and the span of a linearly independent set of two vectors.

**4.1.3. Algebra of subspaces.** We now show that the intersection of two subspaces is again a subspace. However, the union of two subspaces is not, in general, a subspace. The smaller subspace containing the union of two subspaces is precisely the span of the union. We then define the addition of two subspaces as the span of the union of two subspaces.

**Theorem 4.1.10.** *If  $W_1$  and  $W_2$  are subspaces of a vector space  $V$ , then  $W_1 \cap W_2 \subset V$  is also a subspace of  $V$ .*

**Proof of Theorem 4.1.10:** Let  $u$  and  $v$  be any two elements in  $W_1 \cap W_2$ . This means that  $u, v \in W_1$ , which is a subspace, so any linear combination  $(au + bv) \in W_1$ . Since  $u, v$  belong to  $W_1 \cap W_2$  they also belong to  $W_2$ , which is a subspace, so any linear combination  $(au + bv) \in W_2$ . Therefore, any linear combination  $(au + bv) \in W_1 \cap W_2$ . This establishes the Theorem.  $\square$

**EXAMPLE 4.1.10:** The sketch in Fig. 34 shows the intersection of two subspaces in  $\mathbb{R}^3$ , a plane and a line. In this case the intersection is the former line, so the intersection is a subspace.

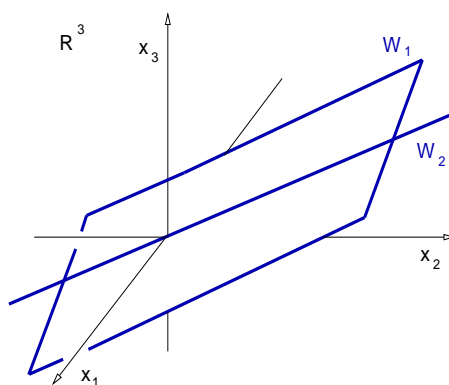


FIGURE 34. Intersection of two subspaces,  $W_1$  and  $W_2$  in  $\mathbb{R}^3$ . Since the line  $W_2$  is included into the plane  $W_1$ , we have that  $W_1 \cap W_2 = W_2$ .

$\triangleleft$

While the intersection of two subspaces is always a subspace, their union is, in general, not a subspace, unless one subspace is contained into the other.

**EXAMPLE 4.1.11:** Consider the vector space  $V = \mathbb{R}^2$ , and the subspaces  $W_1$  and  $W_2$  given by the lines sketched in Fig. 35. Their union is the set formed by these two lines. This set is not a subspace, since the addition of the vectors  $\mathbf{u}_1 \in W_1$  with  $\mathbf{u}_2 \in W_2$  does not belong to  $W_1 \cup W_2$ , as is it shown in Fig. 35.

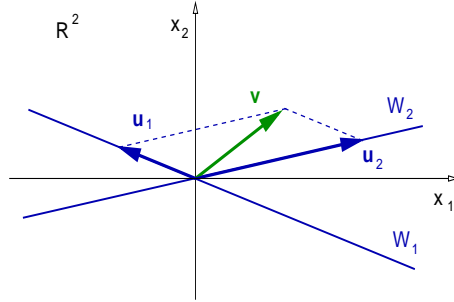


FIGURE 35. The union of the subspaces  $W_1$  and  $W_2$  is the set formed by these two lines. This is not a subspace, since the addition of  $\mathbf{u}_1 \in W_1$  and  $\mathbf{u}_2 \in W_2$  is the vector  $\mathbf{v}$  which does not belong to  $W_1 \cup W_2$ .

◁

Although the union of two subspaces is not always a subspace, it is possible to enlarge the union into a subspace. The idea is to incorporate all possible additions of vectors in the two original subspaces, and the result is called the addition of the two subspaces. Here is the precise definition.

**Definition 4.1.11.** The *addition of the subspaces*  $W_1, W_2$  in a vector space  $V$ , denoted as  $W_1 + W_2$ , is the set given by

$$W_1 + W_2 = \{ \mathbf{u} \in V : \mathbf{u} = \mathbf{w}_1 + \mathbf{w}_2 \text{ with } \mathbf{w}_1 \in W_1, \mathbf{w}_2 \in W_2 \}.$$

The following result remarks that the addition of subspaces is again a subspace.

**Theorem 4.1.12.** If  $W_1$  and  $W_2$  are subspaces of a vector space  $V$ , then the addition  $W_1 + W_2$  is also a subspace of  $V$ .

**Proof of Theorem 4.1.12:** Suppose that  $\mathbf{x} \in W_1 + W_2$  and  $\mathbf{y} \in W_1 + W_2$ . We must show that any linear combination  $a\mathbf{x} + b\mathbf{y}$  also belongs to  $W_1 + W_2$ . This is the case, by the following argument. Since  $\mathbf{x} \in W_1 + W_2$ , there exist  $\mathbf{x}_1 \in W_1$  and  $\mathbf{x}_2 \in W_2$  such that  $\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_2$ . Analogously, since  $\mathbf{y} \in W_1 + W_2$ , there exist  $\mathbf{y}_1 \in W_1$  and  $\mathbf{y}_2 \in W_2$  such that  $\mathbf{y} = \mathbf{y}_1 + \mathbf{y}_2$ . Now any linear combination of  $\mathbf{x}$  and  $\mathbf{y}$  satisfies

$$\begin{aligned} a\mathbf{x} + b\mathbf{y} &= a(\mathbf{x}_1 + \mathbf{x}_2) + b(\mathbf{y}_1 + \mathbf{y}_2) \\ &= (a\mathbf{x}_1 + b\mathbf{y}_1) + (a\mathbf{x}_2 + b\mathbf{y}_2) \end{aligned}$$

Since  $W_1$  and  $W_2$  are subspaces,  $(a\mathbf{x}_1 + b\mathbf{y}_1) \in W_1$ , and  $(a\mathbf{x}_2 + b\mathbf{y}_2) \in W_2$ . Therefore, the equation above says that  $(a\mathbf{x} + b\mathbf{y}) \in W_1 + W_2$ . This establishes the Theorem.  $\square$

**EXAMPLE 4.1.12:** The sketch in Fig. 36 shows the union and the addition of two subspaces in  $\mathbb{R}^3$ , each subspace given by a line through the origin. While the union is not a subspace, their addition is the plane containing both lines, which is a subspace. Given any non-zero vector  $\mathbf{w}_1 \in W_1$  and any other non-zero vector  $\mathbf{w}_2 \in W_2$ , one can verify that the sum of two subspaces is the span of  $\{\mathbf{w}_1, \mathbf{w}_2\}$ , that is,

$$W_1 + W_2 = \text{Span}(\{\mathbf{w}_1\} \cup \{\mathbf{w}_2\}).$$



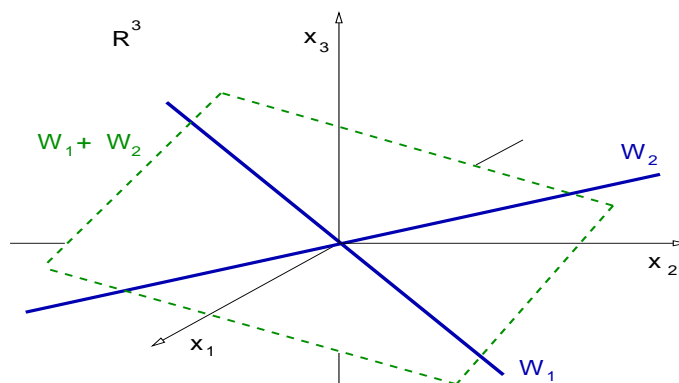


FIGURE 36. Union and addition of the subspaces  $W_1$  and  $W_2$  in  $\mathbb{R}^3$ . The union is not a subspace, while the addition is a subspace of  $\mathbb{R}^3$ .

**4.1.4. Internal direct sums.** This is a particular case of the addition of subspaces. It is called internal direct sum in order to differentiate it from another type of direct sum found in the literature. The latter, also called external direct sum, is a sum of different vector spaces, and it is a way to construct new vector spaces from old ones. We do not discuss this type of direct sums here. From now on, direct sum in these notes means the internal direct sum of subspaces inside a vector space.

**Definition 4.1.13.** Given a vector space  $V$ , we say that  $V$  is the *internal direct sum* of the subspaces  $W_1, W_2 \subset V$ , denoted as  $V = W_1 \oplus W_2$ , iff every vector  $v \in V$  can be written in a *unique way*, except for order, as a sum of vectors from  $W_1$  and  $W_2$ .

A crucial part in the definition above is the uniqueness of the decomposition of every vector  $v \in V$  as a sum of a vector in  $W_1$  plus a vector in  $W_2$ . By uniqueness we mean the following: For every  $v \in V$  exist  $w_1 \in W_1$  and  $w_2 \in W_2$  such that  $v = w_1 + w_2$ , and if  $v = \tilde{w}_1 + \tilde{w}_2$  with  $\tilde{w}_1 \in W_1$  and  $\tilde{w}_2 \in W_2$ , then  $w_1 = \tilde{w}_1$  and  $w_2 = \tilde{w}_2$ . In the case that  $V = W_1 \oplus W_2$  we say that  $W_1$  and  $W_2$  are *direct summands* of  $V$ , and we also say that  $W_1$  is the *direct complement* of  $W_2$  in  $V$ . There is an useful characterization of internal direct sums.

**Theorem 4.1.14.** A vector space  $V$  is the direct sum of subspaces  $W_1$  and  $W_2$  iff holds both  $V = W_1 + W_2$  and  $W_1 \cap W_2 = \{0\}$ .

**Proof of Theorem 4.1.14:**

( $\Rightarrow$ ) If  $V = W_1 \oplus W_2$ , then it implies that  $V = W_1 + W_2$ . Suppose that  $v \in W_1 \cap W_2$ , then on the one hand, there exists  $w_1 \in W_1$  such that  $v = w_1 + 0$ ; on the other hand, there is  $w_2 \in W_2$  such that  $v = 0 + w_2$ . Therefore,  $w_1 = 0$  and  $w_2 = 0$ , so  $W_1 \cap W_2 = \{0\}$ .

( $\Leftarrow$ ) Since  $V = W_1 + W_2$ , for every  $v \in V$  there exist  $w_1 \in W_1$  and  $w_2 \in W_2$  such that  $v = w_1 + w_2$ . Suppose there exists other vectors  $\tilde{w}_1 \in W_1$  and  $\tilde{w}_2 \in W_2$  such that  $v = \tilde{w}_1 + \tilde{w}_2$ . Then,

$$0 = (w_1 - \tilde{w}_1) + (w_2 - \tilde{w}_2) \Leftrightarrow (w_1 - \tilde{w}_1) = -(w_2 - \tilde{w}_2),$$

Therefore  $(w_1 - \tilde{w}_1) \in W_2$  and so  $(w_1 - \tilde{w}_1) \in W_1 \cap W_2$ . Since  $W_1 \cap W_2 = \{0\}$ , we then conclude that  $w_1 = \tilde{w}_1$ , which also says  $w_2 = \tilde{w}_2$ . Then  $V = W_1 \oplus W_2$ . This establishes the Theorem.  $\square$

**EXAMPLE 4.1.13:** Denote by  $\text{Sym}$  and  $\text{SkewSym}$  the sets of all symmetric and all skew-symmetric  $n \times n$  matrices. Show that  $\mathbb{F}^{n,n} = \text{Sym} \oplus \text{SkewSym}$ .

**SOLUTION:** Given any matrix  $A \in \mathbb{F}^{n,n}$ , holds

$$A = A + \frac{1}{2}(A^T - A^T) = \frac{1}{2}(A + A^T) + \frac{1}{2}(A - A^T).$$

We then can decompose matrix as  $A = B + C$ , where matrix  $B = (A + A^T)/2 \in \text{Sym}$  while matrix  $C = (A - A^T)/2 \in \text{SkewSym}$ . That is, we can write any square matrix as a symmetric matrix plus a skew-symmetric matrix, hence  $\mathbb{F}^{n,n} \subset \text{Sym} + \text{SkewSym}$ . The other inclusion is obvious, that is,  $\text{Sym} + \text{SkewSym} \subset \mathbb{F}^{n,n}$ , because each term in the sum is a subset of  $\mathbb{F}^{n,n}$ . So, we conclude that

$$\mathbb{F}^{n,n} = \text{Sym} + \text{SkewSym}.$$

Now we must show that  $\text{Sym} \cap \text{SkewSym} = \{\mathbf{0}\}$ . This is the case, as the following argument shows. If matrix  $D \in \text{Sym} \cap \text{SkewSym}$ , then matrix  $D$  is symmetric,  $D = D^T$ , but matrix  $D$  is also skew-symmetric,  $D = -D^T$ . This implies that  $D = -D$ , that is,  $D = \mathbf{0}$ , proving our assertion that  $\text{Sym} \cap \text{SkewSym} = \{\mathbf{0}\}$ . We then conclude that

$$\mathbb{F}^{n,n} = \text{Sym} \oplus \text{SkewSym}.$$

◁

## 4.1.5. Exercises.

4.1.1.- Determine which of the following subsets of  $\mathbb{R}^n$ , with  $n \geq 1$ , are in fact subspaces. Justify your answers.

- (a)  $\{x \in \mathbb{R}^n : x_i \geq 0 \quad i = 1, \dots, n\}$ ;
- (b)  $\{x \in \mathbb{R}^n : x_1 = 0\}$ ;
- (c)  $\{x \in \mathbb{R}^n : x_1 x_2 = 0 \quad n \geq 2\}$ ;
- (d)  $\{x \in \mathbb{R}^n : x_1 + \dots + x_n = 0\}$ ;
- (e)  $\{x \in \mathbb{R}^n : x_1 + \dots + x_n = 1\}$ ;
- (f)  $\{x \in \mathbb{R}^n : Ax = b, A \neq 0, b \neq 0\}$ .

4.1.2.- Determine which of the following subsets of  $\mathbb{F}^{n,n}$ , with  $n \geq 1$ , are in fact subspaces. Justify your answers.

- (a)  $\{A \in \mathbb{F}^{n,n} : A = A^T\}$ ;
- (b)  $\{A \in \mathbb{F}^{n,n} : A \text{ invertible}\}$ ;
- (c)  $\{A \in \mathbb{F}^{n,n} : A \text{ not invertible}\}$ ;
- (d)  $\{A \in \mathbb{F}^{n,n} : A \text{ upper-triangular}\}$ ;
- (e)  $\{A \in \mathbb{F}^{n,n} : A^2 = A\}$ ;
- (f)  $\{A \in \mathbb{F}^{n,n} : \text{tr}(A) = 0\}$ .
- (g) Given a matrix  $X \in \mathbb{F}^{n,n}$ , define  $\{A \in \mathbb{F}^{n,n} : [A, X] = 0\}$ .

4.1.3.- Find  $W_1 + W_2 \subset \mathbb{R}^3$ , where  $W_1$  is a plane passing through the origin in  $\mathbb{R}^3$  and  $W_2$  is a line passing through the origin in  $\mathbb{R}^3$  not contained in  $W_1$ .

4.1.4.- Sketch a picture of the subspaces spanned by the following vectors:

- (a)  $\left\{ \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 3 \\ 6 \\ 9 \end{bmatrix}, \begin{bmatrix} -2 \\ -4 \\ -6 \end{bmatrix} \right\}$ ;
- (b)  $\left\{ \begin{bmatrix} 0 \\ 2 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 2 \\ 3 \\ 0 \end{bmatrix} \right\}$ ;
- (c)  $\left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right\}$ .

4.1.5.- Given two finite subsets  $S_1, S_2$  in a vector space  $V$ , show that

$$\text{Span}(S_1 \cup S_2) = \text{Span}(S_1) + \text{Span}(S_2).$$

4.1.6.- Let  $W_1 \subset \mathbb{R}^3$  be the subspace

$$W_1 = \text{Span}\left(\left\{ \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \right\}\right).$$

Find a subspace  $W_2 \subset \mathbb{R}^3$  such that  $\mathbb{R}^3 = W_1 \oplus W_2$ .

## 4.2. LINEAR DEPENDENCE

**4.2.1. Linearly dependent sets.** In this Section we present the notion of a linearly dependent set of vectors. If one of the vectors in the set is a linear combination of the other vectors in the set, then the set is called linearly dependent. If this is not the case, the set is called linearly independent. This notion plays a crucial role in Sect. 4.3 to define a basis of a vector space. Bases are very useful in part because every vector in the vector space can be decomposed in a unique way as a linear combination of the basis elements. Bases also provide a precise way to measure the size of a vector space.

**Definition 4.2.1.** A *finite* set of vectors  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  in a vector space is called **linearly dependent** iff there exists a set of scalars  $\{c_1, \dots, c_k\}$ , not all of them zero, such that,

$$c_1 \mathbf{v}_1 + \dots + c_k \mathbf{v}_k = \mathbf{0}. \quad (4.1)$$

On the other hand, the set  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  is called **linearly independent** iff Eq. (4.1) implies that every scalar vanishes,  $c_1 = \dots = c_k = 0$ .

The wording in this definition is carefully chosen to cover the case of the empty set. The result is that the empty set is linearly independent. It might seem strange, but this result fits well with the rest of the theory. On the other hand, the set  $\{\mathbf{0}\}$  is linearly dependent, since  $c_1 \mathbf{0} = \mathbf{0}$  for any  $c_1 \neq 0$ . Moreover, any set containing the zero vector is also linearly dependent.

Linear dependence or independence are properties of a set of vectors. There is no meaning to a vector to be linearly dependent, or independent. And there is no meaning of a set of linearly dependent vectors, as well as a set of linearly independent vectors. What is meaningful is to talk of a linearly dependent or independent set of vectors.

**EXAMPLE 4.2.1:** Show that the set  $S \subset \mathbb{R}^2$  below is linearly dependent,

$$S = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 2 \\ 3 \end{bmatrix} \right\}.$$

**SOLUTION:** It is clear that

$$\begin{bmatrix} 2 \\ 3 \end{bmatrix} = 2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 3 \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad \Rightarrow \quad 2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 3 \begin{bmatrix} 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

Since  $c_1 = 2$ ,  $c_2 = 3$ , and  $c_3 = -1$  are non-zero, the set  $S$  is linearly dependent.  $\triangleleft$

It will be convenient to have the concept of a linearly dependent or independent set containing infinitely many vectors.

**Definition 4.2.2.** An *infinite* set of vectors  $S = \{\mathbf{v}_1, \mathbf{v}_2, \dots\}$  in a vector space  $V$  is called **linearly independent** iff every finite subset of  $S$  is linearly independent. Otherwise, the infinite set  $S$  is called **linearly dependent**.

**EXAMPLE 4.2.2:** Consider the vector space  $V = C^\infty([-\ell, \ell], \mathbb{R})$ , that is, the space of infinitely differentiable real-valued functions defined on the domain  $[-\ell, \ell] \subset \mathbb{R}$  with the usual operation of linear combination of functions. This vector space contains linearly independent sets with infinitely many vectors. One example is the infinite sets  $S_1$  below, which is linearly independent,

$$S_1 = \{1, x, x^2, \dots, x^n, \dots\}.$$

Another example is the infinite set  $S_2$ , which is also linearly independent,

$$S_2 = \left\{ 1, \cos\left(\frac{n\pi x}{\ell}\right), \sin\left(\frac{n\pi x}{\ell}\right) \right\}_{n=1}^{\infty}.$$

$\triangleleft$

**4.2.2. Main properties.** As we have seen in the Example 4.2.1 above, in a linearly dependent set there is always at least one vector that is a linear combination of the other vectors in the set. This is simple to see from the Definition 4.2.1. Since not all the coefficients  $c_i$  are zero in a linearly dependent set, let us suppose that  $c_j \neq 0$ ; then from the Eq. (4.1) we obtain

$$\mathbf{v}_j = -\frac{1}{c_j} [c_1 \mathbf{v}_1 + \cdots + c_{j-1} \mathbf{v}_{j-1} + c_{j+1} \mathbf{v}_{j+1} + \cdots + c_k \mathbf{v}_k],$$

that is,  $\mathbf{v}_j$  is a linear combination of the other vectors in the set.

**Theorem 4.2.3.** *The set  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  is linearly dependent with the vector  $\mathbf{v}_k$  being a linear combination of the remaining  $k - 1$  vectors iff*

$$\text{Span}(\{\mathbf{v}_1, \dots, \mathbf{v}_k\}) = \text{Span}(\{\mathbf{v}_1, \dots, \mathbf{v}_{k-1}\}).$$

This Theorem captures the idea behind the notion of a linearly dependent set: A finite set is linearly dependent iff there exists a smaller set with the same span. In this sense the vector  $\mathbf{v}_k$  in the Proposition above is redundant with respect to linear combinations.

**Proof of Theorem 4.2.3:** Let  $S_k = \{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  and  $S_{k-1} = \{\mathbf{v}_1, \dots, \mathbf{v}_{k-1}\}$ .

On the one hand, if  $\mathbf{v}_k$  is a linear combination of the other vectors in  $S$ , then for every  $\mathbf{x} \in \text{Span}(S_k)$  can be expressed as an element in  $\text{Span}(S_{k-1})$  simply by replacing  $\mathbf{v}_k$  in terms of the vectors in  $\tilde{S}$ . This shows that  $\text{Span}(S_k) \subset \text{Span}(S_{k-1})$ . The other inclusion is trivial, so  $\text{Span}(S_k) = \text{Span}(S_{k-1})$ .

On the other hand, if  $\text{Span}(S_k) = \text{Span}(S_{k-1})$ , this means that  $\mathbf{v}_k$  is a linear combination of the elements in  $S_{k-1}$ . Therefore, the set  $S_k$  is linearly dependent. This establishes the Theorem.  $\square$

**EXAMPLE 4.2.3:** Consider the set  $S \subset \mathbb{R}^3$  given by  $S = \left\{ \begin{bmatrix} -2 \\ 2 \\ -3 \end{bmatrix}, \begin{bmatrix} 4 \\ -6 \\ 8 \end{bmatrix}, \begin{bmatrix} -2 \\ -3 \\ 2 \end{bmatrix}, \begin{bmatrix} -4 \\ 1 \\ -3 \end{bmatrix} \right\}$ .

Find a set  $\tilde{S} \subset S$  having the smallest number of vectors such that  $\text{Span}(\tilde{S}) = \text{Span}(S)$ .

**SOLUTION:** We have to find all the redundant vectors in  $S$  with respect to linear combinations. In other words, we have to find a linearly independent subset of  $\tilde{S} \subset S$  such that  $\text{Span}(\tilde{S}) = \text{Span}(S)$ . The calculation we must do is to find the non-zero coefficients  $c_i$  in the solution of the equation

$$\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = c_1 \begin{bmatrix} -2 \\ 2 \\ -3 \end{bmatrix} + c_2 \begin{bmatrix} 4 \\ -6 \\ 8 \end{bmatrix} + c_3 \begin{bmatrix} -2 \\ -3 \\ 2 \end{bmatrix} + c_4 \begin{bmatrix} -4 \\ 1 \\ -3 \end{bmatrix} = \begin{bmatrix} -2 & 4 & -2 & -4 \\ 2 & -6 & -3 & 1 \\ -3 & 8 & 2 & -3 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ c_4 \end{bmatrix}.$$

Hence, we must find the reduced echelon form of the coefficient matrix above, that is,

$$A = \begin{bmatrix} -2 & 4 & -2 & -4 \\ 2 & -6 & -3 & 1 \\ -3 & 8 & 2 & -3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -2 & 1 & 2 \\ 0 & -2 & -5 & -3 \\ 0 & 2 & 5 & 3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -2 & 1 & 2 \\ 0 & 2 & 5 & 3 \\ 0 & 0 & 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 6 & 5 \\ 0 & 1 & \frac{5}{2} & \frac{3}{2} \\ 0 & 0 & 0 & 0 \end{bmatrix} = \mathbf{E}_A.$$

This means that the solution for the coefficients is

$$c_1 = -6c_3 - 5c_4, \quad c_2 = -\frac{5}{2}c_3 - \frac{3}{2}c_4, \quad c_3, c_4 \text{ free variables.}$$

Choosing:

$$c_4 = 0, \quad c_3 = 2 \quad \Rightarrow \quad c_1 = -12, \quad c_2 = -5 \quad \Rightarrow \quad -12 \begin{bmatrix} -2 \\ 2 \\ -3 \end{bmatrix} - 5 \begin{bmatrix} 4 \\ -6 \\ 8 \end{bmatrix} + 2 \begin{bmatrix} -2 \\ -3 \\ 2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix},$$

$$c_4 = 2, c_3 = 0 \Rightarrow c_1 = -10, c_2 = -3 \Rightarrow -10 \begin{bmatrix} -2 \\ 2 \\ -3 \end{bmatrix} - 3 \begin{bmatrix} 4 \\ -6 \\ 8 \end{bmatrix} + 2 \begin{bmatrix} -4 \\ 1 \\ -3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

We can interpret this result thinking that the third and fourth vectors in matrix  $A$  are linear combination of the first two vectors. Therefore, a linearly independent subset of  $S$  having its same span is given by

$$\tilde{S} = \left\{ \begin{bmatrix} -2 \\ 2 \\ -3 \end{bmatrix}, \begin{bmatrix} 4 \\ -6 \\ 8 \end{bmatrix} \right\}.$$

Notice that all the information to find  $\tilde{S}$  is in matrix  $E_A$ , the reduced echelon form of matrix  $A$ ,

$$A = \begin{bmatrix} -2 & 4 & -2 & -4 \\ 2 & -6 & -3 & 1 \\ -3 & 8 & 2 & -3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 6 & 5 \\ 0 & 1 & \frac{5}{2} & \frac{3}{2} \\ 0 & 0 & 0 & 0 \end{bmatrix} = E_A.$$

The columns with pivots in  $E_A$  determine the column vectors in  $A$  that form a linearly independent set. The non-pivot columns in  $E_A$  determine the column vectors in  $A$  that are linear combination of the vectors in the linearly independent set. The factors of these linear combinations are precisely the component of the non-pivot vectors in  $E_A$ . For example, the last column vector in  $E_A$  has components 5 and  $3/2$ , and these are precisely the coefficients in the linear combination:

$$\begin{bmatrix} -4 \\ 1 \\ -3 \end{bmatrix} = 5 \begin{bmatrix} -2 \\ 2 \\ -3 \end{bmatrix} + \frac{3}{2} \begin{bmatrix} 4 \\ -6 \\ 8 \end{bmatrix}.$$

◁

In Example 4.2.3 we answered a question about the linear independence of a set  $S = \{v_1, \dots, v_n\} \subset \mathbb{F}^n$  by studying the properties of a matrix having these vectors a column vectors, that is,  $A = [v_1, \dots, v_n]$ . It turns out that this is a good idea and the following result summarizes few useful relations.

**Theorem 4.2.4.** *Given  $A = [A_{:1}, \dots, A_{:n}] \in \mathbb{F}^{m,n}$ , the following statements are equivalent:*

- (a) *The column vectors set  $\{A_{:1}, \dots, A_{:n}\} \subset \mathbb{F}^m$  is linearly independent;*
- (b)  $N(A) = \{0\}$ ;
- (c)  $\text{rank}(A) = n$

*In the case  $A \in \mathbb{F}^{n,n}$ , the set  $\{A_{:1}, \dots, A_{:n}\} \subset \mathbb{F}^n$  is linearly independent iff  $A$  is invertible.*

**Proof of Theorem 4.2.4:** Let us denote by  $S = \{v_1, \dots, v_n\} \subset \mathbb{F}^m$  a set of vectors in a vector space, and introduce the matrix  $A = [v_1, \dots, v_n]$ . The set  $S$  is linearly independent iff only solution  $c \in \mathbb{R}^n$  to the equation  $Ac = 0$  is the trivial solution  $c = 0$ . This is equivalent to say that  $N(A) = \{0\}$ . This is equivalent to say that  $E_A$  has  $n$  pivot columns, which is equivalent to say that  $\text{rank}(A) = n$ . The furthermore part is straightforward, since an  $n \times n$  matrix  $A$  is invertible iff  $\text{rank}(A) = n$ . This establishes the Theorem.  $\square$

**Further reading.** See Section 4.3 in Meyer's book [3].

## 4.2.3. Exercises.

4.2.1.- Determine which of the following sets is linearly independent. For those who are linearly dependent, express one vector as a linear combination of the other vectors in the set.

$$(a) \left\{ \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 5 \\ 9 \end{bmatrix} \right\};$$

$$(b) \left\{ \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 0 \\ 4 \\ 5 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 6 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right\};$$

$$(c) \left\{ \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} \right\}.$$

4.2.2.- Let  $A = \begin{bmatrix} 2 & 1 & 1 & 0 \\ 4 & 2 & 1 & 2 \\ 6 & 3 & 2 & 3 \end{bmatrix}$ .

- (a) Find a linearly independent set containing the largest possible number of columns of  $A$ .
- (b) Find how many linearly independent sets can be constructed using any number of column vectors of  $A$ .

4.2.3.- Show that any set containing the zero vector must be linearly dependent.

4.2.4.- Given a vector space  $V$ , prove the following: If the set

$$\{v, w\} \subset V$$

is linearly independent, then so is

$$\{(v+w), (v-w)\}.$$

4.2.5.- Determine whether the set

$$\left\{ \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}, \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix}, \begin{bmatrix} 4 & -1 \\ -1 & 1 \end{bmatrix} \right\} \subset \mathbb{R}^{2,2}$$

is linearly independent or dependent.

4.2.6.- Show that the following set in  $\mathbb{P}_2$  is linearly dependent,

$$\{1, x, x^2, 1+x+x^2\}.$$

4.2.7.- Determine whether  $S \subset \mathbb{P}_2$  is a linearly independent set, where

$$S = \{1+x+x^2, 2x-3x^2, 2+x\}.$$

## 4.3. BASES AND DIMENSION

In this Section we introduce a notion that quantifies the size of a vector space. Before doing that, however, we need to separate two main cases, the vector spaces we call finite dimensional from those called infinite dimensional. In the first case, finite dimensional vector spaces, we introduce the notion of a basis. This is a particular type of set in the vector space that is small enough to be a linearly independent set and big enough to span the whole vector space. A basis of a finite dimensional vector space is not unique. However, every basis contains the same number of vectors. This number, called dimension, quantifies the size of the finite dimensional vector space. In the second case above, infinite dimensional vector spaces, we do not introduce here a concept of basis. More structure is needed in the vector space to be able to determine whether or not an infinite sum of vectors converges. We will not discuss these issues here.

**4.3.1. Basis of a vector space.** A particular type of finite sets in a vector space, small enough to be linearly independent and big enough to span the whole vector space, is called a basis of that vector space. Vector spaces having a finite set with these properties are essentially small, and they are called finite dimensional. When there is no finite set that spans the whole vector space, we call that space infinite dimensional. We now highlight these ideas in a more precise way.

**Definition 4.3.1.** A finite set  $\mathcal{S} \subset V$  is called a **finite basis** of a vector space  $V$  iff  $\mathcal{S}$  is linearly independent and  $\text{Span}(\mathcal{S}) = V$ .

The existence of a finite basis is the property that defines the size of the vector space.

**Definition 4.3.2.** A vector space  $V$  is **finite dimensional** iff  $V$  has a finite basis or  $V$  is one of the following two extreme cases:  $V = \emptyset$  or  $V = \{\mathbf{0}\}$ . Otherwise, the vector space  $V$  is called **infinite dimensional**.

In these notes we will often call a finite basis just simply as a basis, without remarking that they contain a finite number of elements. We only study this type of basis, and we do not introduce the concept of an infinite basis. Why don't we define the notion of an infinite basis, since we have already defined the notion of an infinite linearly independent set? Because we do not have any way to define what is the span of an infinite set of vectors. In a vector space, without any further structure, there is no way to know whether an infinite sum converges or not. The notion of convergence needs further structure in a vector space, for example it needs a notion of distance between vectors. So, only in certain vector spaces with a notion of distance it is possible to introduce an infinite basis. We will discuss this subject in a later Chapter.

**EXAMPLE 4.3.1:** We now present several examples.

- Let  $V = \mathbb{R}^2$ , then the set  $\mathcal{S}_2 = \left\{ \mathbf{e}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \mathbf{e}_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\}$  is a basis of  $\mathbb{R}^2$ . Notice that  $\mathbf{e}_i = \mathbf{l}_{:i}$ , that is,  $\mathbf{e}_i$  is the  $i$ -th column of the identity matrix  $\mathbf{l}_2$ . This basis  $\mathcal{S}_2$  is called the standard basis of  $\mathbb{R}^2$ .
- A vector space can have infinitely many bases. For example, a second basis for  $\mathbb{R}^2$  is the set  $\mathcal{U} = \left\{ \mathbf{u}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \mathbf{u}_2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right\}$ . It is not difficult to verify that this set is a basis of  $\mathbb{R}^2$ , since  $\mathcal{U}$  is linearly independent, and  $\text{Span}(\mathcal{U}) = \mathbb{R}^2$ .
- Let  $V = \mathbb{F}^n$ , then the set  $\mathcal{S}_n = \left\{ \mathbf{e}_1 = \mathbf{l}_{:1}, \dots, \mathbf{e}_n = \mathbf{l}_{:n} \right\}$  is a basis of  $\mathbb{F}^n$ , where  $\mathbf{l}_{:i}$  is the  $i$ -th column of the identity matrix  $\mathbf{l}_n$ . This set  $\mathcal{S}_n$  is called the **standard basis** of  $\mathbb{F}^n$ .



(d) A basis for the vector space  $\mathbb{F}^{2,2}$  of all  $2 \times 2$  matrices is the set  $\mathcal{S}_{2,2}$  given by

$$\mathcal{S}_{2,2} = \left\{ \mathbf{E}_{11} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \mathbf{E}_{12} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}, \mathbf{E}_{21} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}, \mathbf{E}_{22} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \right\};$$

This set is linearly independent and  $\text{Span}(\mathcal{S}_{2,2}) = \mathbb{F}^{2,2}$ , since any element  $\mathbf{A} \in \mathbb{F}^{2,2}$  can be decomposed as follows,

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = A_{11} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} + A_{12} \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} + A_{21} \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix} + A_{22} \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}.$$

(e) A basis for the vector space  $\mathbb{F}^{m,n}$  of all  $m \times n$  matrices is the following:

$$\mathcal{S}_{m,n} = \{ \mathbf{E}_{11}, \mathbf{E}_{12}, \dots, \mathbf{E}_{mn} \},$$

where each  $m \times n$  matrix  $\mathbf{E}_{ij}$  is a matrix with all coefficients zero except the coefficient  $(i, j)$  which is equal to one (see previous example). The set  $\mathcal{S}_{m,n}$  is linearly independent, and  $\text{Span}(\mathcal{S}_{m,n}) = \mathbb{F}^{m,n}$ , since

$$\mathbf{A} = [A_{ij}] = \sum_{i=1}^m \sum_{j=1}^n A_{ij} \mathbf{E}_{ij}.$$

(f) Let  $V = \mathbb{P}_n$ , the set of all polynomials with domain  $\mathbb{R}$  and degree less or equal  $n$ . Any element  $\mathbf{p} \in \mathbb{P}_n$  can be expressed as follows,

$$\mathbf{p}(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n,$$

that is equivalent to say that the set

$$\mathcal{S} = \{ \mathbf{p}_0 = 1, \mathbf{p}_1 = x, \mathbf{p}_2 = x^2, \dots, \mathbf{p}_n = x^n \}$$

satisfies  $\mathbb{P}_n = \text{Span}(\mathcal{S})$ . The set  $\mathcal{S}$  is also linearly independent, since

$$\mathbf{q}(x) = c_0 + c_1x + c_2x^2 + \dots + c_nx^n = 0 \quad \Rightarrow \quad c_0 = \dots = c_n = 0.$$

The proof of the latter statement is simple: Compute the  $n$ -th derivative of  $\mathbf{q}$  above, and obtain the equation  $n!c_n = 0$ , so  $c_n = 0$ . Add this information into the  $(n - 1)$ -th derivative of  $\mathbf{q}$  and we conclude that  $c_{n-1} = 0$ . Continue in this way, and you will prove that all the coefficient  $c$ 's vanish. Therefore,  $\mathcal{S}$  is a basis of  $\mathbb{P}_n$ , and it is also called the standard basis of  $\mathbb{P}_n$ . ◁

**EXAMPLE 4.3.2:** Show that the set  $\mathcal{U} = \left\{ \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} \right\}$  is a basis for  $\mathbb{R}^3$ .

**SOLUTION:** We must show that  $\mathcal{U}$  is a linearly independent set and  $\text{Span}(\mathcal{U}) = \mathbb{R}^3$ . Both properties follow from the fact that matrix  $\mathbf{U}$  below, whose columns are the elements in  $\mathcal{U}$ , is invertible,

$$\mathbf{U} = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 1 & -1 \end{bmatrix} \quad \Rightarrow \quad \mathbf{U}^{-1} = \frac{1}{2} \begin{bmatrix} 1 & -1 & 1 \\ 0 & 2 & 0 \\ 1 & 1 & -1 \end{bmatrix}.$$

Let us show that  $\mathcal{U}$  is a basis of  $\mathbb{R}^3$ : Since matrix  $\mathbf{U}$  is invertible, this implies that its reduced echelon form  $\mathbf{E}_\mathbf{U} = \mathbf{I}_3$ , so its column vectors form a linearly independent set. The existence of  $\mathbf{U}^{-1}$  implies that the system of equations  $\mathbf{U}\mathbf{x} = \mathbf{y}$  has a solution  $\mathbf{x} = \mathbf{U}^{-1}\mathbf{y}$  for every  $\mathbf{y} \in \mathbb{R}^3$ , that is,  $\mathbf{y} \in \text{Col}(\mathbf{U}) = \text{Span}(\mathcal{U})$  for all  $\mathbf{y} \in \mathbb{R}^3$ . This means that  $\text{Span}(\mathcal{U}) = \mathbb{R}^3$ . Hence, the set  $\mathcal{U}$  is a basis of  $\mathbb{R}^3$ . ◁

The following definitions will be useful to establish important properties of a basis.

**Definition 4.3.3.** Let  $V$  be a vector space and  $S_n \subset V$  be a subset with  $n$  elements. The set  $S_n$  is a **maximal linearly independent** set iff  $S_n$  is linearly independent and every other set  $\tilde{S}_m$  with  $m > n$  elements is linearly dependent. The set  $S_n$  is a **minimal spanning** set iff  $\text{Span}(S_n) = V$  and every other set  $\tilde{S}_m$  with  $m < n$  elements satisfies  $\text{Span}(\tilde{S}_m) \subsetneq V$ .

A maximal linearly independent set  $S$  is the biggest set in a vector space that is linearly independent. A set cannot be linearly independent if it is too big, since the bigger the set the more probable that one element in the set is a linear combination of the other elements in the set. A minimal spanning set is the smallest set in a vector space that spans the whole space. A spanning set, that is, a set whose span is the whole space, cannot be too small, since the smaller the set the more probable that an element in the vector space is outside the span of the set. The following result provides a useful characterization of a basis: *A basis is a set in the vector space that is both maximal linearly independent and minimal spanning.* In this sense, a basis is a set with the right size, small enough to be linearly independent and big enough to span the whole vector space.

**Theorem 4.3.4.** Let  $V$  be a vector space. The following statements are equivalent:

- (a)  $\mathcal{U}$  is a basis of  $V$ ;
- (b)  $\mathcal{U}$  is a minimal spanning set in  $V$ .
- (c)  $\mathcal{U}$  is a maximal linearly independent set in  $V$ ;

**EXAMPLE 4.3.3:** We showed in Example 4.3.2 above that the set  $\mathcal{U} = \left\{ \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} \right\}$

is a basis for  $\mathbb{R}^3$ . Since this basis has three elements, Theorem 4.3.4 says that any other spanning set in  $\mathbb{R}^3$  cannot have less than three vectors, and any other linearly independent set in  $\mathbb{R}^3$  cannot have more than three vectors. For example, any subset of  $\mathcal{U}$  containing two elements cannot span  $\mathbb{R}^3$ ; the linear combination of two vectors in  $\mathcal{U}$  span a plane in  $\mathbb{R}^3$ . Another example, any set of four vectors in  $\mathbb{R}^3$  must be linearly dependent.  $\triangleleft$

**Proof of Theorem 4.3.4:** We first show part (a)-(b).

( $\Rightarrow$ ) Assume that  $\mathcal{U}$  is a basis of  $V$ . If the set  $\mathcal{U}$  is not a minimal spanning set of  $V$ , that means there exists  $\tilde{\mathcal{U}} = \{\tilde{\mathbf{u}}_1, \dots, \tilde{\mathbf{u}}_{n-1}\}$  such that  $\text{Span}(\tilde{\mathcal{U}}) = V$ . So, every vector in  $\mathcal{U}$  can be expressed as a linear combination of vectors in  $\tilde{\mathcal{U}}$ . Hence, there exists a set of coefficients  $C_{ij}$  such that

$$\mathbf{u}_j = \sum_{i=1}^{n-1} C_{ij} \tilde{\mathbf{u}}_i, \quad j = 1, \dots, n.$$

The reason to order the coefficients  $C_{ij}$  in this form is that they form a matrix  $C = [C_{ij}]$  which is  $(n-1) \times n$ . This matrix  $C$  defines a function  $C: \mathbb{R}^n \rightarrow \mathbb{R}^{n-1}$ , and since  $\text{rank}(C) \leq (n-1) < n$ , this matrix satisfies that  $N(C)$  is nontrivial as a subset of  $\mathbb{R}^n$ . So there exists a nonzero column vector in  $\mathbb{R}^n$  with components  $\mathbf{z} = [z_j] \in \mathbb{R}^n$ , not all components zero, such that  $\mathbf{z} \in N(C)$ , that is,

$$\sum_{j=1}^n C_{ij} z_j = 0, \quad i = 1, \dots, (n-1).$$

What we have found is that the linear combination

$$z_1 \mathbf{u}_1 + \dots + z_n \mathbf{u}_n = \sum_{j=1}^n z_j \mathbf{u}_j = \sum_{j=1}^n z_j \left( \sum_{i=1}^{n-1} C_{ij} \tilde{\mathbf{u}}_i \right) = \sum_{i=1}^{n-1} \left( \sum_{j=1}^n C_{ij} z_j \right) \tilde{\mathbf{u}}_i = \mathbf{0},$$

with at least one of the coefficients  $z_j$  non-zero. This means that the set  $\mathcal{U}$  is not linearly independent. But this contradicts that  $\mathcal{U}$  is a basis. Therefore, the set  $\mathcal{U}$  is a minimal spanning set of  $V$ .

( $\Leftarrow$ ) Assume that  $\mathcal{U}$  is a minimal spanning set of  $V$ . If  $\mathcal{U}$  is not a basis, that means  $\mathcal{U}$  is not a linearly independent set. At least one element in  $\mathcal{U}$  must be a linear combination of the others. Let us arrange the order of the basis vectors such that the vector  $\mathbf{u}_n$  is a linear combination of the other vectors in  $\mathcal{U}$ . Then, the set  $\tilde{\mathcal{U}} = \{\mathbf{u}_1, \dots, \mathbf{u}_{n-1}\}$  must still span  $V$ , that is  $\text{Span}(\tilde{\mathcal{U}}) = V$ . But this contradicts the assumption that  $\mathcal{U}$  is a minimal spanning set of  $V$ .

We now show part (a)-(c).

( $\Leftarrow$ ) Assume that  $\mathcal{U}$  is a maximal linearly independent set in  $V$ . If  $\mathcal{U}$  is not a basis, that means  $\text{Span}(\mathcal{U}) \subsetneq V$ , so there exists  $\mathbf{u}_{n+1} \in V$  such that  $\mathbf{u}_{n+1} \notin \text{Span}(\mathcal{U})$ . Hence, the set  $\tilde{\mathcal{U}} = \{\mathbf{u}_1, \dots, \mathbf{u}_{n+1}\}$  is a linearly independent set. However, this contradicts the assumption that  $\mathcal{U}$  is a maximal linearly independent set. We conclude that  $\mathcal{U}$  is a basis of  $V$ .

( $\Rightarrow$ ) Assume that  $\mathcal{U}$  is a basis of  $V$ . If the set  $\mathcal{U}$  is not a maximal linearly independent set in  $V$ , then there exists a maximal linearly independent set  $\tilde{\mathcal{U}} = \{\tilde{\mathbf{u}}_1, \dots, \tilde{\mathbf{u}}_k\}$ , with  $k > n$ . By the argument given just above,  $\tilde{\mathcal{U}}$  is a basis of  $V$ . By part (b) the set  $\tilde{\mathcal{U}}$  must be a minimal spanning set of  $V$ . However, this is not true, since  $\mathcal{U}$  is smaller and spans  $V$ . Therefore,  $\mathcal{U}$  must be a maximal linearly independent set in  $V$ .

This establishes the Theorem.  $\square$

**4.3.2. Dimension of a vector space.** The characterization of a basis given in Theorem 4.3.4 above implies that the number of elements in a basis is always the same as in any other basis.

**Theorem 4.3.5.** *The number of elements in any basis of a finite dimensional vector space is the same as in any other basis.*

**Proof of Theorem (4.3.5):** Let  $\mathcal{V}_n$  and  $\mathcal{V}_m$  be two bases of a vector space  $V$  with  $n$  and  $m$  elements, respectively. If  $m > n$ , the property that  $\mathcal{V}_m$  is a minimal spanning set implies that  $\text{Span}(\mathcal{V}_n) \subsetneq \text{Span}(\mathcal{V}_m) = V$ . The former inclusion contradicts that  $\mathcal{V}_n$  is a basis. Therefore,  $n = m$ . (A similar proof can be constructed with the maximal linearly independence property of a basis.) This establishes the Theorem.  $\square$

The number of elements in a basis of a finite dimensional vector space is a characteristic of the vector space, so we give that characteristic a name.

**Definition 4.3.6.** *The **dimension** of a finite dimensional vector space  $V$  with a finite basis, denoted as  $\dim V$ , is the number of elements in any basis of  $V$ . The extreme cases of  $V = \emptyset$  and  $V = \{\mathbf{0}\}$  are defined as zero dimensional.*

From the definition above we see that  $\dim\{\mathbf{0}\} = 0$  and  $\dim \emptyset = 0$ .

**EXAMPLE 4.3.4:** We now present several examples.

(a) The set  $\mathcal{S}_n = \{\mathbf{e}_1 = \mathbf{l}_{:1}, \dots, \mathbf{e}_n = \mathbf{l}_{:n}\}$  is a basis for  $\mathbb{F}^n$ , so  $\dim \mathbb{F}^n = n$ .

(b) A basis for the vector space  $\mathbb{F}^{2,2}$  of all  $2 \times 2$  matrices is the set  $\mathcal{S}_{2,2}$  is given by

$$\mathcal{S}_{2,2} = \left\{ \mathbf{E}_{11} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \mathbf{E}_{12} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}, \mathbf{E}_{21} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}, \mathbf{E}_{22} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \right\},$$

so we conclude that  $\dim \mathbb{F}^{2,2} = 4$ .

(c) A basis for the vector space  $\mathbb{F}^{m,n}$  of all  $m \times n$  matrices is the following:

$$\mathcal{S}_{m,n} = \{\mathbf{E}_{11}, \mathbf{E}_{12}, \dots, \mathbf{E}_{mn}\},$$

where we recall that each  $m \times n$  matrix  $E_{ij}$  is a matrix with all coefficients zero except the coefficient  $(i, j)$  which is equal to one. Since the basis  $\mathcal{S}_{m,n}$  contains  $mn$  elements, we conclude that  $\dim \mathbb{F}^{m,n} = mn$ .

- (d) A basis for the vector space  $\mathbb{P}_n$  of all polynomial with degree less or equal  $n$  is the set  $\mathcal{S}$  given by  $\mathcal{S} = \{\mathbf{p}_0 = 1, \mathbf{p}_1 = x, \mathbf{p}_2 = x^2, \dots, \mathbf{p}_n = x^n\}$ . This set has  $n + 1$  elements, so  $\dim \mathbb{P}_n = n + 1$ .

◁

**REMARK:** Any subspace  $W \subset V$  of a vector space  $V$  is itself a vector space, so the definition of basis also holds for  $W$ . Since  $W \subset V$ , we conclude that  $\dim W \leq \dim V$ .

**EXAMPLE 4.3.5:** Consider the case  $V = \mathbb{R}^3$ . It is simple to see in Fig. 37 that  $\dim U = 1$  and  $\dim W = 2$ , where the subspaces  $U$  and  $W$  are spanned by one vector and by two non-collinear vectors, respectively.

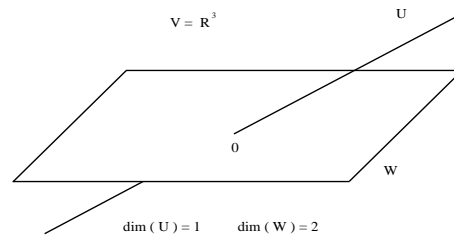


FIGURE 37. Sketch of two subspaces  $U$  and  $W$ , in the vector space  $\mathbb{R}^3$ , of dimension one and two, respectively.

◁

**EXAMPLE 4.3.6:** Find a basis for  $N(\mathbf{A})$  and  $R(\mathbf{A})$ , where matrix  $\mathbf{A} \in \mathbb{R}^{3,4}$  is given by

$$\mathbf{A} = \begin{bmatrix} -2 & 4 & -2 & -4 \\ 2 & -6 & -3 & 1 \\ -3 & 8 & 2 & -3 \end{bmatrix}. \quad (4.2)$$

**SOLUTION:** Since  $\mathbf{A} \in \mathbb{R}^{3,4}$ , then  $\mathbf{A} : \mathbb{R}^4 \rightarrow \mathbb{R}^3$ , which implies that  $N(\mathbf{A}) \subset \mathbb{R}^4$  while  $R(\mathbf{A}) \subset \mathbb{R}^3$ . A basis for  $N(\mathbf{A})$  is found as follows: Find all solution of  $\mathbf{A}\mathbf{x} = \mathbf{0}$  and express these solutions as the span of a linearly independent set of vectors. We first find  $\mathbf{E}_{\mathbf{A}}$ ,

$$\mathbf{A} = \begin{bmatrix} -2 & 4 & -2 & -4 \\ 2 & -6 & -3 & 1 \\ -3 & 8 & 2 & -3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 6 & 5 \\ 0 & 1 & \frac{5}{2} & \frac{3}{2} \\ 0 & 0 & 0 & 0 \end{bmatrix} = \mathbf{E}_{\mathbf{A}} \Rightarrow \begin{cases} x_1 = -6x_3 - 5x_4, \\ x_2 = -\frac{5}{2}x_3 - \frac{3}{2}x_4, \\ x_3, x_4 \text{ free variables.} \end{cases}$$

Therefore, every element in  $N(\mathbf{A})$  can be expressed as follows,

$$\mathbf{x} = \begin{bmatrix} -6x_3 - 5x_4 \\ -\frac{5}{2}x_3 - \frac{3}{2}x_4 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} -6 \\ -\frac{5}{2} \\ 1 \\ 0 \end{bmatrix} x_3 + \begin{bmatrix} -5 \\ -\frac{3}{2} \\ 0 \\ 1 \end{bmatrix} x_4, \Rightarrow N(\mathbf{A}) = \text{Span} \left( \left\{ \begin{bmatrix} -6 \\ -\frac{5}{2} \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} -5 \\ -\frac{3}{2} \\ 0 \\ 1 \end{bmatrix} \right\} \right).$$

Since the vectors in the span above form a linearly independent set, we conclude that a basis for  $N(\mathbf{A})$  is the set  $\mathcal{N}$  given by

$$\mathcal{N} = \left\{ \begin{bmatrix} -6 \\ -\frac{5}{2} \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} -5 \\ -\frac{3}{2} \\ 0 \\ 1 \end{bmatrix} \right\}.$$

We now find a basis for  $R(\mathbf{A})$ . We know that  $R(\mathbf{A}) = \text{Col}(\mathbf{A})$ , that is, the span of the column vectors of matrix  $\mathbf{A}$ . We only need to find a linearly independent subset of column vectors of  $\mathbf{A}$ . This information is given in  $\mathbf{E}_\mathbf{A}$ , since the pivot columns in  $\mathbf{E}_\mathbf{A}$  indicate the columns in  $\mathbf{A}$  which form a linearly independent set. In our case, the pivot columns in  $\mathbf{E}_\mathbf{A}$  are the first and second columns, so we conclude that a basis for  $R(\mathbf{A})$  is the set  $\mathcal{R}$  given by

$$\mathcal{R} = \left\{ \begin{bmatrix} -2 \\ 2 \\ -3 \end{bmatrix}, \begin{bmatrix} 4 \\ -6 \\ 8 \end{bmatrix} \right\}.$$

◀

**4.3.3. Extension of a set to a basis.** We know that a basis of a vector space is not unique, and the following result says that actually any linearly independent set can be extended into a basis of a vector space.

**Theorem 4.3.7.** *If  $S_k = \{\mathbf{u}_1, \dots, \mathbf{u}_k\}$  is a linearly independent set in a vector space  $V$  with basis  $\mathcal{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ , where  $k < n$ , then, there always exists a basis of  $V$  given by an extension of the set  $S_k$  of the form  $\mathcal{S} = \{\mathbf{u}_1, \dots, \mathbf{u}_k, \mathbf{v}_{i_1}, \dots, \mathbf{v}_{i_{n-k}}\}$ .*

The statement above says that a linearly independent set  $S_k$  can be extended into a basis  $\mathcal{S}$  of a vector space  $V$  simply incorporating appropriate vectors from any basis of  $V$ . If a basis  $\mathcal{V}$  of  $V$  has  $n$  vectors, and the set  $S_k$  has  $k < n$  vectors, then one can always select  $n - k$  vectors from the basis  $\mathcal{V}$  to enlarge the set  $S_k$  into a basis of  $V$ .

**Proof of Theorem 4.3.7:** Introduce the set  $S_{k+n}$

$$S_{k+n} = \{\mathbf{u}_1, \dots, \mathbf{u}_k, \mathbf{v}_1, \dots, \mathbf{v}_n\}.$$

We know that  $\text{Span}(S_{k+n}) = V$  since  $\mathcal{V} \subset S_{k+n}$ . We also know that  $S_{k+n}$  is linearly dependent, since the maximal linearly independent set contains  $n$  elements and  $S_{k+n}$  contains  $n + k > n$  elements. The idea is to eliminate the  $\mathbf{v}_i$  such that  $S_k \cup \{\mathbf{v}_i\}$  is linearly dependent. Since the maximal linearly independent set contains  $n$  elements and the  $S_k$  is linearly independent, there are  $k$  elements in  $\mathcal{V}$  that will be eliminated. The resulting set is  $\mathcal{S}$ , which is a basis of  $V$  containing  $S_k$ . This establishes the Theorem.  $\square$

**EXAMPLE 4.3.7:** Given the  $3 \times 4$  matrix  $\mathbf{A}$  defined in Eq. (4.2) in Example 4.3.6 above, extend the basis of  $N(\mathbf{A}) \subset \mathbb{R}^4$  into a basis of  $\mathbb{R}^4$ .

**SOLUTION:** We know from Example 4.3.6 that a basis for the  $N(\mathbf{A})$  is the set  $\mathcal{N}$  given by

$$\mathcal{N} = \left\{ \begin{bmatrix} -6 \\ -\frac{5}{2} \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} -5 \\ -\frac{3}{2} \\ 0 \\ 1 \end{bmatrix} \right\}.$$

Following the idea in the proof of Theorem 4.3.7, we look for a linear independent set of vectors among the columns of the matrix

$$M = \begin{bmatrix} -6 & -5 & 1 & 0 & 0 & 0 \\ -\frac{5}{2} & -\frac{3}{2} & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \end{bmatrix}.$$

That is, matrix  $M$  include the basis vectors of  $N(A)$  and the four vectors  $e_i$  of the standard basis of  $\mathbb{R}^4$ . It is important to place the basis vectors of  $N(A)$  in the first columns of  $M$ . In this way, the Gauss method will select these first vectors as part of the linearly independent set of vectors. Find now the reduced echelon form matrix  $E_M$ ,

$$E_M = \begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 6 & 5 \\ 0 & 0 & 0 & 1 & 5 & 3 \end{bmatrix}.$$

Therefore, the first four vectors in  $M$  are form a linearly independent set, so a basis of  $\mathbb{R}^4$  that includes  $\mathcal{N}$  is given by

$$\mathcal{V} = \left\{ \begin{bmatrix} -6 \\ -\frac{5}{2} \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} -5 \\ -\frac{3}{2} \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \right\}.$$

◁

**4.3.4. The dimension of subspace addition.** Recall that the sum and the intersection of two subspaces is again a subspace in a given vector space. The following result relates the dimension of a sum of subspaces with the dimension of the individual subspaces and the dimension of their intersection.

**Theorem 4.3.8.** *If  $W_1, W_2 \subset V$  are subspaces of a vector space  $V$ , then holds*

$$\dim(W_1 + W_2) = \dim W_1 + \dim W_2 - \dim(W_1 \cap W_2).$$

**Proof of Theorem 4.3.8:** We find the dimension of  $W_1 + W_2$  finding a basis of this sum. The key idea is to start with a basis of  $W_1 \cap W_2$ . Let  $\mathcal{B}_0 = \{z_1, \dots, z_l\}$  be a basis for  $W_1 \cap W_2$ . Enlarge that basis into basis  $\mathcal{B}_1$  for  $W_1$  and  $\mathcal{B}_2$  for  $W_2$  as follows,

$$\mathcal{B}_1 = \{z_1, \dots, z_l, x_1, \dots, x_n\}, \quad \mathcal{B}_2 = \{z_1, \dots, z_l, y_1, \dots, y_m\}.$$

We use the notation  $l = \dim(W_1 \cap W_2)$ ,  $l + n = \dim W_1$  and  $l + m = \dim W_2$ . We now propose as basis for  $W_1 + W_2$  the set

$$\mathcal{B} = \{z_1, \dots, z_l, x_1, \dots, x_n, y_1, \dots, y_m\}.$$

By construction this set satisfies that  $\text{Span}(\mathcal{B}) = W_1 + W_2$ . We only need to show that  $\mathcal{B}$  is linearly independent. Assume that the set  $\mathcal{B}$  is linearly dependent. This means that there is non-zero constants  $a_i, b_j$  and  $c_k$  solutions of the equation

$$\sum_{i=1}^n a_i x_i + \sum_{j=1}^m b_j y_j + \sum_{k=1}^l c_k z_k = \mathbf{0}. \quad (4.3)$$

This implies that the vector  $\sum_{i=1}^n a_i \mathbf{x}_i$ , which by definition belongs to  $W_1$ , also belongs to  $W_2$ , since

$$\sum_{i=1}^n a_i \mathbf{x}_i = -\left(\sum_{j=1}^m b_j \mathbf{y}_j + \sum_{k=1}^l c_k \mathbf{z}_k\right) \in W_2.$$

Therefore,  $\sum_{i=1}^n a_i \mathbf{x}_i$  belongs to  $W_1 \cap W_2$ , and so is a linear combination of the elements of  $\mathcal{B}_0$ , that is, there exists scalars  $d_k$  such that

$$\sum_{i=1}^n a_i \mathbf{x}_i = \sum_{k=1}^l d_k \mathbf{z}_k.$$

Since  $\mathcal{B}_1$  is a basis of  $W_1$ , this implies that all the coefficients  $a_i$  and  $d_k$  vanish. Introduce this information into Eq. (4.3) and we conclude that

$$\sum_{j=1}^m b_j \mathbf{y}_j + \sum_{k=1}^l c_k \mathbf{z}_k = \mathbf{0}.$$

Analogously, the set  $\mathcal{B}_2$  is a basis, so all the coefficients  $b_j$  and  $c_k$  must vanish. This implies that the set  $\mathcal{B}$  is linearly independent, hence a basis of  $W_1 + W_2$ . Therefore, the dimension of the sum is given by

$$\dim(W_1 + W_2) = n + m + k = (n + k) + (m + k) - k = \dim W_1 + \dim W_2 - \dim(W_1 \cap W_2).$$

This establishes the Theorem.  $\square$

The following corollary is immediate.

**Corollary 4.3.9.** *If a vector space can be decomposed as  $V = W_1 \oplus W_2$ , then*

$$\dim(W_1 \oplus W_2) = \dim W_1 + \dim W_2.$$

The proof is straightforward from Theorem 4.3.8, since the condition of subspaces direct sum,  $W_1 \cap W_2 = \{\mathbf{0}\}$ , says that  $\dim(W_1 \cap W_2) = 0$ .

## 4.3.5. Exercises.

4.3.1.- Find a basis for each of the spaces  $N(A)$ ,  $R(A)$ ,  $N(A^T)$ ,  $R(A^T)$ , where

$$A = \begin{bmatrix} 1 & 2 & 2 & 3 \\ 2 & 4 & 1 & 3 \\ 3 & 6 & 1 & 4 \end{bmatrix}.$$

4.3.2.- Find the dimension of the space spanned by

$$\left\{ \begin{bmatrix} 1 \\ 2 \\ -1 \\ 3 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \\ 2 \end{bmatrix}, \begin{bmatrix} 2 \\ 8 \\ -4 \\ 8 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 3 \\ 3 \\ 0 \\ 6 \end{bmatrix} \right\}.$$

4.3.3.- Find the dimension of the following spaces:

- The space  $\mathbb{P}_n$  of polynomials of degree less or equal  $n$ .
- The space  $\mathbb{F}^{m,n}$  of  $m \times n$  matrices.
- The space of real symmetric  $n \times n$  matrices.
- The space of real skew-symmetric  $n \times n$  matrices.

4.3.4.- Find an example to show that the following statement is false: Given a basis  $\{v_1, v_2\}$  of  $\mathbb{R}^2$ , then every subspace  $W \subset \mathbb{R}^2$  has a basis containing at least one of the basis vectors  $v_1, v_2$ .

4.3.5.- Given the matrix  $A$  and vector  $v$ ,

$$A = \begin{bmatrix} 1 & 2 & 2 & 0 & 5 \\ 2 & 4 & 3 & 1 & 8 \\ 3 & 6 & 1 & 5 & 5 \end{bmatrix}, \quad v = \begin{bmatrix} -8 \\ 1 \\ 3 \\ 3 \\ 0 \end{bmatrix},$$

verify that  $v \in N(A)$ , and then find a basis of  $N(A)$  containing  $v$ .

4.3.6.- Determine whether or not the set

$$\mathcal{B} = \left\{ \begin{bmatrix} 2 \\ 3 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix} \right\}$$

is a basis for the subspace

$$\text{Span} \left( \left\{ \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 5 \\ 8 \\ 7 \end{bmatrix}, \begin{bmatrix} 3 \\ 4 \\ 1 \end{bmatrix} \right\} \right) \subset \mathbb{R}^3.$$



## 4.4. VECTOR COMPONENTS

**4.4.1. Ordered bases.** In the previous Section we introduced a finite basis in a vector space. Although a vector space can have different bases, every basis has the same number of elements, which provides a measure of the vector space size, called the dimension of the vector space. In this Section we study another property of a basis. Every vector in a finite dimensional vector space can be expressed in a unique way as a linear combination of the basis vectors. This property can be clearly stated in an ordered basis, which is a basis with the basis vectors given in a specific order.

**Definition 4.4.1.** An *ordered basis* of an  $n$ -dimensional vector space  $V$  is a sequence  $(\mathbf{v}_1, \dots, \mathbf{v}_n)$  of vectors such that the set  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is a basis of  $V$ .

Recall that a sequence is an ordered set, that is, a set with elements given in a particular order.

**EXAMPLE 4.4.1:** The following four ordered basis of  $\mathbb{R}^3$  are all different,

$$\left( \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right), \quad \left( \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right), \quad \left( \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \right), \quad \left( \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \right),$$

however, they determine the same basis  $\mathcal{S}_3 = \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right\}$ . ◁

**4.4.2. Vector components in a basis.** The following result states that given a vector space with an ordered basis, there exists a correspondence between vectors and certain sequences of scalars.

**Theorem 4.4.2.** Let  $V$  be an  $n$ -dimensional vector space over the scalar field  $\mathbb{F}$  with an ordered basis  $(\mathbf{u}_1, \dots, \mathbf{u}_n)$ . Then, every vector  $\mathbf{v} \in V$  determines a unique scalars' sequence  $(v_1, \dots, v_n) \subset \mathbb{F}$  such that

$$\mathbf{v} = v_1 \mathbf{u}_1 + \dots + v_n \mathbf{u}_n. \quad (4.4)$$

And every scalars' sequence  $(v_1, \dots, v_n) \subset \mathbb{F}$  determines a unique vector  $\mathbf{v} \in V$  by Eq. (4.4).

**Proof of Theorem 4.4.2:** Denote by  $\mathcal{U} = (\mathbf{u}_1, \dots, \mathbf{u}_n)$  the ordered basis of  $V$ . Since  $\mathcal{U}$  is a basis,  $\text{Span}(\mathcal{U}) = V$  and  $\mathcal{U}$  is linearly independent. The first property implies that for every  $\mathbf{v} \in V$  there exist scalars  $v_1, \dots, v_n$  such that  $\mathbf{v}$  is a linear combination of the basis vectors, that is,

$$\mathbf{v} = v_1 \mathbf{u}_1 + \dots + v_n \mathbf{u}_n.$$

The second property of a basis implies that the linear combination above is unique. Indeed, if there exists another linear combination

$$\mathbf{v} = \nu_1 \mathbf{u}_1 + \dots + \nu_n \mathbf{u}_n,$$

then  $\mathbf{0} = \mathbf{v} - \mathbf{v} = (v_1 - \nu_1) \mathbf{u}_1 + \dots + (v_n - \nu_n) \mathbf{u}_n$ . Since  $\mathcal{U}$  is linearly independent, this implies that each coefficient above vanishes, so  $v_1 = \nu_1, \dots, v_n = \nu_n$ .

The converse statement is simple to show, since the scalars are given in a specific order. Every scalars' sequence  $(v_1, \dots, v_n)$  determines a unique linear combination with an ordered basis  $(\mathbf{u}_1, \dots, \mathbf{u}_n)$  given by  $v_1 \mathbf{u}_1 + \dots + v_n \mathbf{u}_n$ . This unique linear combination determines a unique vector in the vector space. This establishes the Theorem. □

Theorem 4.4.2 says that there exists a correspondence between vectors in a vector space with an ordered basis and scalars' sequences. This correspondence is called a coordinate map and the scalars are called vector components in the basis. Here is a precise definition.

**Definition 4.4.3.** Let  $V$  be an  $n$ -dimensional vector space over  $\mathbb{F}$  with an ordered basis  $\mathcal{U} = (\mathbf{u}_1, \dots, \mathbf{u}_n)$ . The **coordinate map** is the function  $[\ ]_u : V \rightarrow \mathbb{F}^n$ , with  $[\mathbf{v}]_u = \mathbf{v}_u$ , and

$$\mathbf{v}_u = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \Leftrightarrow \mathbf{v} = v_1 \mathbf{u}_1 + \dots + v_n \mathbf{u}_n.$$

The scalars  $v_1, \dots, v_n$  are called the **vector components** of  $\mathbf{v}$  in the ordered basis  $\mathcal{U}$ .

Therefore, we use the notation  $[\mathbf{v}]_u = \mathbf{v}_u \in \mathbb{F}^n$  for the components of a vector  $\mathbf{v} \in V$  in an ordered basis  $\mathcal{U}$ . We remark that the coordinate map is defined only after an ordered basis is fixed in  $V$ . Different ordered bases on  $V$  determine different coordinate maps between  $V$  and  $\mathbb{F}^n$ . When the situation under study involves only one ordered basis, we suppress the basis subindex. The coordinate map will be denoted by  $[\ ] : V \rightarrow \mathbb{F}^n$  and the vector components by  $\mathbf{v} = [\mathbf{v}]$ . In the particular case that  $V = \mathbb{F}^n$  and the basis is the standard basis  $\mathcal{S}_n$ , then the coordinate map  $[\ ]_s$  is the identity map, so  $\mathbf{v} = [\mathbf{v}]_s = \mathbf{v}$ . In this case we follow the convention established in the first Chapters, that is, we denote vectors in  $\mathbb{F}^n$  by  $\mathbf{v}$  instead of  $\mathbf{v}$ . When the situation under study involves more than one ordered basis we keep the sub-indices in the coordinate map, like  $[\ ]_u$ , and in the vector components, like  $\mathbf{v}_u$ , to keep track of the basis attached to these expressions.

**EXAMPLE 4.4.2:** Let  $V$  be the set of points on the plane with a preferred origin. Let  $\mathcal{S} = (\mathbf{e}_1, \mathbf{e}_2)$  be an ordered basis, pictured in Fig. 38.

- Find the components  $\mathbf{v}_s = [\mathbf{v}]_s \in \mathbb{R}^2$  of the vector  $\mathbf{v} = \mathbf{e}_1 + 3\mathbf{e}_2$  in the ordered basis  $\mathcal{S}$ .
- Find the components  $\mathbf{v}_u = [\mathbf{v}]_u \in \mathbb{R}^2$  of the same vector  $\mathbf{v}$  given in part (a) but now in the ordered basis  $\mathcal{U} = (\mathbf{u}_1 = \mathbf{e}_1 + \mathbf{e}_2, \mathbf{u}_2 = -\mathbf{e}_1 + \mathbf{e}_2)$ .

**SOLUTION:** Part (a) is straightforward to compute, since the definition of component of a vector says that the numbers multiplying the basis vectors in the equation  $\mathbf{v} = \mathbf{e}_1 + 3\mathbf{e}_2$  are the components of the vector, that is,

$$\mathbf{v} = \mathbf{e}_1 + 3\mathbf{e}_2 \Leftrightarrow \mathbf{v}_s = \begin{bmatrix} 1 \\ 3 \end{bmatrix}.$$

Part (b) is more involved. We are looking for numbers  $\tilde{v}_1$  and  $\tilde{v}_2$  such that

$$\mathbf{v} = \tilde{v}_1 \mathbf{u}_1 + \tilde{v}_2 \mathbf{u}_2 \Leftrightarrow \mathbf{v}_u = \begin{bmatrix} \tilde{v}_1 \\ \tilde{v}_2 \end{bmatrix}. \quad (4.5)$$

From the definition of the basis  $\mathcal{U}$  we know the components of the basis vectors in  $\mathcal{U}$  in terms of the standard basis, that is,

$$\begin{aligned} \mathbf{u}_1 = \mathbf{e}_1 + \mathbf{e}_2 &\Leftrightarrow \mathbf{u}_{1s} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \\ \mathbf{u}_2 = -\mathbf{e}_1 + \mathbf{e}_2 &\Leftrightarrow \mathbf{u}_{2s} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}. \end{aligned}$$

In other words, we can write the ordered basis  $\mathcal{U}$  as the column vectors of the matrix  $\mathbf{U}_s = [\mathcal{U}]_s = [[\mathbf{u}_1]_s, [\mathbf{u}_2]_s]$  given by

$$\mathbf{U}_s = [\mathbf{u}_{1s}, \mathbf{u}_{2s}] = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}.$$

Expressing Eq. (4.5) in the standard basis means

$$\mathbf{e}_1 + 3\mathbf{e}_2 = \mathbf{v} = \tilde{v}_1(\mathbf{e}_1 + \mathbf{e}_2) + \tilde{v}_2(-\mathbf{e}_1 + \mathbf{e}_2) \Leftrightarrow \begin{bmatrix} 1 \\ 3 \end{bmatrix} = \tilde{v}_1 \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \tilde{v}_2 \begin{bmatrix} -1 \\ 1 \end{bmatrix}.$$

The last equation on the right is a matrix equation for the unknowns  $\mathbf{v}_u = \begin{bmatrix} \tilde{v}_1 \\ \tilde{v}_2 \end{bmatrix}$ ,

$$\begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \tilde{v}_1 \\ \tilde{v}_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix} \Leftrightarrow \mathbf{U}_s \mathbf{v}_u = \mathbf{v}_s.$$

We find the solution using the Gauss method,

$$\left[ \begin{array}{cc|c} 1 & -1 & 1 \\ 1 & 1 & 3 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 0 & 2 \\ 0 & 1 & 1 \end{array} \right] \Rightarrow \mathbf{v}_u = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \Leftrightarrow \mathbf{v} = 2\mathbf{u}_1 + \mathbf{u}_2.$$

A sketch of what has been computed is in Fig. 38. In this Figure is clear that the vector  $\mathbf{v}$  is fixed, and we have only expressed this fixed vector in as a linear combination of two different bases. It is clear in this Fig. 38 that one has to stretch the vector  $\mathbf{u}_1$  by two and add the result to the vector  $\mathbf{u}_2$  to obtain  $\mathbf{v}$ .  $\triangleleft$

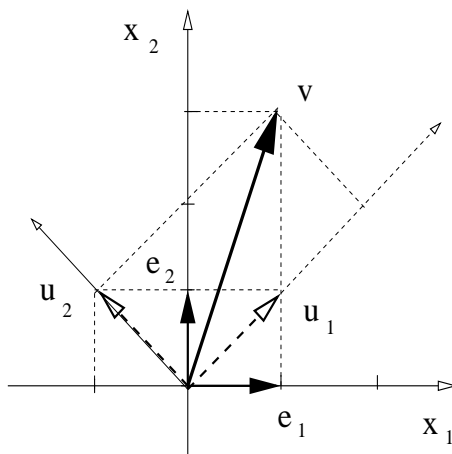


FIGURE 38. The vector  $\mathbf{v} = \mathbf{e}_1 + 3\mathbf{e}_2$  expressed in terms of the basis  $\mathcal{U} = \{\mathbf{u}_1 = \mathbf{e}_1 + \mathbf{e}_2, \mathbf{u}_2 = -\mathbf{e}_1 + \mathbf{e}_2\}$  is given by  $\mathbf{v} = 2\mathbf{u}_1 + \mathbf{u}_2$ .

**EXAMPLE 4.4.3:** Consider the vector space  $\mathbb{P}_2$  of all polynomials of degree less or equal two, and let us consider the case of  $\mathbb{F} = \mathbb{R}$ . An ordered basis is  $\mathcal{S} = (\mathbf{p}_0 = 1, \mathbf{p}_1 = x, \mathbf{p}_2 = x^2)$ . The coordinate map is  $[\ ]_s : \mathbb{P}_2 \rightarrow \mathbb{R}^3$  defined as follows,  $[\mathbf{p}]_s = \mathbf{p}_s$ , where

$$\mathbf{p}_s = \begin{bmatrix} a \\ b \\ c \end{bmatrix} \Leftrightarrow \mathbf{p}(x) = a + bx + cx^2.$$

The column vector  $\mathbf{p}_s$  represents the components of the vector  $\mathbf{p}$  in the ordered basis  $\mathcal{S}$ . The equation above defines a correspondence between every element in  $\mathbb{P}_2$  and every element in  $\mathbb{R}^3$ . The coordinate map depend on the choice of the ordered basis. For example, choosing the ordered basis  $\tilde{\mathcal{S}} = (\mathbf{p}_0 = x^2, \mathbf{p}_1 = x, \mathbf{p}_2 = 1)$ , the corresponding coordinate map is  $[\ ]_{\tilde{s}} : \mathbb{P}_2 \rightarrow \mathbb{R}^3$  defined by  $[\mathbf{p}]_{\tilde{s}} = \mathbf{p}_{\tilde{s}}$ , where

$$\mathbf{p}_{\tilde{s}} = \begin{bmatrix} c \\ b \\ a \end{bmatrix} \Leftrightarrow \mathbf{p}(x) = a + bx + cx^2.$$

The coordinate maps above generalize to the spaces  $\mathbb{P}_n$  and  $\mathbb{R}^{n+1}$  for all  $n \in \mathbb{N}$ . Given the ordered basis  $\mathcal{S} = (\mathbf{p}_0 = 1, \mathbf{p}_1 = x, \dots, \mathbf{p}_n = x^n)$ , the corresponding coordinate map  $[\ ]_s : \mathbb{P}_n \rightarrow \mathbb{R}^{n+1}$  defined by  $[\mathbf{p}]_s = \mathbf{p}_s$ , where

$$\mathbf{p}_s = \begin{bmatrix} a_0 \\ \vdots \\ a_n \end{bmatrix} \Leftrightarrow p(x) = a_0 + \dots + a_n x^n.$$

◁

**EXAMPLE 4.4.4:** Consider  $V = \mathbb{P}_2$  with ordered basis  $\mathcal{S} = (\mathbf{p}_0 = 1, \mathbf{p}_1 = x, \mathbf{p}_2 = x^2)$ .

- (a) Find  $\mathbf{r}_s = [\mathbf{r}]_s$ , the components of  $\mathbf{r}(x) = 3 + 2x + 4x^2$  in the ordered basis  $\mathcal{S}$ .  
 (b) Find  $\mathbf{r}_q = [\mathbf{r}]_q$ , the components of the same polynomial  $\mathbf{r}$  given in part (a) but now in the ordered basis  $\mathcal{Q} = (\mathbf{q}_0 = 1, \mathbf{q}_1 = 1 + x, \mathbf{q}_2 = 1 + x + x^2)$ .

**SOLUTION:**

**Part (a):** This is straightforward to compute, since  $\mathbf{r}(x) = 3 + 2x + 4x^2$  implies that

$$\mathbf{r}(x) = 3\mathbf{p}_0(x) + 2\mathbf{p}_1(x) + 4\mathbf{p}_2(x) \Leftrightarrow \mathbf{r}_s = \begin{bmatrix} 3 \\ 2 \\ 4 \end{bmatrix}.$$

**Part (b):** This is more involved, as in Example 4.4.2. We look for numbers  $\tilde{r}_1, \tilde{r}_2, \tilde{r}_3$  such that

$$\mathbf{r}(x) = \tilde{r}_0 \mathbf{q}_0(x) + \tilde{r}_1 \mathbf{q}_1(x) + \tilde{r}_2 \mathbf{q}_2(x) \Leftrightarrow \mathbf{r}_q = \begin{bmatrix} \tilde{r}_0 \\ \tilde{r}_1 \\ \tilde{r}_2 \end{bmatrix}. \quad (4.6)$$

From the definition of the basis  $\mathcal{Q}$  we know the components of the basis vectors in  $\mathcal{Q}$  in terms of the  $\mathcal{S}$  basis, that is,

$$\begin{aligned} \mathbf{q}_0(x) = \mathbf{p}_0(x) &\Leftrightarrow \mathbf{q}_{0s} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \\ \mathbf{q}_1(x) = \mathbf{p}_0(x) + \mathbf{p}_1(x) &\Leftrightarrow \mathbf{q}_{1s} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \\ \mathbf{q}_2(x) = \mathbf{p}_0(x) + \mathbf{p}_1(x) + \mathbf{p}_2(x) &\Leftrightarrow \mathbf{q}_{2s} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}. \end{aligned}$$

Now we can write the ordered basis  $\mathcal{Q}$  in terms of the column vectors of the matrix  $\mathbf{Q}_s = [\mathcal{Q}]_s = [\mathbf{q}_{0s}, \mathbf{q}_{1s}, \mathbf{q}_{2s}]$ , as follows,

$$\mathbf{Q}_s = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}.$$

Expressing Eq. (4.6) in the standard basis means

$$\begin{bmatrix} 3 \\ 2 \\ 4 \end{bmatrix} = \tilde{r}_0 \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + \tilde{r}_1 \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} + \tilde{r}_2 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}.$$

The last equation on the right is a matrix equation for the unknowns  $\tilde{r}_0$ ,  $\tilde{r}_1$ , and  $\tilde{r}_2$

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \tilde{r}_0 \\ \tilde{r}_1 \\ \tilde{r}_2 \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \\ 4 \end{bmatrix} \Leftrightarrow \mathbf{Q}_s \mathbf{r}_q = \mathbf{r}_s.$$

We find the solution using the Gauss method,

$$\left[ \begin{array}{ccc|c} 1 & 1 & 1 & 3 \\ 0 & 1 & 1 & 2 \\ 0 & 0 & 1 & 4 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & -2 \\ 0 & 0 & 1 & 4 \end{array} \right],$$

hence the solution is

$$\mathbf{r}_q = \begin{bmatrix} 1 \\ -2 \\ 4 \end{bmatrix} \Leftrightarrow \mathbf{r}(x) = \mathbf{q}_0(x) - 2\mathbf{q}_1(x) + 4\mathbf{q}_2(x).$$

We can verify that this is the solution, since

$$\begin{aligned} \mathbf{r}(x) &= \mathbf{q}_0(x) - 2\mathbf{q}_1(x) + 4\mathbf{q}_2(x) \\ &= 1 - 2(1+x) + 4(1+x+x^2) \\ &= (1-2+4) + (-2+4)x + 4x^2 \\ &= 3 + 2x + 4x^2 \\ &= 3\mathbf{p}_0(x) + 2\mathbf{p}_1(x) + 4\mathbf{p}_2(x). \end{aligned}$$

◁

**EXAMPLE 4.4.5:** Given any ordered basis  $\mathcal{U} = (\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3)$  of a 3-dimensional vector space  $V$ , find  $\mathbf{U}_u = [\mathcal{U}]_u \subset \mathbb{F}^{3,3}$ , that is, find  $\mathbf{u}_i = [\mathbf{u}_i]_u$  for  $i = 1, 2, 3$ , the components of the basis vectors  $\mathbf{u}_i$  in its own basis  $\mathcal{U}$ .

**SOLUTION:** The answer is simple: The definition of vector components in a basis says that

$$\begin{aligned} \mathbf{u}_1 &= \mathbf{u}_1 + 0\mathbf{u}_2 + 0\mathbf{u}_3 & \Leftrightarrow & \mathbf{u}_{1u} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \mathbf{e}_1, \\ \mathbf{u}_2 &= 0\mathbf{u}_1 + \mathbf{u}_2 + 0\mathbf{u}_3 & \Leftrightarrow & \mathbf{u}_{2u} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \mathbf{e}_2, \\ \mathbf{u}_3 &= 0\mathbf{u}_1 + 0\mathbf{u}_2 + \mathbf{u}_3 & \Leftrightarrow & \mathbf{u}_{3u} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = \mathbf{e}_3. \end{aligned}$$

In other words, using the coordinate map  $\phi_u : V \rightarrow \mathbb{F}^3$ , we can always write any basis  $\mathcal{U}$  as components in its own basis as follows,  $\mathbf{U}_u = [\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3] = \mathbf{I}_3$ . This example says that there is nothing special about the standard basis  $\mathcal{S} = \{\mathbf{e}_1, \dots, \mathbf{e}_n\}$  of  $\mathbb{F}^n$ , where  $\mathbf{e}_i = \mathbf{l}_i$  is the  $i$ -th column of the identity matrix  $\mathbf{I}_n$ . Given any  $n$ -dimensional vector space  $V$  over  $\mathbb{F}$  with any ordered basis  $\mathcal{V}$ , the components of the basis vectors expressed on its own basis is always the standard basis of  $\mathbb{F}^n$ , that is, the result is always  $[\mathcal{V}]_v = \mathbf{I}_n$ . ◁

## 4.4.3. Exercises.

4.4.1.- Let  $\mathcal{S} = (\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3)$  be the standard basis of  $\mathbb{R}^3$ . Find the components of the vector  $\mathbf{v} = \mathbf{e}_1 + \mathbf{e}_2 + 2\mathbf{e}_3$  in the ordered basis  $\mathcal{U}$

$$\left( \mathbf{u}_{1s} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, \mathbf{u}_{2s} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}, \mathbf{u}_{3s} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \right).$$

4.4.2.- Let  $\mathcal{S} = (\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3)$  be the standard basis of  $\mathbb{R}^3$ . Find the components of the vector

$$\mathbf{v}_s = \begin{bmatrix} 8 \\ 7 \\ 4 \end{bmatrix}$$

in the ordered basis  $\mathcal{U}$  given by

$$\left( \mathbf{u}_{1s} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \mathbf{u}_{2s} = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}, \mathbf{u}_{3s} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \right).$$

4.4.3.- Consider the vector space  $V = \mathbb{P}_2$  with the ordered basis  $\mathcal{S}$  given by

$$\mathcal{S} = (\mathbf{p}_0 = 1, \mathbf{p}_1 = x, \mathbf{p}_2 = x^2).$$

- (a) Find the components of the polynomial  $\mathbf{r}(x) = 2 + 3x - x^2$  in the ordered basis  $\mathcal{S}$ .
- (b) Find the components of the same polynomial  $\mathbf{r}$  given in part (a) but now in the ordered basis  $\mathcal{Q}$  given by  $(\mathbf{q}_0 = 1, \mathbf{q}_1 = 1 - x, \mathbf{q}_2 = x + x^2)$ .

4.4.4.- Let  $\mathcal{S}$  be the standard ordered basis of  $\mathbb{R}^{2,2}$ , that is,

$$\mathcal{S} = (\mathbf{E}_{11}, \mathbf{E}_{12}, \mathbf{E}_{21}, \mathbf{E}_{22}) \subset \mathbb{R}^{2,2},$$

with

$$\mathbf{E}_{11} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{E}_{12} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix},$$

$$\mathbf{E}_{21} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}, \quad \mathbf{E}_{22} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}.$$

- (a) Show that the ordered set  $\mathcal{M}$  below is a basis of  $\mathbb{R}^{2,2}$ , where

$$\mathcal{M} = (\mathbf{M}_1, \mathbf{M}_2, \mathbf{M}_3, \mathbf{M}_4) \subset \mathbb{R}^{2,2},$$

with

$$\mathbf{M}_1 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathbf{M}_2 = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix},$$

$$\mathbf{M}_3 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \mathbf{M}_4 = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix},$$

where the matrices above are written in the standard basis.

- (b) Consider the matrix  $\mathbf{A}$  written in the standard basis  $\mathcal{S}$ ,

$$\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}.$$

Find the components of the matrix  $\mathbf{A}$  in the ordered basis  $\mathcal{M}$ .

## CHAPTER 5. LINEAR TRANSFORMATIONS

## 5.1. LINEAR TRANSFORMATIONS

We now introduce the notion of a linear transformation between vector spaces. Vector spaces are defined not by the elements they are made of but by the relation among these elements. They are defined by the properties of the operation we called linear combinations of the vector space elements. Linear transformations are a very special type of functions that preserve the linear combinations on both vector spaces where they are defined. Linear transformations generalize the definition of a linear function introduced in Sect. 2.1 between the spaces  $\mathbb{F}^n$  and  $\mathbb{F}^m$  to any two vector spaces. Linear transformations are also called in the literature as linear maps or linear mappings.

**Definition 5.1.1.** *Given the vector spaces  $V$  and  $W$  over  $\mathbb{F}$ , the function  $\mathbf{T} : V \rightarrow W$  is called a **linear transformation** iff for all  $\mathbf{u}, \mathbf{v} \in V$  and all scalars  $a, b \in \mathbb{F}$  holds*

$$\mathbf{T}(a\mathbf{u} + b\mathbf{v}) = a\mathbf{T}(\mathbf{u}) + b\mathbf{T}(\mathbf{v}).$$

*In the case that  $V = W$  a linear transformation  $\mathbf{T} : V \rightarrow V$  is called a **linear operator**.*

Taking  $a = b = 0$  in the definition above we see that every linear transformation satisfies that  $\mathbf{T}(\mathbf{0}) = \mathbf{0}$ . As we said in Sect. 2.1, if a function does not satisfy this condition, then it cannot be a linear transformation. We now consider several examples of linear transformations. We do not prove that these examples are indeed linear transformations; the proofs are left to the reader.

**EXAMPLE 5.1.1:** We give several examples of linear transformations.

- (a) The transformation  $\mathbf{I}_V : V \rightarrow V$  defined by  $\mathbf{I}_V(\mathbf{v}) = \mathbf{v}$  for all  $\mathbf{v} \in V$  is called the *identity transformation*. The transformation  $\mathbf{0} : V \rightarrow W$  defined by  $\mathbf{0}(\mathbf{v}) = \mathbf{0}$  for all  $\mathbf{v} \in V$  is called the *zero transformation*.
- (b) Every example of a matrix as a linear function given in Sect. 2.1 is a linear transformation. More precisely, if  $V = \mathbb{F}^n$  and  $W = \mathbb{F}^m$ , then any  $m \times n$  matrix  $\mathbf{A}$  defines a linear transformation  $\mathbf{A} : \mathbb{F}^n \rightarrow \mathbb{F}^m$  by  $\mathbf{A}(\mathbf{x}) = \mathbf{A}\mathbf{x}$ . Therefore, rotations, reflections, projections, dilations are linear transformations in the sense given in Def. 5.1.1 above. In particular, *square matrices are now called linear operators*.
- (c) The vector spaces  $V$  and  $W$  are part of the definition of the linear transformation. For example, a matrix  $\mathbf{A}$  alone does not determine a linear transformation, since the vector spaces must also be specified. For example, and  $m \times n$  matrix  $\mathbf{A}$  defines two different linear transformations, the first one  $\mathbf{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , and the second one  $\mathbf{A} : \mathbb{C}^n \rightarrow \mathbb{C}^m$ . Although the action of these transformations is the same, the matrix-vector product  $\mathbf{A}(\mathbf{x}) = \mathbf{A}\mathbf{x}$ , the linear transformations are different. We will see in Chapter 9 that in the case  $m = n$  these transformations may have different eigenvalues and eigenvectors.
- (d) Let  $V = \mathbb{P}_3$ ,  $W = \mathbb{P}_2$  and let  $\mathbf{D} : \mathbb{P}_3 \rightarrow \mathbb{P}_2$  be the differentiation transformation

$$\mathbf{D}(\mathbf{p}) = \frac{d\mathbf{p}}{dx}.$$

That is, given a polynomial  $\mathbf{p} \in \mathbb{P}_3$ , the transformation  $\mathbf{D}$  acting on  $\mathbf{p}$  is a polynomial one degree less than  $\mathbf{p}$  given by the derivative of  $\mathbf{p}$ . For example,

$$\mathbf{p}(x) = 2 + 3x + x^2 - x^3 \in \mathbb{P}_3 \quad \Rightarrow \quad \mathbf{D}(\mathbf{p})(x) = 3 + 2x - 3x^2 \in \mathbb{P}_2.$$

The transformation  $\mathbf{D}$  is linear, since for all polynomials  $\mathbf{p}, \mathbf{q} \in \mathbb{P}_3$  and scalars  $a, b \in \mathbb{F}$  holds

$$\mathbf{D}(a\mathbf{p} + b\mathbf{q})(x) = \frac{d}{dx} [a\mathbf{p}(x) + b\mathbf{q}(x)] = a \frac{d\mathbf{p}(x)}{dx} + b \frac{d\mathbf{q}(x)}{dx} = a\mathbf{D}(\mathbf{p})(x) + b\mathbf{D}(\mathbf{q})(x).$$

- (e) Notice that the differentiation transformation introduced above can also be defined as a linear operator: Let  $V = W = \mathbb{P}_3$ , and introduce  $\mathbf{D} : \mathbb{P}_3 \rightarrow \mathbb{P}_3$  with the same action as above, that is,  $\mathbf{D}(\mathbf{p}) = \frac{d\mathbf{p}}{dx}$ . The vector spaces used to define the transformation are important. The transformation defined in this example  $\mathbf{D} : \mathbb{P}_3 \rightarrow \mathbb{P}_3$  is different for the one defined above  $\mathbf{D} : \mathbb{P}_3 \rightarrow \mathbb{P}_2$ . Although the action is the same, since the action of both transformations is to take the derivative of their arguments, the transformations are different. We comment on these issues later on.
- (f) Let  $V = \mathbb{P}_2$ ,  $W = \mathbb{P}_3$  and let  $\mathbf{S} : \mathbb{P}_2 \rightarrow \mathbb{P}_3$  be the integral transformation

$$\mathbf{S}(\mathbf{p})(x) = \int_0^x \mathbf{p}(t) dt.$$

For example,

$$\mathbf{p}(x) = 2 + 3x + x^2 \in \mathbb{P}_2 \quad \Rightarrow \quad \mathbf{S}(\mathbf{p})(x) = 2x + \frac{3}{2}x^2 + \frac{1}{3}x^3 \in \mathbb{P}_3.$$

The transformation  $\mathbf{S}$  is linear, since for all polynomials  $\mathbf{p}, \mathbf{q} \in \mathbb{P}_2$  and scalars  $a, b \in \mathbb{F}$  holds

$$\mathbf{S}(a\mathbf{p} + b\mathbf{q})(x) = \int_0^x [a\mathbf{p}(t) + b\mathbf{q}(t)] dt = a \int_0^x \mathbf{p}(t) dt + b \int_0^x \mathbf{q}(t) dt = a\mathbf{S}(\mathbf{p})(x) + b\mathbf{S}(\mathbf{q})(x).$$

- (g) Let  $V = C^1(\mathbb{R}, \mathbb{R})$ , the space of all functions  $\mathbf{f} : \mathbb{R} \rightarrow \mathbb{R}$  having one continuous derivative, that is  $\mathbf{f}'$  is continuous. Let  $W = C^0(\mathbb{R}, \mathbb{R})$ , the space of all continuous functions  $\mathbf{f} : \mathbb{R} \rightarrow \mathbb{R}$ . Notice that  $V \subsetneq W$ , since  $\mathbf{f}(x) = |x|$ , the absolute value function, belongs to  $W$  but it does not belong to  $V$ . The differentiation transformation  $\mathbf{D} : V \rightarrow W$ ,

$$\mathbf{D}(\mathbf{f})(x) = \frac{d\mathbf{f}}{dx}(x),$$

is a linear transformation. The integral transformation  $\mathbf{S} : W \rightarrow V$ ,

$$\mathbf{S}(\mathbf{f})(x) = \int_0^x \mathbf{f}(t) dt,$$

is also a linear transformation.

◁

**5.1.1. The null and range spaces.** The range and null spaces of an  $m \times n$  matrix introduced in Sect. 2.5 can be also be defined for linear transformations between arbitrary vector spaces.

**Definition 5.1.2.** Let  $V, W$  be vector spaces and  $\mathbf{T} : V \rightarrow W$  be a linear transformation. The **null space** of the transformation  $\mathbf{T}$  is the set  $N(\mathbf{T}) \subset V$  given by

$$N(\mathbf{T}) = \{ \mathbf{x} \in V : \mathbf{T}(\mathbf{x}) = \mathbf{0} \}$$

The **range space** of the transformation  $\mathbf{T}$  is the set  $R(\mathbf{T}) \subset W$  given by

$$R(\mathbf{T}) = \{ \mathbf{y} \in W : \mathbf{y} = \mathbf{T}(\mathbf{x}) \text{ for all } \mathbf{x} \in V \}.$$

The null space of a linear transformation is also called the **kernel** of the transformation and denoted as  $\ker(\mathbf{T})$ . The range space of a linear transformation is also called the **image** of the transformation.

**EXAMPLE 5.1.2:** In the case of  $V = \mathbb{R}^n$ ,  $W = \mathbb{R}^m$  and  $\mathbf{A} = \mathbf{A}$  an  $m \times n$  matrix, we have seen many examples of the sets  $N(\mathbf{A})$  and  $R(\mathbf{A})$  in Sect. 2.5. Additional examples in the case of more general linear transformations are the following:



(a) Let  $V = \mathbb{P}_3$ ,  $W = \mathbb{P}_2$ , and let  $\mathbf{D} : \mathbb{P}_3 \rightarrow \mathbb{P}_2$  be the differentiation transformation

$$\mathbf{D}(\mathbf{p}) = \frac{d\mathbf{p}}{dx}.$$

Recalling that  $\mathbf{D}(\mathbf{p}) = 0$  iff  $\mathbf{p}(x) = c$ , with  $c$  constant, we conclude that  $N(\mathbf{D})$  is the set of constant polynomials, that is

$$N(\mathbf{D}) = \text{Span}(\{\mathbf{p}_0 = 1\}) \subset \mathbb{P}_3.$$

Let us now find the set  $R(\mathbf{D})$ . Notice that the derivative of a degree three polynomial is a degree two polynomial, never a degree three. We conclude that  $R(\mathbf{D}) = \mathbb{P}_2$ .

(b) Let  $V = \mathbb{P}_2$ ,  $W = \mathbb{P}_3$ , and let  $\mathbf{S} : \mathbb{P}_2 \rightarrow \mathbb{P}_3$  be the integration transformation

$$\mathbf{S}(\mathbf{p})(x) = \int_0^x \mathbf{p}(t) dt.$$

Recalling that  $\mathbf{S}(\mathbf{p}) = 0$  iff  $\mathbf{p}(x) = 0$ , we conclude that

$$N(\mathbf{S}) = \{\mathbf{p}(x) = 0\} \subset \mathbb{P}_2.$$

Let us now find the set  $R(\mathbf{S})$ . Notice that the integral of a nonzero constant polynomial is a polynomial of degree one. We conclude that

$$R(\mathbf{S}) = \{\mathbf{p}(x) \in \mathbb{P}_3 : \mathbf{p}(x) = a_1x + a_2x^2 + a_3x^3\}.$$

Therefore, non-zero constant polynomials do not belong to  $R(\mathbf{S})$ , so  $R(\mathbf{S}) \subsetneq \mathbb{P}_3$ .

◁

Given a linear transformation  $\mathbf{T} : V \rightarrow W$  between vector spaces  $V$  and  $W$ , it is not difficult to show that the sets  $N(\mathbf{T}) \subset V$  and  $R(\mathbf{T}) \subset W$  are also subspaces of their respective vector spaces.

**Theorem 5.1.3.** *Let  $V$  and  $W$  be vector spaces and  $\mathbf{T} : V \rightarrow W$  be a linear transformation. Then, the sets  $N(\mathbf{T}) \subset V$  and  $R(\mathbf{T}) \subset W$  are subspaces of  $V$  and  $W$ , respectively.*

The proof is formally the same as the proof of Theorem 2.5.3 in Sect. 2.5.

**Proof of Theorem 5.1.3:** The sets  $N(\mathbf{T})$  and  $R(\mathbf{T})$  are subspaces because the transformation  $\mathbf{T}$  is linear. Consider two arbitrary elements  $\mathbf{x}_1, \mathbf{x}_2 \in N(\mathbf{T})$ , that is,  $\mathbf{T}(\mathbf{x}_1) = \mathbf{0}$  and  $\mathbf{T}(\mathbf{x}_2) = \mathbf{0}$ . Then, for any  $a, b \in \mathbb{F}$  holds

$$\mathbf{T}(a\mathbf{x}_1 + b\mathbf{x}_2) = a\mathbf{T}(\mathbf{x}_1) + b\mathbf{T}(\mathbf{x}_2) = \mathbf{0} \Rightarrow (a\mathbf{x}_1 + b\mathbf{x}_2) \in N(\mathbf{T}).$$

Therefore,  $N(\mathbf{T}) \subset V$  is a subspace. Analogously, consider two arbitrary elements  $\mathbf{y}_1, \mathbf{y}_2 \in R(\mathbf{T})$ , that is, there exist  $\mathbf{x}_1, \mathbf{x}_2 \in V$  such that  $\mathbf{y}_1 = \mathbf{T}(\mathbf{x}_1)$  and  $\mathbf{y}_2 = \mathbf{T}(\mathbf{x}_2)$ . Then, for any  $a, b \in \mathbb{F}$  holds

$$(a\mathbf{y}_1 + b\mathbf{y}_2) = a\mathbf{T}(\mathbf{x}_1) + b\mathbf{T}(\mathbf{x}_2) = \mathbf{T}(a\mathbf{x}_1 + b\mathbf{x}_2) \Rightarrow (a\mathbf{y}_1 + b\mathbf{y}_2) \in R(\mathbf{T}).$$

Therefore,  $R(\mathbf{T}) \subset W$  is a subspace. This establishes the Theorem.  $\square$

**5.1.2. Injections, surjections and bijections.** We now specify useful properties a linear transformation might have.

**Definition 5.1.4.** *Let  $V$  and  $W$  be vector spaces and  $\mathbf{T} : V \rightarrow W$  be a linear transformation.*

(a)  $\mathbf{T}$  is **injective** (or *one-to-one*) iff for all vectors  $\mathbf{v}_1, \mathbf{v}_2 \in V$  holds

$$\mathbf{v}_1 \neq \mathbf{v}_2 \Rightarrow \mathbf{T}(\mathbf{v}_1) \neq \mathbf{T}(\mathbf{v}_2).$$

(b)  $\mathbf{T}$  is **surjective** (or *onto*) iff  $R(\mathbf{T}) = W$ .

(c)  $\mathbf{T}$  is **bijective** (or *an isomorphism*) iff  $\mathbf{T}$  is both injective and surjective.

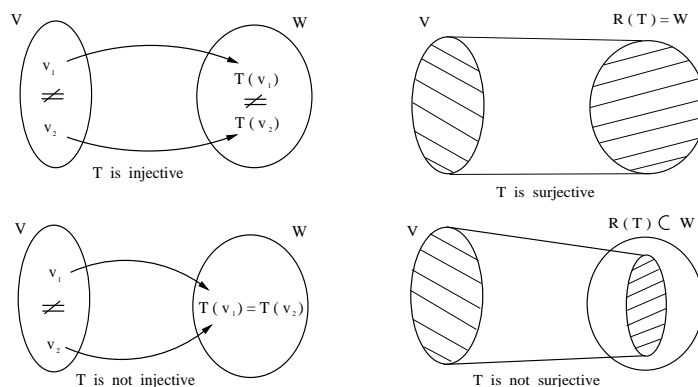


FIGURE 39. Sketch representing an injective and a non-injective function, as well as a surjective and a non-surjective function.

In Fig. 39 we sketch the meaning of these definitions using standard pictures from set theory. Notice that a given transformation can be only injective, or it can be only surjective, or it can be both injective and surjective (bijective), or it can be neither. That a transformation is injective does not imply anything about whether it is surjective or not. That a transformation is surjective does not imply anything about whether it is injective or not. Before we present examples of injective and/or surjective linear transformations it is useful to introduce a result to characterize those transformations that are injective. This can be done in terms of the null space of the transformation.

**Theorem 5.1.5.** *The linear transformation  $T: V \rightarrow W$  is injective iff  $N(T) = \{\mathbf{0}\}$ .*

**Proof of Theorem 5.1.5:**

( $\Rightarrow$ ) Since the transformation  $T$  is injective, given two elements  $\mathbf{v} \neq \mathbf{0}$ , we know that  $T(\mathbf{v}) \neq T(\mathbf{0}) = \mathbf{0}$ , where the last equality comes from the fact that  $T$  is linear. We conclude that the null space of  $T$  contains only the zero vector.

( $\Leftarrow$ ) Since  $N(T) = \{\mathbf{0}\}$ , given any two different elements  $\mathbf{v}_1, \mathbf{v}_2 \in V$ , that is,  $\mathbf{v}_1 - \mathbf{v}_2 \neq \mathbf{0}$ , we know that  $T(\mathbf{v}_1 - \mathbf{v}_2) \neq \mathbf{0}$ , because the only element in  $N(T)$  is the zero vector. Since  $T(\mathbf{v}_1 - \mathbf{v}_2) = T(\mathbf{v}_1) - T(\mathbf{v}_2)$ , we obtain that  $T(\mathbf{v}_1) \neq T(\mathbf{v}_2)$ . We conclude that  $T$  is injective. This establishes the Theorem.  $\square$

**EXAMPLE 5.1.3:** Let  $V = \mathbb{R}^3$ ,  $W = \mathbb{R}^2$ , and consider the linear transformation  $A$  defined by the  $2 \times 3$  matrix  $A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 1 \end{bmatrix}$ . Is  $A$  injective? Is  $A$  surjective?

**SOLUTION:** A simple way to answer these questions is to find bases for the  $N(A)$  and  $R(A)$  spaces. Both bases can be obtained from the information given in  $E_A$ , the reduced echelon form of  $A$ . A simple calculation shows

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 0 \\ 0 & 0 & 1 \end{bmatrix} = E_A.$$

This information gives a basis for  $N(A)$ , since all solutions  $A\mathbf{x} = \mathbf{0}$  are given by

$$\left. \begin{array}{l} x_1 = -2x_2, \\ x_2 \text{ free variable,} \\ x_3 = 0, \end{array} \right\} \Rightarrow \mathbf{x} = \begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix} x_2 \Rightarrow N(A) = \text{Span}\left(\left\{ \begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix} \right\}\right).$$

Since  $N(\mathbf{A}) \neq \{\mathbf{0}\}$ , the transformation associated with  $\mathbf{A}$  is not injective. The reduced echelon form above also says that the first and third column vectors in  $\mathbf{A}$  form a linearly independent set. Therefore, the set

$$\mathcal{R} = \left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 3 \\ 1 \end{bmatrix} \right\}$$

is a basis for  $R(\mathbf{A})$ , so  $\dim R(\mathbf{A}) = 2$ , and then  $R(\mathbf{A}) = \mathbb{R}^2$ . Therefore, the transformation associated with  $\mathbf{A}$  is surjective.  $\triangleleft$

**EXAMPLE 5.1.4:** Consider the differentiation and integration transformations

$$\mathbf{D}: \mathbb{P}_3 \rightarrow \mathbb{P}_2 \quad \mathbf{D}(\mathbf{p})(x) = \frac{d\mathbf{p}}{dx}(x), \quad \mathbf{S}: \mathbb{P}_2 \rightarrow \mathbb{P}_3 \quad \mathbf{S}(\mathbf{p})(x) = \int_0^x \mathbf{p}(t) dt.$$

Show that  $\mathbf{D}$  above is surjective but not injective, and  $\mathbf{S}$  above is injective but not surjective.

**SOLUTION:** Let us start with the differentiation transformation. In Example 5.1.2 we found that  $N(\mathbf{D}) \neq \{\mathbf{0}\}$  and  $R(\mathbf{D}) = \mathbb{P}_2$ . Therefore,  $\mathbf{D}$  above is not injective but it is surjective. Regarding the integration transformation, we have seen in the same Example 5.1.2 that  $N(\mathbf{S}) = \{\mathbf{0}\}$  and  $R(\mathbf{S}) \subsetneq \mathbb{P}_3$ . Therefore,  $\mathbf{S}$  above is injective but it is not surjective.  $\triangleleft$

**5.1.3. Nullity-Rank Theorem.** The following result relates the dimensions of the null and range spaces of a linear transformation on finite dimensional vector spaces. This result is usually called Nullity-Rank Theorem, where the *nullity* of a linear transformation is the dimension of its null space, and the rank is the dimension of its range space. This result is also called the Dimension Theorem.

**Theorem 5.1.6.** Every linear transformation  $\mathbf{T}: V \rightarrow W$  between finite dimensional vector spaces  $V$  and  $W$  satisfies that

$$\dim N(\mathbf{T}) + \dim R(\mathbf{T}) = \dim V. \quad (5.1)$$

**Proof of Theorem 5.1.6:** Let us denote by  $n = \dim V$  and by  $\mathcal{N} = \{\mathbf{u}_1, \dots, \mathbf{u}_k\}$  a basis for  $N(\mathbf{T})$ , where  $0 \leq k \leq n$ , with  $k = 0$  representing the case  $N(\mathbf{T}) = \{\mathbf{0}\}$ . By Theorem 4.3.7 we know we can increase the set  $\mathcal{N}$  into a basis of  $V$ , so let us denote this basis as

$$\mathcal{V} = \{\mathbf{u}_1, \dots, \mathbf{u}_k, \mathbf{v}_1, \dots, \mathbf{v}_l\}, \quad k + l = n.$$

In order to prove Eq. (5.1) we now show that the set

$$\mathcal{R} = \{\mathbf{T}(\mathbf{v}_1), \dots, \mathbf{T}(\mathbf{v}_l)\}$$

is a basis of  $R(\mathbf{T})$ . We first show that  $\text{Span}(\mathcal{R}) = R(\mathbf{T})$ : Given any vector  $\mathbf{v} \in V$ , we can express it in the basis  $\mathcal{V}$  as follows,

$$\mathbf{v} = x_1 \mathbf{u}_1 + \dots + x_k \mathbf{u}_k + y_1 \mathbf{v}_1 + \dots + y_l \mathbf{v}_l.$$

Since  $R(\mathbf{T})$  is the set of vectors of the form  $\mathbf{T}(\mathbf{v})$  for all  $\mathbf{v} \in V$ , therefore  $\mathbf{T}(\mathbf{v}) \in R(\mathbf{T})$  iff

$$\mathbf{T}(\mathbf{v}) = y_1 \mathbf{T}(\mathbf{v}_1) + \dots + y_l \mathbf{T}(\mathbf{v}_l) \in \text{Span}(\mathcal{R}).$$

But this just says that  $R(\mathbf{T}) = \text{Span}(\mathcal{R})$ . Now we show that  $\mathcal{R}$  is linearly independent: Given any linear combination of the form

$$\mathbf{0} = c_1 \mathbf{T}(\mathbf{v}_1) + \dots + c_l \mathbf{T}(\mathbf{v}_l) = \mathbf{T}(c_1 \mathbf{v}_1 + \dots + c_l \mathbf{v}_l),$$

we conclude that  $c_1 \mathbf{v}_1 + \dots + c_l \mathbf{v}_l \in N(\mathbf{T})$ , so there exist  $d_1, \dots, d_k$  such that

$$c_1 \mathbf{v}_1 + \dots + c_l \mathbf{v}_l = d_1 \mathbf{u}_1 + \dots + d_k \mathbf{u}_k,$$

which is equivalent to

$$d_1 \mathbf{u}_1 + \cdots + d_k \mathbf{u}_k - c_1 \mathbf{v}_1 - \cdots - c_l \mathbf{v}_l = \mathbf{0}.$$

Since  $\mathcal{V}$  is a basis we conclude that all coefficients  $c_i$  and  $d_j$  must vanish, for  $i = 1, \dots, l$  and  $j_1, \dots, k$ . This shows that  $\mathcal{R}$  is a linearly independent set. Then  $\mathcal{R}$  is a basis of  $R(\mathbf{T})$  and so,

$$\dim N(\mathbf{T}) = k, \quad \dim R(\mathbf{T}) = l = n - k.$$

This establishes the Theorem.  $\square$

One simple consequence of the Dimension Theorem is that an injective linear operator is also surjective, and vice versa.

**Corollary 5.1.7.** *A linear operator on a finite dimensional vector space is injective iff it is surjective.*

The proof is left as an exercise, Exercise 5.1.7. In the case that the linear transformation is given by an  $m \times n$  matrix, Theorem 5.1.6 establishes a relation between the nullity and the rank of the matrix.

**Corollary 5.1.8.** *Every  $m \times n$  matrix  $A$  satisfies that  $\dim N(A) + \text{rank}(A) = n$ .*

**Proof of Corollary 5.1.8:** This Corollary is just a particular case of Theorem 5.1.6. Indeed, an  $m \times n$  matrix  $A$  defines a linear transformation given by  $\mathbf{A} : \mathbb{F}^n \rightarrow \mathbb{F}^m$  by  $\mathbf{A}(\mathbf{x}) = A\mathbf{x}$ . Since  $\text{rank}(A) = \dim R(A)$  and  $\dim \mathbb{F}^n = n$ , Eq. (5.1) implies that  $\dim N(A) + \text{rank}(A) = n$ . This establishes the Corollary.

An alternative wording of this result uses  $E_A$ , the reduced echelon form of  $A$ . Since  $N(A) = N(E_A)$ , the  $\dim N(A)$  is the number of non-pivot columns in  $E_A$ . We also know that  $\text{rank}(A)$  is the number of pivot columns in  $E_A$ . Therefore  $\dim N(A) + \text{rank}(A)$  is the total number of columns in  $A$ , which is  $n$ .  $\square$

## 5.1.4. Exercises.

5.1.1.- Consider the operator  $\mathbf{A} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  given by the matrix  $\mathbf{A} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ , which projects a vector  $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$  onto the line  $x_1 = x_2$  on  $\mathbb{R}^2$ . Is  $\mathbf{A}$  injective? Is  $\mathbf{A}$  surjective?

5.1.2.- Fix any real number  $\theta \in [0, 2\pi)$ , and define the operator  $\mathbf{R}(\theta) : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  by the matrix  $\mathbf{R}(\theta) = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}$ , which is a rotation by an angle  $\theta$  counterclockwise. Is  $\mathbf{R}(\theta)$  injective? Is  $\mathbf{R}(\theta)$  surjective?

5.1.3.- Let  $\mathbb{F}^{n,n}$  be the vector space of all  $n \times n$  matrices, and fix  $\mathbf{A} \in \mathbb{F}^{n,n}$ . Determine which of the following transformations  $\mathbf{T} : \mathbb{F}^{n,n} \rightarrow \mathbb{F}^{n,n}$  is linear.

- (a)  $\mathbf{T}(\mathbf{X}) = \mathbf{A}\mathbf{X} - \mathbf{X}\mathbf{A}$ ;
- (b)  $\mathbf{T}(\mathbf{X}) = \mathbf{X}^T$ ;
- (c)  $\mathbf{T}(\mathbf{X}) = \mathbf{X}^T + \mathbf{A}$ ;
- (d)  $\mathbf{T}(\mathbf{X}) = \mathbf{A} \operatorname{tr}(\mathbf{X})$ .
- (e)  $\mathbf{T}(\mathbf{X}) = \mathbf{X} + \mathbf{X}^T$ .

5.1.4.- Fix a vector  $\mathbf{v} \in \mathbb{F}^n$  and then define the function  $\mathbf{T} : \mathbb{F}^n \rightarrow \mathbb{F}$  by  $\mathbf{T}(\mathbf{x}) = \mathbf{v}^T \mathbf{x}$ . Show that  $\mathbf{T}$  is a linear transformation. Is  $\mathbf{T}$  a linear operator? Is  $\mathbf{T}$  a linear functional?

5.1.5.- Show that the mapping  $\Delta : \mathbb{P}_3 \rightarrow \mathbb{P}_1$  is a linear transformation, where

$$\Delta(\mathbf{p})(x) = \frac{d^2 \mathbf{p}}{dx^2}(x)$$

Is  $\Delta$  injective? Is  $\Delta$  surjective?

5.1.6.- If the following statement is true, give a proof; if it is false, show it with an example. If a linear transformation on vector spaces  $\mathbf{T} : V \rightarrow W$  is injective, then the image of a linearly independent set in  $V$  is a linearly independent set in  $W$ .

5.1.7.- Prove the following statement: A linear operator  $\mathbf{T} : V \rightarrow V$  on a finite dimensional vector space  $V$  is injective iff  $\mathbf{T}$  is surjective.

5.1.8.- Prove the following statement: If  $V$  and  $W$  are finite dimensional vector spaces with  $\dim V > \dim W$ , then every linear transformation  $\mathbf{T} : V \rightarrow W$  is not injective.

## 5.2. PROPERTIES OF LINEAR TRANSFORMATIONS

**5.2.1. The inverse transformation.** Bijective transformations are invertible. We now introduce the inverse of a bijective transformation and we then show that if the invertible transformation is linear, then the inverse is again a linear transformation.

**Definition 5.2.1.** The *inverse* of a bijective transformation  $\mathbf{T} : V \rightarrow W$  is the transformation  $\mathbf{T}^{-1} : W \rightarrow V$  defined for all  $\mathbf{w} \in W$  and  $\mathbf{v} \in V$  as follows,

$$\mathbf{T}^{-1}(\mathbf{w}) = \mathbf{v} \quad \Leftrightarrow \quad \mathbf{T}(\mathbf{v}) = \mathbf{w}. \quad (5.2)$$

The inverse transformation defined above makes sense only in the case that the original transformation is bijective, hence bijective transformations are called *invertible*.

**EXAMPLE 5.2.1:** Given a real number  $\theta \in (0, \pi)$ , define the transformation  $\mathbf{T} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$

$$\mathbf{T}(\mathbf{x}) = \mathbf{R}_\theta \mathbf{x}, \quad \mathbf{R}_\theta = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}.$$

Find the inverse transformation.

**SOLUTION:** Since the transformation above is a rotation counterclockwise by an angle  $\theta$ , the inverse transformation is a rotation clockwise by the same angle  $\theta$ . In other words,  $(\mathbf{R}_\theta)^{-1} = \mathbf{R}_{-\theta}$ . Therefore, we conclude that,

$$\mathbf{T}^{-1} = (\mathbf{R}_\theta)^{-1} = \mathbf{R}_{-\theta} \quad \Rightarrow \quad (\mathbf{R}_\theta)^{-1} = \begin{bmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{bmatrix},$$

where we have used the relations  $\cos(-\theta) = \cos(\theta)$  and  $\sin(-\theta) = -\sin(\theta)$ .  $\triangleleft$

It is common in the literature to define the inverse matrix in an equivalent way, which we state here as a Theorem.

**Theorem 5.2.2.** The transformation  $\mathbf{T}^{-1} : W \rightarrow V$  is the inverse of a bijective transformation  $\mathbf{T} \in L(V)$  iff holds

$$\mathbf{T}^{-1} \circ \mathbf{T} = \mathbf{I}_V, \quad \mathbf{T} \circ \mathbf{T}^{-1} = \mathbf{I}_W.$$

**Proof of Theorem 5.2.2:**

( $\Rightarrow$ ) Replace  $\mathbf{w} = \mathbf{T}(\mathbf{v})$  into the expression of  $\mathbf{T}^{-1}$ , that is,  $\mathbf{T}^{-1}(\mathbf{T}(\mathbf{v})) = \mathbf{v}$ , for all  $\mathbf{v} \in V$ . This says that  $\mathbf{T}^{-1} \circ \mathbf{T}$  is the identity operator in  $V$ . In a similar way, replace  $\mathbf{v} = \mathbf{T}^{-1}(\mathbf{w})$  into the expression of  $\mathbf{T}$ , that is,  $\mathbf{T}(\mathbf{T}^{-1}(\mathbf{w})) = \mathbf{w}$ , for all  $\mathbf{w} \in W$ . This says that  $\mathbf{T} \circ \mathbf{T}^{-1}$  is the identity in  $W$ .

( $\Leftarrow$ ) If  $\mathbf{w} = \mathbf{T}(\mathbf{v})$ , then  $\mathbf{T}^{-1}(\mathbf{w}) = \mathbf{T}^{-1}(\mathbf{T}(\mathbf{v})) = \mathbf{v}$ . Conversely, if  $\mathbf{v} = \mathbf{T}^{-1}(\mathbf{w})$ , then  $\mathbf{T}(\mathbf{v}) = \mathbf{T}(\mathbf{T}^{-1}(\mathbf{w})) = \mathbf{w}$ . We conclude that

$$\mathbf{T}^{-1}(\mathbf{w}) = \mathbf{v} \quad \Leftrightarrow \quad \mathbf{T}(\mathbf{v}) = \mathbf{w}.$$

This establishes the Theorem.  $\square$

Using this theorem is simple to prove that the result in the previous example can be generalized to any linear transformation defined by an invertible matrix.

**EXAMPLE 5.2.2:** Prove the following statement: If  $\mathbf{T} : \mathbb{F}^n \rightarrow \mathbb{F}^n$  is given by  $\mathbf{T}(\mathbf{x}) = \mathbf{A}\mathbf{x}$  with  $\mathbf{A}$  invertible, then  $\mathbf{T}$  is invertible and  $\mathbf{T}^{-1}(\mathbf{y}) = \mathbf{A}^{-1}\mathbf{y}$ .

**SOLUTION:** Denote  $\mathbf{S}(\mathbf{y}) = \mathbf{A}^{-1}\mathbf{y}$ . Then  $\mathbf{S}(\mathbf{T}(\mathbf{x})) = \mathbf{A}^{-1}\mathbf{A}\mathbf{x} = \mathbf{x}$  for all  $\mathbf{x} \in \mathbb{F}^n$ . This shows that  $\mathbf{S} \circ \mathbf{T} = \mathbf{I}_n$ . The transformation  $\mathbf{S}$  also satisfies  $\mathbf{T}(\mathbf{S}(\mathbf{y})) = \mathbf{A}\mathbf{A}^{-1}\mathbf{y} = \mathbf{y}$  for all  $\mathbf{y} \in \mathbb{F}^n$ . This shows that  $\mathbf{T} \circ \mathbf{S} = \mathbf{I}_n$ . Therefore,  $\mathbf{T}^{-1} = \mathbf{S}$ , that is,  $\mathbf{T}^{-1}(\mathbf{y}) = \mathbf{A}^{-1}\mathbf{y}$ .  $\triangleleft$

**EXAMPLE 5.2.3:** Let  $V = \mathbb{P}_3$ ,  $W = \mathbb{P}_2$  and let  $D : V \rightarrow W$  be the differentiation transformation

$$D(\mathbf{p}) = \frac{d\mathbf{p}}{dx}.$$

Is this transformation invertible?

**SOLUTION:** For every constant polynomial  $\mathbf{p}(x) = c \in \mathbb{P}_3$ , with  $c \in \mathbb{F}$ , holds  $D(c) = 0$ . So the null space of the differentiation transformation is non-trivial, that is,  $N(D) \neq \{\mathbf{0}\}$ . Since the transformation is not injective, it is not bijective. We conclude that **the differentiation transformation above is not invertible.**  $\triangleleft$

**EXAMPLE 5.2.4:** Let  $V = \{\mathbf{p} \in \mathbb{P}_3 : \mathbf{p}(0) = 0\}$ ,  $W = \mathbb{P}_2$  and let  $D : V \rightarrow W$  be the differentiation transformation

$$D(\mathbf{p}) = \frac{d\mathbf{p}}{dx}.$$

Is this transformation invertible?

**SOLUTION:** First, notice that  $V$  is a vector space, since given any two elements  $\mathbf{p}, \mathbf{q} \in V$ , the linear combination satisfies  $(a\mathbf{p} + b\mathbf{q})(0) = a\mathbf{p}(0) + b\mathbf{q}(0) = 0$  for all  $a, b \in \mathbb{F}$ , so  $(a\mathbf{p} + b\mathbf{q}) \in V$ . Also notice that the only constant polynomial  $\mathbf{p}(x) = c$  belonging to  $V$  is the trivial one  $c = 0$ . So, the  $N(D) = \{\mathbf{0}\}$  and the differentiation operator is injective. It is simple to see that this operator is also surjective, since every element in  $W$  also belongs to  $R(D)$ . Indeed, for every constants  $a_2, a_1, a_0 \in \mathbb{F}$  the arbitrary polynomial  $\mathbf{r}(x) = a_2x^2 + a_1x + a_0 \in W$  is the image under  $D$  of the polynomial  $\mathbf{p}(x) = a_2x^3/3 + a_1x^2/2 + a_0x \in V$ , that is,  $D(\mathbf{p}) = \mathbf{r}$ . So,  $W = R(D)$ , the differentiation transformation is surjective, so it is bijective. We conclude that  **$D : V \rightarrow W$  is invertible.**  $\triangleleft$

We now show a couple of properties of the inverse transformation. We start showing that the inverse of a bijective linear transformation is also linear.

**Theorem 5.2.3.** *If the bijective transformation  $T : V \rightarrow W$  is linear, then the inverse transformation  $T^{-1} : W \rightarrow V$  is also a linear transformation.*

**Proof of Theorem 5.2.3:** Since the linear transformation  $T$  is bijective, hence invertible, we now that for every pair of vectors  $\mathbf{v}_1, \mathbf{v}_2 \in V$  there exist vectors  $\mathbf{w}_1, \mathbf{w}_2 \in W$  such that

$$\begin{aligned} T(\mathbf{v}_1) = \mathbf{w}_1 &\Leftrightarrow T^{-1}(\mathbf{w}_1) = \mathbf{v}_1, \\ T(\mathbf{v}_2) = \mathbf{w}_2 &\Leftrightarrow T^{-1}(\mathbf{w}_2) = \mathbf{v}_2. \end{aligned}$$

Since  $T : V \rightarrow W$  is linear, for all  $a_1, a_2 \in \mathbb{F}$  holds

$$\begin{aligned} T(a_1\mathbf{v}_1 + a_2\mathbf{v}_2) &= a_1T(\mathbf{v}_1) + a_2T(\mathbf{v}_2), \\ &= a_1\mathbf{w}_1 + a_2\mathbf{w}_2. \end{aligned}$$

By definition of the inverse transformation, the last equation above is equivalent to

$$\begin{aligned} T^{-1}(a_1\mathbf{w}_1 + a_2\mathbf{w}_2) &= a_1\mathbf{v}_1 + a_2\mathbf{v}_2, \\ &= a_1T^{-1}(\mathbf{w}_1) + a_2T^{-1}(\mathbf{w}_2). \end{aligned}$$

The last equation above says that  $T^{-1} : W \rightarrow V$  is a linear transformation. This establishes the Theorem.  $\square$

We now show that the inverse of a linear transformation is not only a linear transformation but it is also a bijection.

**Theorem 5.2.4.** *If the linear transformation  $T : V \rightarrow W$  is a bijection, then the inverse linear transformation  $T^{-1} : W \rightarrow V$  also is a bijection.*

**Proof of Theorem 5.2.4:** Since  $\mathbf{T}$  is a bijection, it is invertible so, for every  $\mathbf{w}_1, \mathbf{w}_2 \in W$  there exist  $\mathbf{v}_1, \mathbf{v}_2 \in V$  such that

$$\begin{aligned}\mathbf{T}^{-1}(\mathbf{w}_1) = \mathbf{v}_1 &\Leftrightarrow \mathbf{T}(\mathbf{v}_1) = \mathbf{w}_1, \\ \mathbf{T}^{-1}(\mathbf{w}_2) = \mathbf{v}_2 &\Leftrightarrow \mathbf{T}(\mathbf{v}_2) = \mathbf{w}_2.\end{aligned}$$

In order to show that  $\mathbf{T}^{-1}$  is injective pick up  $\mathbf{w}_1$  and  $\mathbf{w}_2$  with  $\mathbf{w}_1 \neq \mathbf{w}_2$  and show that this implies  $\mathbf{v}_1 \neq \mathbf{v}_2$ . This is indeed the case, since subtracting the first line from the second line above we get  $\mathbf{0} \neq \mathbf{w}_2 - \mathbf{w}_1 = \mathbf{T}(\mathbf{v}_2 - \mathbf{v}_1)$ . Since  $\mathbf{T}$  is injective, its null space is trivial, so  $\mathbf{v}_2 - \mathbf{v}_1 \neq \mathbf{0}$ . We then conclude that  $\mathbf{T}^{-1}$  is injective. This inverse transformation is also surjective, which is simple to see from the definition of inverse transformation: For all  $\mathbf{v} \in V$  there exists  $\mathbf{w} \in W$  such that

$$\mathbf{T}(\mathbf{v}) = \mathbf{w} \Leftrightarrow \mathbf{T}^{-1}(\mathbf{w}) = \mathbf{v}.$$

We conclude that  $R(\mathbf{T}^{-1}) = V$ , so  $\mathbf{T}^{-1}$  is surjective. This establishes the Theorem.  $\square$

A simple corollary of Theorem 5.2.4 is that  $\mathbf{T}^{-1}$  is invertible. It is not difficult to show that  $(\mathbf{T}^{-1})^{-1} = \mathbf{T}$ . We have mentioned in Sect. 5.1 that a bijective linear transformation is also called and *isomorphism*. This is the reason for the following definition.

**Definition 5.2.5.** *The vector spaces  $V$  and  $W$  are called **isomorphic** iff there exist an isomorphism  $T : V \rightarrow W$ .*

Different vector spaces which are isomorphic are the same space from the point of view of linear algebra. This means that any linear combination performed in one of the spaces can be translated to the other space through the isomorphism, and any linear operator defined on one of these vector spaces can be translated into a unique linear operator on the other space. This translation is not unique, since there exist infinitely many isomorphisms between isomorphic vector spaces.

**EXAMPLE 5.2.5:** Show that the vector space  $V = \{\mathbf{p} \in \mathbb{P}_3 : \mathbf{p}(0) = 0\}$  is isomorphic to  $\mathbb{P}_2$ .

**SOLUTION:** We have seen in Example 5.2.4 that the differentiation  $\mathbf{D} : V \rightarrow \mathbb{P}_2$  is a bijection, that is, an isomorphism. Therefore, **the spaces  $V$  and  $\mathbb{P}_2$  are isomorphic**. We mention that this isomorphism relates the polynomials in  $V$  and  $\mathbb{P}_2$  as follows

$$V \ni \mathbf{p}(x) = a_3x^3 + a_2x^2 + a_1x \mapsto \frac{d\mathbf{p}}{dx} = 3a_3x^2 + 2a_2x + a_1 \in \mathbb{P}_2.$$

$\triangleleft$

**EXAMPLE 5.2.6:** Show whether the vector spaces  $\mathbb{F}^{2,2}$  and  $\mathbb{F}^4$  are isomorphic or not.

**SOLUTION:** In order to see if these vector spaces are isomorphic, we need to find a bijective transformation between the spaces. We propose the following transformation:  $\mathbf{T} : \mathbb{F}^{2,2} \rightarrow \mathbb{F}^4$  given by

$$\mathbf{T}\left(\begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix}\right) = \begin{bmatrix} x_{11} \\ x_{12} \\ x_{21} \\ x_{22} \end{bmatrix}.$$

It is simple to see that  $\mathbf{T}$  is a bijection, so we conclude that  **$\mathbb{F}^{2,2}$  is isomorphic to  $\mathbb{F}^4$** . We mention that this is not the only isomorphism between these vector spaces. Another



isomorphism is the following:

$$\mathcal{S}\left(\begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix}\right) = \begin{bmatrix} x_{22} \\ x_{21} \\ x_{12} \\ x_{11} \end{bmatrix}.$$

◁

**EXAMPLE 5.2.7:** Show whether the vector spaces  $\mathbb{P}_2(\mathbb{F}, \mathbb{F})$  and  $\mathbb{F}^3$  are isomorphic or not.

**SOLUTION:** In order to see if these vector spaces are isomorphic, we need to find a bijective transformation between the spaces. We propose the following transformation:  $\mathbf{T}: \mathbb{P}_2 \rightarrow \mathbb{F}^3$  given by

$$\mathbf{T}(a_2x^2 + a_1x + a_0) = \begin{bmatrix} a_2 \\ a_1 \\ a_0 \end{bmatrix}.$$

It is simple to see that  $\mathbf{T}$  is a bijection, so we conclude that  $\mathbb{P}_2$  is isomorphic to  $\mathbb{F}^3$ . We mention that this is not the only isomorphism between these vector spaces. Another isomorphism is the following:

$$\mathcal{S}(a_2x^2 + a_1x + a_0) = \begin{bmatrix} a_0 \\ a_2 \\ a_1 \end{bmatrix}.$$

◁

**Theorem 5.2.6.** Every  $n$ -dimensional vector space over the field  $\mathbb{F}$  is isomorphic to  $\mathbb{F}^n$ .

The proof is left as Exercise 5.2.5. This result says that nothing essentially different from  $\mathbb{F}^n$  should be expected from any finite dimensional vector space. On the other hand, infinite dimensional vector spaces are essentially different from  $\mathbb{F}^n$ . Several results that hold for finite dimensional vector spaces do not hold for infinite dimensional ones. An example is the Nullity-Rank Theorem stated in Section 5.1. One could say that finite dimensional vector spaces are all alike; every infinite dimensional vector space is special in its own way.

**5.2.2. The vector space of linear transformations.** Linear transformations from  $V$  to  $W$  can be combined in different ways to obtain new linear transformations. Two linear transformations can be added together and a linear transformation can be multiplied by a scalar. The space of all linear transformations with these operations is a vector space. This is why a linear transformation can be seen either as a function between vector spaces or as a vector in the space of all linear transformations. As an example of these ideas recall the set of all  $m \times n$  matrices, denoted as  $\mathbb{F}^{m,n}$ . This set is a vector space, since two matrices can be added together and a matrix can be multiplied by a scalar. Since any element in this set, an  $m \times n$  matrix  $\mathbf{A}$ , is also a linear transformation  $\mathbf{A}: \mathbb{F}^n \rightarrow \mathbb{F}^m$ , the set of all linear transformations from  $\mathbb{F}^n$  to  $\mathbb{F}^m$  is a vector space. This is why an  $m \times n$  matrix  $\mathbf{A}$  can be interpreted either as a linear transformation between vector spaces  $\mathbf{A}: \mathbb{F}^n \rightarrow \mathbb{F}^m$  or as a vector in the space of matrices,  $\mathbf{A} \in \mathbb{F}^{m,n}$ .

**Definition 5.2.7.** Given the vector spaces  $V$  and  $W$  over  $\mathbb{F}$ , we denote by  $L(V, W)$  the set of all linear transformations from  $V$  to  $W$ . The set  $L(V, V)$  containing all linear operators from  $V$  to  $V$  is denoted as  $L(V)$ . Furthermore, given  $\mathbf{T}, \mathbf{S} \in L(V, W)$  and scalars  $a, b \in \mathbb{F}$ , introduce the **linear combination** of linear transformations as follows

$$(a\mathbf{T} + b\mathbf{S})(\mathbf{v}) = a\mathbf{T}(\mathbf{v}) + b\mathbf{S}(\mathbf{v}) \quad \text{for all } \mathbf{v} \in V.$$

Notice that the zero transformation  $\mathbf{0} : V \rightarrow W$ , given by  $\mathbf{0}(\mathbf{v}) = \mathbf{0}$  for all  $\mathbf{v} \in V$ , belongs to  $L(V, W)$ . Moreover, the space  $L(V, W)$  with the addition and scalar multiplication operations above is indeed a vector space.

**Theorem 5.2.8.** *The set  $L(V, W)$  of linear transformations from  $V$  to  $W$  together with the addition and scalar multiplication operations in Definition 5.2.7 is a vector space.*

The proof is to verify all properties in the Definition 4.1.1, and we left it to the reader. We conclude that a linear transformations can be interpreted either as linear maps between vector spaces or as vectors in the space  $L(V, W)$ . Furthermore, for finite dimensional vector spaces we have the following result.

**Theorem 5.2.9.** *If  $V$  and  $W$  are finite dimensional vector spaces, then  $L(V, W)$  is also finite dimensional and  $\dim L(V, W) = (\dim V)(\dim W)$ .*

**Proof of Theorem 5.2.9:** We introduce a set of linear transformations and we show that this set is indeed a basis for  $L(V, W)$ . Before that let us denote  $n = \dim V$ ,  $m = \dim W$ , and let us introduce both a basis  $\mathcal{V} = \{\mathbf{v}_i\}$  of  $V$  and a basis  $\mathcal{W} = \{\mathbf{w}_j\}$  of  $W$ , where the indices take values  $i = 1, \dots, n$  and  $j = 1, \dots, m$ . The proposal for a basis of  $L(V, W)$  is the set  $\mathcal{S} = \{\mathbf{S}_{ij}\}$ , with elements defined as follows,

$$\mathbf{S}_{ij}(\mathbf{v}_k) = \begin{cases} \mathbf{w}_j & \text{for } i = k, \\ \mathbf{0} & \text{for } i \neq k, \end{cases} \quad \text{where } i = 1, \dots, n \text{ and } j = 1, \dots, m.$$

So, the map  $\mathbf{S}_{ij}$  transforms the basis vector  $\mathbf{v}_i$  into the basis vector  $\mathbf{w}_j$  and the rest of the basis vectors in  $\mathcal{V}$  into the zero vector in  $W$ . We now show that  $\mathcal{S}$  is a basis for  $L(V, W)$ . Let us first prove that  $\mathcal{S}$  spans  $L(V, W)$ . Indeed, for every linear transformation  $\mathbf{T} \in L(V, W)$  and every vector  $\mathbf{v} \in V$  holds,

$$\mathbf{T}(\mathbf{v}) = \sum_{i=1}^n v_i \mathbf{T}(\mathbf{v}_i),$$

where we used the decomposition  $\mathbf{v} = \sum_{i=1}^n v_i \mathbf{v}_i$  of vector  $\mathbf{v}$  in the basis  $\mathcal{V}$ . Since  $\mathbf{T}(\mathbf{v}_i) \in W$ , it can be decomposed in the  $\mathcal{W}$  basis as follows,

$$\mathbf{T}(\mathbf{v}_i) = \sum_{j=1}^m T_{ij} \mathbf{w}_j.$$

Replacing the last equation on the previous one we get,

$$\mathbf{T}(\mathbf{v}) = \sum_{i=1}^n \left( \sum_{j=1}^m v_i T_{ij} \mathbf{w}_j \right),$$

The basis vector  $\mathbf{w}_j$  can be written in terms of the maps  $\mathbf{S}_{ij}$  as  $\mathbf{w}_j = \mathbf{S}_{kj}(\mathbf{v}_k)$ . Choosing index  $k$  to be the value of index  $i$  in the sums above we get,

$$\mathbf{T}(\mathbf{v}) = \sum_{i=1}^n \left( \sum_{j=1}^m v_i T_{ij} \mathbf{S}_{ij}(\mathbf{v}_i) \right) = \sum_{i=1}^n \left( \sum_{j=1}^m T_{ij} \mathbf{S}_{ij}(v_i \mathbf{v}_i) \right).$$

Here is the critical step: Since  $\mathbf{S}_{ij}(\mathbf{v}_k) = \mathbf{0}$  for  $k \neq i$ , the following equation holds,  $\mathbf{S}_{ij}(v_i \mathbf{v}_i) = \mathbf{S}_{ij}(\mathbf{v})$ . Introducing this property in the equation above we get

$$\mathbf{T}(\mathbf{v}) = \sum_{i=1}^n \left( \sum_{j=1}^m T_{ij} \mathbf{S}_{ij}(\mathbf{v}) \right) \quad \text{for all } \mathbf{v} \in V,$$

that is,

$$\mathbf{T} = \sum_{i=1}^n \sum_{j=1}^m T_{ij} \mathbf{S}_{ij} \Rightarrow \mathbf{T} \in \text{Span}(\mathcal{S}) \subset L(V, W).$$

Since the map  $\mathbf{T}$  is an arbitrary element in  $L(V, W)$ , we conclude that  $L(V, W) \subset \text{Span}(\mathcal{S})$ , hence  $\text{Span}(\mathcal{S}) = L(V, W)$ . We now show that  $\mathcal{S}$  is a linearly independent set. Indeed, suppose there exist scalars  $c_{ij}$  satisfying

$$\sum_{i=1}^n \sum_{j=1}^m c_{ij} \mathbf{S}_{ij} = \mathbf{0}.$$

Interchange the sums and evaluate this expression at the basis vector  $\mathbf{v}_k$ , we obtain,

$$\mathbf{0} = \sum_{j=1}^m \left( \sum_{i=1}^n c_{ij} \mathbf{S}_{ij}(\mathbf{v}_k) \right) = \sum_{j=1}^m c_{kj} \mathbf{w}_j.$$

Since the set  $\mathcal{W}$  is a basis of  $W$ , this implies that the coefficients  $c_{kj} = 0$  for  $j = 1, \dots, m$ . The same argument holds for every  $k = 1, \dots, n$ , so we conclude that every coefficient  $c_{ij}$  vanishes. The set  $\mathcal{S}$  is linearly independent, and so, it is a basis for  $L(V, W)$ . Since the set contains  $nm$  elements, we conclude that  $\dim L(V, W) = (\dim V)(\dim W)$ . This establishes the Theorem.  $\square$

**5.2.3. Linear functionals and the dual space.** An important example of linear transformations is the case where the target vector space  $W$  is the set of scalars  $\mathbb{F}$ , the latter, let us recall, is always a vector space. We often use boldface Greek letters to denote linear functionals.

**Definition 5.2.10.** Given the vector space  $V$  over  $\mathbb{F}$ , the linear transformation  $\phi : V \rightarrow \mathbb{F}$  is called a **linear functional**. The vector space of linear functionals  $L(V, \mathbb{F})$  is denoted as  $V^*$  and called the **dual space of  $V$** .

**EXAMPLE 5.2.8:** Show that the projection onto the first component  $\pi_1 : \mathbb{F}^3 \rightarrow \mathbb{F}$  is a linear functional, where

$$\pi_1 \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) = x_1.$$

**SOLUTION:** We only need to show that  $\pi_1$  is a linear transformation. This is indeed the case, since for all  $\mathbf{x}, \mathbf{y} \in \mathbb{F}^3$  and all  $a, b \in \mathbb{F}$  holds

$$\pi_1(a\mathbf{x} + b\mathbf{y}) = ax_1 + by_1 = a\pi_1(\mathbf{x}) + b\pi_1(\mathbf{y}).$$

Since the map is defined from  $\mathbb{F}^3$  into  $\mathbb{F}$ , we conclude that  $\pi_1$  is a linear functional.  $\triangleleft$

**EXAMPLE 5.2.9:** Consider the vector space  $V = C([a, b], \mathbb{R})$  of continuous real-valued functions with domain in the interval  $[a, b] \subset \mathbb{R}$ . Show that the function  $\phi : V \rightarrow \mathbb{R}$  given below is a linear functional, where

$$\phi(\mathbf{f}) = \int_a^b \mathbf{f}(x) dx.$$

**SOLUTION:** We only need to show that the function  $\phi : V \rightarrow \mathbb{R}$  is a linear transformation. This is the case, since for all  $\mathbf{f}, \mathbf{g} \in V$  and all scalars  $c, d \in \mathbb{R}$  holds,

$$\begin{aligned}\phi(c\mathbf{f} + d\mathbf{g}) &= \int_a^b (c\mathbf{f}(x) + d\mathbf{g}(x)) dx \\ &= c \int_a^b \mathbf{f}(x) dx + d \int_a^b \mathbf{g}(x) dx \\ &= c\phi(\mathbf{f}) + d\phi(\mathbf{g}).\end{aligned}$$

We conclude that  $\phi \in V^*$ . ◁

**EXAMPLE 5.2.10:** Show that the trace map  $\text{tr} : \mathbb{F}^{n,n} \rightarrow \mathbb{F}$  is a linear functional.

**SOLUTION:** We have already seen that the trace map is a linear map. Since the trace of a matrix is a scalar, we conclude that the trace map is a linear functional, so  $\text{tr} \in (\mathbb{F}^{n,n})^*$ . ◁

The following result says that the dual space of a finite dimensional vector space is isomorphic to the original vector space. This means that the vector space and its dual are the same from the point of view of linear algebra. This result is not true for infinite dimensional vector spaces. This is one reason why infinite dimensional vector spaces, like function spaces, have more structure and are more complicated to study than finite dimensional vector spaces.

**Theorem 5.2.11.** *Every finite dimensional vector space is isomorphic to its dual space.*

We give two proofs of the same result. In the first proof, we show that an  $n$ -dimensional vector space  $V$  is isomorphic to  $\mathbb{F}^n$ , and that  $V^*$  is isomorphic to  $\mathbb{F}^n$ . We conclude that  $V$  is isomorphic to  $V^*$ . The second proof is more abstract, and makes use of what is called a dual basis of  $V^*$ .

**Proof of Theorem 5.2.11:** (First proof.) We already know that the ordered basis  $\mathcal{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$  of an  $n$ -dimensional vector space  $V$  determines the isomorphism  $[\ ]_v : V \rightarrow \mathbb{F}^n$ , called the coordinate map, as follows

$$[\mathbf{x}]_v = \mathbf{x}_v = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}_v \Leftrightarrow \mathbf{x} = x_1\mathbf{v}_1 + \dots + x_n\mathbf{v}_n.$$

So, any vector  $\mathbf{x} \in V$  can be associated with a column vector  $\mathbf{x}_v \in \mathbb{F}^n$ . We now show that the dual space  $V^*$  is also isomorphic to the vector space  $\mathbb{F}^n$ . Once this is shown we can conclude that  $V$  is isomorphic to  $V^*$ . An arbitrary linear functional  $\phi \in V^*$  satisfies that

$$\begin{aligned}\phi(\mathbf{x}) &= \phi(x_1\mathbf{v}_1 + \dots + x_n\mathbf{v}_n) \\ &= x_1\phi(\mathbf{v}_1) + \dots + x_n\phi(\mathbf{v}_n) \\ &= [\phi(\mathbf{v}_1), \dots, \phi(\mathbf{v}_n)]_v \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}_v \in \mathbb{F}\end{aligned}$$

for all  $\mathbf{x} \in V$ . Denoting the scalars  $\phi_i = \phi(\mathbf{v}_i)$ , for  $i = 1, \dots, n$ , we introduce the coordinate map  $[\ ]_v : V^* \rightarrow \mathbb{F}^n$  as follows

$$[\phi]_v = [\phi_1, \dots, \phi_n]_v \Leftrightarrow \phi(\mathbf{x}) = [\phi_1, \dots, \phi_n]_v \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}_v.$$

The coordinate map defined on the space  $V^*$  of linear functionals associates a linear functional with a row vector. We have chosen a row vector instead of a column vector so we can use the matrix-vector product to express the action of  $\phi$  on a vector  $\mathbf{x}$ . We conclude that, once a basis is fixed in a finite dimensional vector space, vectors can be associated with column vectors, while linear functionals can be associated with row vectors. This establishes the Theorem.  $\square$

**Proof of Theorem 5.2.11:** (Second proof.) Since  $V^* = L(V, \mathbb{F})$ , by Theorem 5.2.9 we know that  $\dim V^* = \dim V$ , where we used the fact the space of scalars,  $\mathbb{F}$ , is itself a vector space and  $\dim \mathbb{F} = 1$ . From the proof of Theorem 5.2.9 we know that a basis of  $V^*$  is given by the set of linear functionals  $\Phi = \{\phi_i\}$  defined as follows: Given a basis  $\mathcal{V} = \{\mathbf{v}_i\}$  of  $V$ , for  $i = 1, \dots, n = \dim V$ , the linear functionals  $\phi_i$  satisfy

$$\phi_i : V \rightarrow \mathbb{F}, \quad \phi_i(\mathbf{v}_k) = \begin{cases} 1 & \text{for } i = k, \\ 0 & \text{for } i \neq k, \end{cases} \quad i, k = 1, \dots, n.$$

Such a basis  $\Phi$  is called the dual basis of  $\mathcal{V}$ . Now that we have a basis in  $V$  and a basis in  $V^*$  it is simple to find an isomorphism between these spaces. We propose the map  $\mathbf{R} : V \rightarrow V^*$  that changes one basis into the other,

$$\mathbf{R}(\mathbf{v}_i) = \phi_i.$$

We now show that this linear transformation  $\mathbf{R}$  is an isomorphism. It is injective because of the following argument. Consider an arbitrary element  $\mathbf{v} \in N(\mathbf{R})$ , that is,

$$\mathbf{R}(\mathbf{v}) = \mathbf{0},$$

where  $\mathbf{0} : V \rightarrow \mathbb{F}$  is the zero map. Introducing the basis decomposition  $\mathbf{v} = \sum_{i=1}^n v_i \mathbf{v}_i$  in the expression above we get

$$\mathbf{0} = \sum_{i=1}^n v_i \mathbf{R}(\mathbf{v}_i) = \sum_{i=1}^n v_i \phi_i.$$

Since the set  $\Phi$  is a basis of  $V^*$ , then  $\Phi$  is linearly independent, which implies that all coefficients  $v_i$  above vanish. We conclude that the vector  $\mathbf{v} = \sum_{i=1}^n v_i \mathbf{v}_i = \mathbf{0}$ , that is, the null space of  $\mathbf{R}$  is trivial. Hence  $\mathbf{R}$  is injective. This property together with the Nullity-Rank Theorem and the fact that  $\dim V^* = \dim V$  imply that  $\mathbf{R}$  is also surjective. We conclude that  $V$  is isomorphic to  $V^*$ . This establishes the Theorem.  $\square$

In the second proof above we introduced a particular basis of the dual space, called the dual basis.

**Definition 5.2.12.** Given an  $n$ -dimensional vector space  $V$  with a basis  $\mathcal{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ , the **dual basis of  $\mathcal{V}$**  is a particular basis  $\Phi = \{\phi_1, \dots, \phi_n\}$  of the dual space  $V^*$  that for  $i, j = 1, \dots, n$  the basis linear functionals  $\phi_i$  satisfy

$$\phi_i(\mathbf{v}_j) = \begin{cases} 1 & \text{for } i = j, \\ 0 & \text{for } i \neq j. \end{cases}$$

**EXAMPLE 5.2.11:** Find the dual basis of the basis  $\mathcal{U} \subset \mathbb{R}^2$ , where

$$\mathcal{U} = \left\{ \mathbf{u}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \mathbf{u}_2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right\}.$$

**SOLUTION:** We need to find a set  $\Phi = \{\phi_1, \phi_2\}$  of linear functionals satisfying the equations

$$\begin{aligned} \phi_1(\mathbf{u}_1) &= 1, & \phi_2(\mathbf{u}_1) &= 0, \\ \phi_1(\mathbf{u}_2) &= 0, & \phi_2(\mathbf{u}_2) &= 1. \end{aligned}$$

The two equations on the left can be solved independently of the two equations on the right. Recall that the most general expression for the linear functionals  $\phi_1 : \mathbb{R}^2 \rightarrow \mathbb{R}$  and  $\phi_2 : \mathbb{R}^2 \rightarrow \mathbb{R}$  is given by

$$\phi_1\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = a_1x_1 + a_2x_2, \quad \phi_2\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = b_1x_1 + b_2x_2.$$

We only need to find the components  $a_1, a_2$  and  $b_1, b_2$ . The  $a$ 's can be obtained from the two equations on the left above, that is,

$$\left. \begin{array}{l} \phi_1(\mathbf{u}_1) = 1 \quad \Rightarrow \quad \phi_1\left(\begin{bmatrix} 1 \\ 1 \end{bmatrix}\right) = a_1 + a_2 = 1 \\ \phi_1(\mathbf{u}_2) = 0 \quad \Rightarrow \quad \phi_1\left(\begin{bmatrix} -1 \\ 1 \end{bmatrix}\right) = -a_1 + a_2 = 0. \end{array} \right\} \Rightarrow \begin{array}{l} a_1 = 1/2, \\ a_2 = 1/2. \end{array}$$

A similar calculation for  $\phi_2$  implies,

$$\left. \begin{array}{l} \phi_2(\mathbf{u}_1) = 0 \quad \Rightarrow \quad \phi_2\left(\begin{bmatrix} 1 \\ 1 \end{bmatrix}\right) = b_1 + b_2 = 0 \\ \phi_2(\mathbf{u}_2) = 1 \quad \Rightarrow \quad \phi_2\left(\begin{bmatrix} -1 \\ 1 \end{bmatrix}\right) = -b_1 + b_2 = 1. \end{array} \right\} \Rightarrow \begin{array}{l} b_1 = -1/2, \\ b_2 = 1/2. \end{array}$$

So we conclude that the dual basis of  $\mathcal{U}$  is the basis  $\Phi = \{\phi_1, \phi_2\}$  defined as follows

$$\phi_1\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \frac{1}{2}(x_1 + x_2), \quad \phi_2\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \frac{1}{2}(-x_1 + x_2).$$

If we associate the linear functionals above with row vectors in  $\mathbb{R}^2$  with the isomorphism

$$[\phi] = [a_1, a_2] \quad \Leftrightarrow \quad \phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = a_1x_1 + a_2x_2,$$

then, the answer of the Example is given by the row vectors

$$[\phi_1] = \frac{1}{2}[1, 1], \quad [\phi_2] = \frac{1}{2}[-1, 1].$$

Associating linear functionals components with row vector is useful, because the matrix-vector product can be used to represent the action of a functional onto a vector. The scalar obtained by the action of a linear functional onto a vector is the scalar given by the matrix-vector product of a row vector times a column vector. For example,

$$\begin{aligned} \phi_1(\mathbf{u}_1) &= \frac{1}{2}[1, 1] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1, & \phi_1(\mathbf{u}_2) &= \frac{1}{2}[1, 1] \begin{bmatrix} -1 \\ 1 \end{bmatrix} = 0, \\ \phi_2(\mathbf{u}_1) &= \frac{1}{2}[-1, 1] \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 0, & \phi_2(\mathbf{u}_2) &= \frac{1}{2}[-1, 1] \begin{bmatrix} -1 \\ 1 \end{bmatrix} = 1. \end{aligned}$$

◁

## 5.2.4. Exercises.

5.2.1.- Show that the linear transformation  $T : \mathbb{F}^2 \rightarrow \mathbb{F}^2$  given below is invertible and find the inverse transformation, where

$$T\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1 + x_2 \\ 3x_1 - 2x_2 \end{bmatrix}.$$

5.2.2.- Consider the vector space

$$V = \{\mathbf{p} \in \mathbb{P}_4 : \mathbf{p}(0) = 0, \frac{d\mathbf{p}}{dx}(0) = 0\}.$$

Then show that the linear transformation  $\Delta : V \rightarrow \mathbb{P}_2$  given below is invertible and find the inverse transformation, where

$$\Delta(\mathbf{p}) = \frac{d^2\mathbf{p}}{dx^2}.$$

5.2.3.- Use the Nullity-Rank Theorem 5.1.6 to show that if the finite dimensional vector spaces  $V$  and  $W$  are isomorphic, then they have the same dimension.

5.2.4.- Show that the spaces  $\mathbb{P}_n$  and  $\mathbb{F}^{n+1}$  are isomorphic.

5.2.5.- Use the coordinate map to prove Theorem 5.2.6: Every  $n$ -dimensional vector space over the field  $\mathbb{F}$  is isomorphic to  $\mathbb{F}^n$ .

5.2.6.- Show that the projection onto the  $i$ -component  $\pi_i : \mathbb{F}^n \rightarrow \mathbb{F}$  is a linear functional, where

$$\pi_i\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}\right) = x_i.$$

5.2.7.- Let  $V = \mathbb{P}_2([0, 1])$  be the vector space of polynomials up to degree two on the domain  $[0, 1] \subset \mathbb{R}$ . Show that the function  $\phi : V \rightarrow \mathbb{R}$  given below is a linear functional, where

$$\phi(\mathbf{p}) = \int_0^1 \mathbf{p}(x) dx.$$

5.2.8.- Given the basis  $\mathcal{U} \subset \mathbb{R}^2$  below, find its dual basis, where

$$\mathcal{U} = \left\{ \mathbf{u}_1 = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \mathbf{u}_2 = \begin{bmatrix} 3 \\ 1 \end{bmatrix} \right\}.$$

## 5.3. THE ALGEBRA OF LINEAR OPERATORS

A vector space where one can multiply two vectors to obtain a third vector is called an algebra. We have seen that the set  $L(V, W)$  is a vector space. A multiplication can be defined between linear transformations in many different ways, but the most interesting one from the physical point of view is the composition of linear transformations. If  $T: U \rightarrow V$  and  $S: V \rightarrow W$  are linear transformations, then  $S \circ T: U \rightarrow W$  defined as

$$(S \circ T)(\mathbf{u}) = S(T(\mathbf{u})) \quad \text{for all } \mathbf{u} \in U,$$

is also a linear transformation. This operation is not defined on a single vector space, but on three different spaces, namely,  $L(U, V)$ ,  $L(V, W)$  and  $L(U, W)$ . However, in the particular case that  $U = V = W$  the composition of linear operators is defined on a single vector space  $L(V, V)$ , denoted simply as  $L(V)$ . The vector space  $L(V)$  is an algebra of linear operators. We start introducing the definition of an algebra over a field.

**Definition 5.3.1.** An *algebra*  $V$  over the field  $\mathbb{F}$  is a vector space  $V$  over the field  $\mathbb{F}$  equipped with an operation  $V \times V \rightarrow V$  called *multiplication*, satisfying,

$$\mathbf{u}(a\mathbf{v} + b\mathbf{w}) = a\mathbf{u}\mathbf{v} + b\mathbf{u}\mathbf{w} \quad \text{for all } \mathbf{u}, \mathbf{v}, \mathbf{w} \in V \text{ and } a, b \in \mathbb{F}.$$

The algebra is called *associative* iff the product satisfies

$$\mathbf{u}(\mathbf{v}\mathbf{w}) = (\mathbf{u}\mathbf{v})\mathbf{w} \quad \text{for all } \mathbf{u}, \mathbf{v}, \mathbf{w} \in V.$$

The algebra is called *commutative* iff the product satisfies

$$\mathbf{u}\mathbf{v} = \mathbf{v}\mathbf{u} \quad \text{for all } \mathbf{u}, \mathbf{v} \in V.$$

An algebra is called with *identity* iff there is an element  $\mathbf{i} \in V$  satisfying

$$\mathbf{i}\mathbf{u} = \mathbf{u}\mathbf{i} = \mathbf{u} \quad \text{for all } \mathbf{u} \in V.$$

Maybe the best known example of an algebra is the vector space  $\mathbb{R}^3$  equipped with the cross product, also called vector product. This is a non-associative, non-commutative algebra.

**EXAMPLE 5.3.1:** Let  $\mathcal{S} = (\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3)$  be the standard ordered basis of  $\mathbb{R}^3$  and introduce the cross product of vectors  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^3$  as follows

$$\mathbf{u} \times \mathbf{v} = (u_2v_3 - u_3v_2)\mathbf{e}_1 - (u_1v_3 - u_3v_1)\mathbf{e}_2 + (u_1v_2 - u_2v_1)\mathbf{e}_3,$$

where  $\mathbf{u} = u_1\mathbf{e}_1 + u_2\mathbf{e}_2 + u_3\mathbf{e}_3$  and  $\mathbf{v} = v_1\mathbf{e}_1 + v_2\mathbf{e}_2 + v_3\mathbf{e}_3$ . Show that  $\mathbb{R}^3$  equipped with the cross product is a non-associative, non-commutative algebra.

**SOLUTION:** It is convenient to express the cross product of two vectors using the determinant notation. Indeed, the cross product above is equal to the following expression

$$\mathbf{u} \times \mathbf{v} = \begin{vmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \mathbf{e}_3 \\ u_1 & u_2 & u_3 \\ v_1 & v_2 & v_3 \end{vmatrix},$$

where in the first row of the array above contains the basis vectors and the other two rows contain vector components. So, the array above is not a matrix, but the definition of determinant of a matrix when used on this array produces the cross product of two vectors. This notation is convenient since the cross product shares all the properties that the determinant has. For example, the cross product satisfies the following two properties:

$$\mathbf{u} \times (a\mathbf{v}) = \begin{vmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \mathbf{e}_3 \\ u_1 & u_2 & u_3 \\ av_1 & av_2 & av_3 \end{vmatrix} = a \begin{vmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \mathbf{e}_3 \\ u_1 & u_2 & u_3 \\ v_1 & v_2 & v_3 \end{vmatrix} = a(\mathbf{u} \times \mathbf{v});$$



$$\begin{aligned}
\mathbf{u} \times (\mathbf{v} + \mathbf{w}) &= \begin{vmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \mathbf{e}_3 \\ u_1 & u_2 & u_3 \\ (v_1 + w_1) & (v_2 + w_2) & (v_3 + w_3) \end{vmatrix} \\
&= \begin{vmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \mathbf{e}_3 \\ u_1 & u_2 & u_3 \\ v_1 & v_2 & v_3 \end{vmatrix} + \begin{vmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \mathbf{e}_3 \\ u_1 & u_2 & u_3 \\ w_1 & w_2 & w_3 \end{vmatrix} \\
&= \mathbf{u} \times \mathbf{v} + \mathbf{u} \times \mathbf{w}.
\end{aligned}$$

These two properties show that  $\mathbb{R}^3$  with the cross product form an algebra. This algebra is not commutative, since

$$\mathbf{u} \times \mathbf{v} = \begin{vmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \mathbf{e}_3 \\ u_1 & u_2 & u_3 \\ v_1 & v_2 & v_3 \end{vmatrix} = - \begin{vmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \mathbf{e}_3 \\ v_1 & v_2 & v_3 \\ u_1 & u_2 & u_3 \end{vmatrix} = -(\mathbf{v} \times \mathbf{u}).$$

In particular, notice that the equation above implies that  $\mathbf{u} \times \mathbf{u} = \mathbf{0}$ . Finally, this algebra is not associative as the following argument shows. First, from the definition of cross product it is simple to see that

$$\mathbf{e}_1 \times \mathbf{e}_2 = \mathbf{e}_3, \quad \mathbf{e}_3 \times \mathbf{e}_1 = \mathbf{e}_2, \quad \mathbf{e}_2 \times \mathbf{e}_3 = \mathbf{e}_1.$$

Using the first two equations we obtain the following relations,

$$\begin{aligned}
\mathbf{e}_1 \times (\mathbf{e}_1 \times \mathbf{e}_2) &= \mathbf{e}_1 \times \mathbf{e}_3 = -\mathbf{e}_2, \\
(\mathbf{e}_1 \times \mathbf{e}_1) \times \mathbf{e}_2 &= \mathbf{0} \times \mathbf{e}_2 = \mathbf{0},
\end{aligned}$$

showing that the algebra is not associative. ◁

In these notes we concentrate on only one example, the algebra of linear operators on a vector space, which we now describe in the following result.

**Theorem 5.3.2.** *The vector space  $L(V)$  of linear operators on a vector space  $V$  over the field  $\mathbb{F}$  equipped with the composition operation is an **associative algebra with identity**. This algebra is not commutative. Furthermore,  $\dim L(V) = (\dim V)^2$ .*

**Proof of Theorem 5.3.2:** The vector space  $L(V)$  equipped with the composition operation is an algebra. Indeed, given arbitrary linear operators  $\mathbf{R}, \mathbf{S}, \mathbf{T} \in L(V)$ , the following equations holds for all  $\mathbf{v} \in V$ , and all scalars all  $a, b \in \mathbb{F}$ ,

$$\begin{aligned}
(\mathbf{T} \circ (a\mathbf{S} + b\mathbf{R}))(\mathbf{v}) &= \mathbf{T}(a\mathbf{S}(\mathbf{v}) + b\mathbf{R}(\mathbf{v})) \\
&= a\mathbf{T}(\mathbf{S}(\mathbf{v})) + b\mathbf{T}(\mathbf{R}(\mathbf{v})) \\
&= a\mathbf{T} \circ \mathbf{S}(\mathbf{v}) + b\mathbf{T} \circ \mathbf{R}(\mathbf{v}).
\end{aligned}$$

Therefore, we established that the composition operation satisfies

$$\mathbf{T} \circ (a\mathbf{S} + b\mathbf{R}) = a\mathbf{T} \circ \mathbf{S} + b\mathbf{T} \circ \mathbf{R},$$

proving that the vector space  $L(V)$  equipped with the composition operation is and algebra. This algebra is associative, since for all linear operators  $\mathbf{R}, \mathbf{S}, \mathbf{T} \in L(V)$ , the following equations holds for all  $\mathbf{v} \in V$ ,

$$\begin{aligned}
(\mathbf{T} \circ (\mathbf{S} \circ \mathbf{R}))(\mathbf{v}) &= \mathbf{T}(\mathbf{S} \circ \mathbf{R})(\mathbf{v}) \\
&= \mathbf{T}(\mathbf{S}(\mathbf{R}(\mathbf{v}))) \\
&= (\mathbf{T} \circ \mathbf{S})(\mathbf{R}(\mathbf{v})) \\
&= ((\mathbf{T} \circ \mathbf{S}) \circ \mathbf{R})(\mathbf{v}),
\end{aligned}$$

so we conclude that

$$\mathbf{T} \circ (\mathbf{S} \circ \mathbf{R}) = (\mathbf{T} \circ \mathbf{S}) \circ \mathbf{R}.$$

This algebra has identity, since the identity transformation  $\mathbf{I}_V \in L(V)$  satisfies the equation  $\mathbf{I}_V \circ \mathbf{T} = \mathbf{T} = \mathbf{T} \circ \mathbf{I}_V$  for all  $\mathbf{T} \in L(V)$ . Finally, this algebra is not commutative. We just need to show one example with this property. We choose  $V = \mathbb{R}^2$  and the transformations given by matrices

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}.$$

We have seen in Example 2.1.3 and 2.1.4 that  $\mathbf{A}$  is a reflection along the axis  $x_1 = x_2$  while  $\mathbf{B}$  is a rotation by  $\pi/2$  counterclockwise. It is simple to see that

$$\begin{aligned} \mathbf{A} \circ \mathbf{B} &= \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \\ \mathbf{B} \circ \mathbf{A} &= \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}. \end{aligned}$$

Therefore, the algebra is not commutative. Finally, the result that  $\dim L(V) = n^2$  follows from Theorem 5.2.9. This establishes the Theorem.  $\square$

**EXAMPLE 5.3.2:** Show that the vector space  $\mathbb{F}^{n,n}$  equipped with the matrix multiplication operation is an associative algebra with identity.

**SOLUTION:** We know that  $\mathbb{F}^{n,n} = L(\mathbb{F}^n)$ , and we also know that matrix multiplication is equal to matrix composition, that is, for all  $\mathbf{A}, \mathbf{B} \in \mathbb{F}^{n,n}$  and all  $\mathbf{x} \in \mathbb{F}^n$  holds

$$(\mathbf{A}\mathbf{B})\mathbf{x} = \mathbf{A}(\mathbf{B}\mathbf{x}) = (\mathbf{A} \circ \mathbf{B})(\mathbf{x}) \quad \Rightarrow \quad \mathbf{A}\mathbf{B} = \mathbf{A} \circ \mathbf{B}.$$

Therefore, Theorem 5.3.2 implies that  $\mathbb{F}^{n,n}$  equipped with the matrix multiplication operation is an associative algebra with identity.  $\triangleleft$

**5.3.1. Polynomial functions of linear operators.** Polynomial functions of vectors in an algebra can be defined using linear combinations and multiplications of vectors. In the particular case of the algebra of linear operators, we use the notation  $\mathbf{T} \circ \mathbf{S} = \mathbf{TS}$  for the composition of linear operators in  $L(V)$ , and also  $\mathbf{T}^n = \mathbf{TT}^{n-1}$  for the  $n$ -th power of the operator  $\mathbf{T}$ , which is defined for all positive integers  $n \geq 1$ . The consistency of this equation for  $n = 1$  demands the definition  $\mathbf{T}^0 = \mathbf{I}_V$ . It follows that functions like the operator-valued polynomial below can be defined.

**Definition 5.3.3.** An *operator-valued polynomial* of degree  $n$  defined by the scalars  $a_0, \dots, a_n \in \mathbb{F}$  is the function  $p : L(V) \rightarrow L(V)$  given by

$$p(\mathbf{T}) = a_0 \mathbf{I}_V + a_1 \mathbf{T} + a_2 \mathbf{T}^2 + \dots + a_n \mathbf{T}^n.$$

**EXAMPLE 5.3.3:** Consider the vector space  $V = \mathbb{R}^2$  and given any real number  $\theta$  introduce the linear operator

$$\mathbf{R}_\theta = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}.$$

Find the explicit expression of the operator  $(\mathbf{R}_\theta)^n$  for any positive integer  $n$ .

**SOLUTION:** We have seen in Example 2.1.5 that the operator  $\mathbf{R}_\theta$  is a rotation by an angle  $\theta$  counterclockwise. We have also seen in Example 2.3.6 that given two operators  $\mathbf{R}_{\theta_1}$  and  $\mathbf{R}_{\theta_2}$ , the following formula holds,

$$\mathbf{R}_{\theta_1} \mathbf{R}_{\theta_2} = \mathbf{R}_{\theta_1 + \theta_2}.$$

In the particular case of  $\theta_1 = \theta_2$  we obtain  $(R_\theta)^2 = R_{2\theta}$ . Analogously,

$$(R_\theta)^3 = R_\theta(R_\theta)^2 = R_\theta R_{2\theta} = R_{3\theta}.$$

This calculation suggests the formula  $(R_\theta)^n = R_{n\theta}$  for all positive integer  $n$ . We now prove this formula using induction in  $n$ . Assume that the formula holds for an arbitrary value of  $n$  and then prove that the formula also holds for  $(n + 1)$ . Indeed,

$$(R_\theta)^{(n+1)} = R_\theta(R_\theta)^n = R_\theta R_{n\theta} = R_{(n+1)\theta}.$$

We then conclude that  $(R_\theta)^n = R_{n\theta}$  holds for every positive integer  $n$ . ◁

**5.3.2. Functions of linear operators.** Negative powers of an operator can be defined in the case that the operator is invertible. In that case, if one defines  $\mathbf{T}^{-n} = (\mathbf{T}^{-1})^n$  for any positive integer  $n$ , then the exponents satisfy the usual well known rules that hold when the base is a real number. That is, for an invertible operator  $\mathbf{T}$  holds that  $\mathbf{T}^m \mathbf{T}^n = \mathbf{T}^{(m+n)}$  and  $(\mathbf{T}^m)^n = \mathbf{T}^{mn}$ , for all integers  $m, n$ .

It is much more complicated to define further generalizations of these operator-valued power functions to include fractional exponents and real number exponents. The general idea behind such generalizations is the same used to define an infinitely differentiable function of an operator. Such functions are defined using Taylor expansions similar to those used on real-valued functions. Let us recall that the Taylor expansion centered at  $x_0 \in \mathbb{R}$  of an infinitely differentiable real-valued function is given by

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(x_0)}{n!} (x - x_0)^n = f(x_0) + f'(x_0)(x - x_0) + f''(x_0)(x - x_0)^2 + \dots .$$

**EXAMPLE 5.3.4:** Find the Taylor series expansion of the real valued exponential function  $f(x) = e^x$  centered  $x_0 = 0$ .

**SOLUTION:** Since  $x_0 = 0$ , the Taylor expansion formula above has the form

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(0)}{n!} x^n = f(0) + f'(0)x + f''(0)x^2 + \dots .$$

Since the exponential function  $f(x) = e^x$  satisfies that  $f^{(n)}(x) = f(x)$  for all positive integer  $n$ , and  $f(0) = e^0 = 1$ , we obtain the formula

$$e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots .$$

◁

Coming back to linear operators on a vector space, one can use the right hand side of the Taylor expansion formulas to define a function of a linear operator. However, such idea cannot be applied to linear operators until one understands the meaning of an infinite sum of linear operators. We discuss these ideas in Chapter 9. What we can do right now is to truncate the Taylor series and define a particular type of polynomial functions of linear operators as follows. Given the real-valued function  $f : \mathbb{R} \rightarrow \mathbb{R}$ , a vector space  $V$ , and a linear operator  $\mathbf{T} \in L(V)$ , introduce the Taylor polynomial  $f_N(\mathbf{T})$  for any positive integer  $N$  as follows

$$f_N(\mathbf{T}) = \sum_{n=0}^N \frac{f^{(n)}(x_0)}{n!} (\mathbf{T} - x_0 \mathbf{I}_V)^n,$$

that is,

$$f_N(\mathbf{T}) = f(x_0)\mathbf{I}_V + f'(x_0)(\mathbf{T} - x_0 \mathbf{I}_V) + \dots + \frac{f^{(N)}(x_0)}{N!} (\mathbf{T} - x_0 \mathbf{I}_V)^N.$$

The Taylor polynomial of degree  $N$  simplify in the case that  $x_0 = 0$ , that is

$$f_N(\mathbf{T}) = \sum_{n=0}^N \frac{f^{(n)}(0)}{n!} \mathbf{T}^n = f(0)\mathbf{I}_V + f'(0)\mathbf{T} + \cdots + \frac{f^{(N)}(0)}{N!} \mathbf{T}^N.$$

In order to define the limit of  $f_N(\mathbf{T})$  as  $N \rightarrow \infty$  it is needed a notion of distance between linear operators. We come back to this point in Chapter 9.

**5.3.3. The commutator of linear operators.** We know that the composition of linear maps depends on the order the operators appear. We say that function composition is not commutative. The commutator of two linear operators measures the lack of commutativity of these two operators.

**Definition 5.3.4.** The *commutator*,  $[\mathbf{T}, \mathbf{S}]$ , of the operators  $\mathbf{T}, \mathbf{S} \in L(V)$  is given by,

$$[\mathbf{T}, \mathbf{S}] = \mathbf{TS} - \mathbf{ST}.$$

In the case where  $L(V) = \mathbb{F}^{n,n}$  the commutator of linear operators defined above agrees with the definition of matrix commutators introduced at the end of Section 2.3. An immediate consequence of the definition above are the following properties.

**Theorem 5.3.5.** For  $\mathbf{S}, \mathbf{T}, \mathbf{U} \in L(V)$  and  $a, b \in \mathbb{F}$ , we have:

- |  |                                      |
|--|--------------------------------------|
| (a) $[\mathbf{T}, \mathbf{S}] = -[\mathbf{S}, \mathbf{T}]$ ,   | <i>antisymmetry;</i>                 |
| (b) $[a\mathbf{T}, b\mathbf{S}] = ab[\mathbf{T}, \mathbf{S}]$ ,  | <i>linearity;</i>                    |
| (c) $[\mathbf{U}, (\mathbf{T} + \mathbf{S})] = [\mathbf{U}, \mathbf{T}] + [\mathbf{U}, \mathbf{S}]$ ,                                | <i>linearity in the right entry;</i> |
| (d) $[(\mathbf{U} + \mathbf{T}), \mathbf{S}] = [\mathbf{U}, \mathbf{S}] + [\mathbf{T}, \mathbf{S}]$ ,                                | <i>linearity in the left entry;</i>  |
| (e) $[\mathbf{UT}, \mathbf{S}] = \mathbf{U}[\mathbf{T}, \mathbf{S}] + [\mathbf{U}, \mathbf{S}]\mathbf{T}$ ,                          | <i>left derivation property;</i>     |
| (f) $[\mathbf{U}, \mathbf{TS}] = \mathbf{T}[\mathbf{U}, \mathbf{S}] + [\mathbf{U}, \mathbf{T}]\mathbf{S}$ ,                          | <i>right derivation property;</i>    |
| (g) $[[\mathbf{S}, \mathbf{T}], \mathbf{U}] + [[\mathbf{U}, \mathbf{S}], \mathbf{T}] + [[\mathbf{T}, \mathbf{U}], \mathbf{S}] = 0$ , | <i>Jacobi identity.</i>              |

**Proof of Theorem 5.3.5:** Properties (a)-(d) and (g) follow from the definition of commutator in a straightforward way. We just show here the left derivation property:

$$\begin{aligned} [\mathbf{UT}, \mathbf{S}] &= \mathbf{UTS} - \mathbf{SUT} \\ &= \mathbf{UTS} - \mathbf{SUT} + \mathbf{UST} - \mathbf{UST} \\ &= (\mathbf{UTS} - \mathbf{UST}) + (\mathbf{UST} - \mathbf{SUT}) \\ &= \mathbf{U}[\mathbf{T}, \mathbf{S}] + [\mathbf{U}, \mathbf{S}]\mathbf{T}. \end{aligned}$$

The proof of the right derivation property is similar. This establishes the Theorem.  $\square$

The commutator of two linear operators in a vector space is an important concept in quantum mechanics, since it indicates how well two properties of the physical system described by these operators can be measured simultaneously. The Heisenberg uncertainty relations are statements about the commutators of linear operators.

## 5.3.4. Exercises.

5.3.1.- Let  $T, S \in L(\mathbb{F}^2)$  be the linear operators given by

$$T\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_2 \\ x_1 \end{bmatrix}, \quad S\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} 0 \\ x_1 \end{bmatrix}.$$

Compute explicitly the operators

- (a)  $3T - 2S$ ;
- (b)  $T \circ S$  and  $S \circ T$ ;
- (c)  $T^2$  and  $S^2$ .

5.3.2.- Find the dimension of the vector spaces  $L(\mathbb{F}^4)$ ,  $L(\mathbb{P}_2)$  and  $L(\mathbb{F}^{3,2})$ .

5.3.3.- Let  $T: \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be the linear operator

$$T\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} 2x_1 \\ 4x_1 - x_2 \end{bmatrix}.$$

Find the operator  $T^{-2}$ .

5.3.4.- Let  $T: \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be the linear operator

$$T\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1 + 2x_2 \\ 3x_1 + 4x_2 \end{bmatrix}.$$

Find the operator

$$p(T) = T^2 - 2T - 3I_2.$$

5.3.5.- Use the definition of the rotation matrix  $R_\theta$  given in Example 5.3.2 and the formula  $(R_\theta)^n = R_{n\theta}$  to explicitly verify that  $(R_\theta)^n (R_\theta)^{-n} = I_2$ .

5.3.6.- Prove the properties (a), (b) and (c) in Theorem 5.3.5.

## 5.4. TRANSFORMATION COMPONENTS

**5.4.1. The matrix of a linear transformation.** We have seen in Sections 2.1 and 5.1 that every  $m \times n$  matrix defines a linear transformation between the vector spaces  $\mathbb{F}^n$  and  $\mathbb{F}^m$  equipped with standard bases. We now show that every linear transformation  $T : V \rightarrow W$  between finite dimensional vector spaces  $V$  and  $W$  has associated a matrix  $\mathbf{T}$ . The matrix associated with such linear transformation is not unique. In order to compute the matrix of a linear transformation a basis must be chosen both in  $V$  and in  $W$ . The matrix associated with a linear transformation depends on the choice of bases done in  $V$  and  $W$ . We can say that the matrix  $\mathbf{T}$  of a linear transformation  $T : V \rightarrow W$  are the components of  $T$  in the given bases for  $V$  and  $W$ , analogously to the components  $\mathbf{v}$  of a vector  $\mathbf{v} \in V$  in a basis for  $V$  studied in Section 4.4.

**Definition 5.4.1.** Let both  $V$  and  $W$  be finite dimensional vector spaces over the field  $\mathbb{F}$  with respective ordered bases  $\mathcal{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$  and  $\mathcal{W} = (\mathbf{w}_1, \dots, \mathbf{w}_m)$ , and let the map  $[\ ]_w : W \rightarrow \mathbb{F}^m$  be the coordinate map on  $W$ . The **matrix of the linear transformation  $T : V \rightarrow W$**  in the ordered bases  $\mathcal{V}$  and  $\mathcal{W}$  is the  $m \times n$  matrix

$$\mathbf{T}_{vw} = [ [\mathbf{T}(\mathbf{v}_1)]_w, \dots, [\mathbf{T}(\mathbf{v}_n)]_w ]. \quad (5.3)$$

Let us explain the notation in Eq. (5.3). Given any basis vector  $\mathbf{v}_i \in \mathcal{V} \subset V$ , where  $i = 1, \dots, n$ , we know that  $\mathbf{T}(\mathbf{v}_i)$  is a vector in  $W$ . Like any vector in a vector space,  $\mathbf{T}(\mathbf{v}_i)$  can be decomposed in a unique way in terms of the basis vectors in  $\mathcal{W}$ , and following Sect. 4.4 we denote them as  $[\mathbf{T}(\mathbf{v}_i)]_w$ . When there is no possibility of confusion, we denote  $\mathbf{T}_{vw}$  simply as  $\mathbf{T}$ .

**EXAMPLE 5.4.1:** Consider the vector spaces  $V = W = \mathbb{R}^2$ , both with the standard basis, that is,  $\mathcal{S} = (\mathbf{e}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \mathbf{e}_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix})$ . Let the linear operator  $\mathbf{T} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be given by

$$\left[ \mathbf{T} \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) \right]_s = \begin{bmatrix} 3x_1 + 2x_2 \\ 4x_1 - x_2 \end{bmatrix}_s. \quad (5.4)$$

- (a) Find the matrix  $\mathbf{T}_{ss}$  associated with the linear operator above.  
 (b) Consider the ordered basis  $\mathcal{U}$  for  $\mathbb{R}^2$  given by  $\mathcal{U} = (\mathbf{u}_{1s} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \mathbf{u}_{2s} = \begin{bmatrix} -1 \\ 1 \end{bmatrix})$ . Find the matrix  $\mathbf{T}_{uu}$  associated with the linear operator above.

**SOLUTION:**

**Part (a):** The definition of  $\mathbf{T}_{ss}$  implies  $\mathbf{T}_{ss} = [[\mathbf{T}(\mathbf{e}_1)]_s, [\mathbf{T}(\mathbf{e}_2)]_s]$ . From Eq. (5.4) we know that

$$\left[ \mathbf{T} \left( \begin{bmatrix} 1 \\ 0 \end{bmatrix} \right) \right]_s = \begin{bmatrix} 3 \\ 4 \end{bmatrix}_s, \quad \left[ \mathbf{T} \left( \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right) \right]_s = \begin{bmatrix} 2 \\ -1 \end{bmatrix}_s.$$

Therefore,

$$\mathbf{T}_{ss} = \begin{bmatrix} 3 & 2 \\ 4 & -1 \end{bmatrix}.$$

**Part (b):** By definition  $\mathbf{T}_{uu} = [[\mathbf{T}(\mathbf{u}_1)]_u, [\mathbf{T}(\mathbf{u}_2)]_u]$ . Notice that from the definition of basis  $\mathcal{U}$  we have  $\mathbf{u}_{is} = [\mathbf{u}_i]_s$ , the column vector form with the components of the basis vectors  $\mathbf{u}_i$  in the standard basis  $\mathcal{S}$ . So, we can use the definition of  $\mathbf{T}$  in Eq. (5.4), that is,

$$[\mathbf{T}(\mathbf{u}_1)]_s = \left[ \mathbf{T} \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right) \right]_s = \begin{bmatrix} 5 \\ 3 \end{bmatrix}_s, \quad [\mathbf{T}(\mathbf{u}_2)]_s = \left[ \mathbf{T} \left( \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right) \right]_s = \begin{bmatrix} -1 \\ -5 \end{bmatrix}_s.$$

The results above are vector components in the standard basis  $\mathcal{S}$ , so we need to translate these results into components on the  $\mathcal{U}$  basis. We translate them in the usual way, solving a linear system for each of these two vectors:

$$\begin{bmatrix} 5 \\ 3 \end{bmatrix}_s = y_1 \begin{bmatrix} 1 \\ 1 \end{bmatrix}_s + y_2 \begin{bmatrix} -1 \\ 1 \end{bmatrix}_s, \quad \begin{bmatrix} -1 \\ -5 \end{bmatrix}_s = z_1 \begin{bmatrix} 1 \\ 1 \end{bmatrix}_s + z_2 \begin{bmatrix} -1 \\ 1 \end{bmatrix}_s.$$

The solutions are the components we are looking for, since

$$[\mathbf{T}(\mathbf{u}_1)]_u = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}_u, \quad [\mathbf{T}(\mathbf{u}_2)]_u = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}_u.$$

We can solve both systems above at the same time using the augmented matrix

$$\left[ \begin{array}{cc|cc} 1 & -1 & 5 & -1 \\ 1 & 1 & 3 & -5 \end{array} \right] \rightarrow \left[ \begin{array}{cc|cc} 1 & 0 & 4 & -3 \\ 0 & 1 & -1 & -2 \end{array} \right].$$

We have obtained that

$$[\mathbf{T}(\mathbf{u}_1)]_u = \begin{bmatrix} 4 \\ -1 \end{bmatrix}_u, \quad [\mathbf{T}(\mathbf{u}_2)]_u = \begin{bmatrix} -3 \\ -2 \end{bmatrix}_u.$$

So we conclude that

$$\mathsf{T}_{uu} = \begin{bmatrix} 4 & -3 \\ -1 & -2 \end{bmatrix}.$$

◁

From the Definition 5.4.1 it is clear that the matrix associated with a linear transformation  $\mathbf{T}: V \rightarrow W$  depends on the choice of bases  $\mathcal{V}$  and  $\mathcal{W}$  for the spaces  $V$  and  $W$ , respectively. In the case that  $V = W$  the matrix of the linear operator  $\mathbf{T}: V \rightarrow V$  usually means  $\mathsf{T}_{vv}$ , that is, to choose the same basis  $\mathcal{V}$  for the domain space  $V$  and the range space  $V$ . But this is not the only choice. One can choose different bases  $\mathcal{V}$  and  $\tilde{\mathcal{V}}$  for the domain and range spaces, respectively. In this case, the matrix associated with the linear operator  $\mathbf{T}$  is  $\mathsf{T}_{v\tilde{v}}$ , that is,

$$\mathsf{T}_{v\tilde{v}} = [[\mathbf{T}(\mathbf{v}_1)]_{\tilde{v}}, \dots, [\mathbf{T}(\mathbf{v}_n)]_{\tilde{v}}].$$

**EXAMPLE 5.4.2:** Consider the linear operator defined in Eq. (5.4) in Example 5.4.1. Find the associated matrices  $\mathsf{T}_{us}$  and  $\mathsf{T}_{su}$ , where  $\mathcal{S}$  and  $\mathcal{U}$  are the bases defined in that Example 5.4.1.

**SOLUTION:** The first matrix is simple to obtain, since

$$\mathsf{T}_{us} = [[\mathbf{T}(\mathbf{u}_1)]_s, [\mathbf{T}(\mathbf{u}_2)]_s]$$

and it is straightforward to compute

$$\left[ \mathbf{T} \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix}_s \right) \right]_s = \begin{bmatrix} 5 \\ 3 \end{bmatrix}_s, \quad \left[ \mathbf{T} \left( \begin{bmatrix} -1 \\ 1 \end{bmatrix}_s \right) \right]_s = \begin{bmatrix} -1 \\ -5 \end{bmatrix}_s,$$

so, we then conclude that

$$\mathsf{T}_{us} = \begin{bmatrix} 5 & -1 \\ 3 & -5 \end{bmatrix}_{us}.$$

We now compute

$$\mathsf{T}_{su} = [[\mathbf{T}(\mathbf{e}_1)]_u, [\mathbf{T}(\mathbf{e}_2)]_u].$$

From the definition of  $\mathbf{T}$  we know that

$$\left[ \mathbf{T} \left( \begin{bmatrix} 1 \\ 0 \end{bmatrix}_s \right) \right]_s = \begin{bmatrix} 3 \\ 4 \end{bmatrix}_s, \quad \left[ \mathbf{T} \left( \begin{bmatrix} 0 \\ 1 \end{bmatrix}_s \right) \right]_s = \begin{bmatrix} 2 \\ -1 \end{bmatrix}_s.$$

The results are expressed in the standard basis  $\mathcal{S}$ , so we need to translate them into the  $\mathcal{U}$  basis, as follows

$$\begin{bmatrix} 3 \\ 4 \end{bmatrix}_s = y_1 \begin{bmatrix} 1 \\ 1 \end{bmatrix}_s + y_2 \begin{bmatrix} -1 \\ 1 \end{bmatrix}_s, \quad \begin{bmatrix} 2 \\ -1 \end{bmatrix}_s = z_1 \begin{bmatrix} 1 \\ 1 \end{bmatrix}_s + z_2 \begin{bmatrix} -1 \\ 1 \end{bmatrix}_s.$$

The solutions are the components we are looking for, since

$$[\mathbf{T}(\mathbf{e}_1)]_u = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}_u, \quad [\mathbf{T}(\mathbf{e}_2)]_u = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}_u.$$

We can solve both systems above at the same time using the augmented matrix

$$\left[ \begin{array}{cc|cc} 1 & -1 & 3 & 2 \\ 1 & 1 & 4 & -1 \end{array} \right] \rightarrow \left[ \begin{array}{cc|cc} 2 & 0 & 7 & 1 \\ 0 & 2 & 1 & -3 \end{array} \right].$$

We have obtained that

$$[\mathbf{T}(\mathbf{e}_1)]_u = \frac{1}{2} \begin{bmatrix} 7 \\ 1 \end{bmatrix}_u, \quad [\mathbf{T}(\mathbf{e}_2)]_u = \frac{1}{2} \begin{bmatrix} 1 \\ -3 \end{bmatrix}_u.$$

So we conclude that

$$\mathsf{T}_{su} = \frac{1}{2} \begin{bmatrix} 7 & 1 \\ 1 & -3 \end{bmatrix}_{su}.$$

◁

**5.4.2. Action as matrix-vector product.** An important property of the matrix associated with a linear transformation  $\mathbf{T}$  is that the action of the transformation onto a vector can be represented as the matrix-vector product between the transformation matrix  $\mathsf{T}$  the vector components in the appropriate bases.

**Theorem 5.4.2.** *Let  $V$  and  $W$  be finite dimensional vector spaces with ordered bases  $\mathcal{V}$  and  $\mathcal{W}$ , respectively. Let  $\mathbf{T} : V \rightarrow W$  be a linear transformation with associated matrix  $\mathsf{T}_{vw}$ . Then, the components of the vector  $\mathbf{T}(\mathbf{x}) \in W$  in the basis  $\mathcal{W}$  can be expressed as the matrix-vector product*

$$[\mathbf{T}(\mathbf{x})]_w = \mathsf{T}_{vw} \mathbf{x}_v,$$

where  $\mathbf{x}_v = [\mathbf{x}]_v$  are the components of the vector  $\mathbf{x} \in V$  in the basis  $\mathcal{V}$ .

**Proof of Theorem 5.4.2:** Given any vector  $\mathbf{x} \in V$ , the definition of its vector components in the basis  $\mathcal{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  implies that

$$\mathbf{x}_v = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}_v \Leftrightarrow \mathbf{x} = x_1 \mathbf{v}_1 + \dots + x_n \mathbf{v}_n.$$

Since  $\mathbf{T}$  is a linear transformation we know that

$$\mathbf{T}(\mathbf{x}) = x_1 \mathbf{T}(\mathbf{v}_1) + \dots + x_n \mathbf{T}(\mathbf{v}_n).$$

The equation above holds in any basis of  $W$ , in particular in  $\mathcal{W}$ , that is,

$$\begin{aligned} [\mathbf{T}(\mathbf{x})]_w &= x_1 [\mathbf{T}(\mathbf{v}_1)]_w + \dots + x_n [\mathbf{T}(\mathbf{v}_n)]_w \\ &= \left[ [\mathbf{T}(\mathbf{v}_1)]_w, \dots, [\mathbf{T}(\mathbf{v}_n)]_w \right] \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}_v \\ &= \mathsf{T}_{vw} \mathbf{x}_v. \end{aligned}$$

This establishes the Theorem. □



**EXAMPLE 5.4.3:** Consider the linear operator  $\mathbf{T} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  given in Example 5.4.1 above together with the bases  $\mathcal{S}$  and  $\mathcal{U}$  defined in that example.

- (a) Use the matrix vector product to express the action of the operator  $T$  when the standard basis  $\mathcal{S}$  is used in both domain and range spaces  $\mathbb{R}^2$ .
- (b) Use the matrix vector product to express the action of the operator  $\mathbf{T}$  when the basis  $\mathcal{U}$  is used in both domain and range spaces  $\mathbb{R}^2$ .

**SOLUTION:**

**Part (a):** We need to find  $[\mathbf{T}(\mathbf{x})]_s$  and express the result using  $\mathbb{T}_{ss}$ . We use the notation  $\mathbf{x}_s = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_s$ , and we repeat the steps followed in the proof of Theorem 5.4.2.

$$[\mathbf{T}(\mathbf{x})]_s = [\mathbf{T}(x_1 \mathbf{e}_1 + x_2 \mathbf{e}_2)]_s = x_1 [\mathbf{T}(\mathbf{e}_1)]_s + x_2 [\mathbf{T}(\mathbf{e}_2)]_s = [[\mathbf{T}(\mathbf{e}_1)]_s, [\mathbf{T}(\mathbf{e}_2)]_s] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_s,$$

so we conclude that  $[\mathbf{T}(\mathbf{x})]_s = \mathbb{T}_{ss} \mathbf{x}_s$  and the matrix on the far right-hand side agrees with the matrix we found in the first half of Example 5.4.1, where we have found that

$$\mathbb{T}_{ss} = \begin{bmatrix} 3 & -2 \\ 4 & 1 \end{bmatrix}_{ss}.$$

Therefore, the action of  $\mathbf{T}$  on  $\mathbf{x}$  when expressed in the standard basis  $\mathcal{S}$  is given by the following matrix-vector product,

$$[\mathbf{T}(\mathbf{x})]_s = \begin{bmatrix} 3 & 2 \\ 4 & -1 \end{bmatrix}_{ss} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_s.$$

**Part (b):** We need to find  $[\mathbf{T}(\mathbf{x})]_u$  and express the result using  $\mathbb{T}_{uu}$ . We use the notation  $\mathbf{x}_u = \begin{bmatrix} \tilde{x}_1 \\ \tilde{x}_2 \end{bmatrix}_u$ , and we repeat the steps followed in the proof of Theorem 5.4.2.

$$[\mathbf{T}(\mathbf{x})]_u = [\mathbf{T}(\tilde{x}_1 \mathbf{u}_1 + \tilde{x}_2 \mathbf{u}_2)]_u = \tilde{x}_1 [\mathbf{T}(\mathbf{u}_1)]_u + \tilde{x}_2 [\mathbf{T}(\mathbf{u}_2)]_u = [[\mathbf{T}(\mathbf{u}_1)]_u, [\mathbf{T}(\mathbf{u}_2)]_u] \begin{bmatrix} \tilde{x}_1 \\ \tilde{x}_2 \end{bmatrix}_u,$$

so we conclude that  $[\mathbf{T}(\mathbf{x})]_u = \mathbb{T}_{uu} \mathbf{x}_u$ . At the end of Example 5.4.1 we have found that

$$\mathbb{T}_{uu} = \begin{bmatrix} 4 & -3 \\ -1 & -2 \end{bmatrix}_{uu}.$$

Therefore, the action of  $\mathbf{T}$  on  $\mathbf{x}$  when expressed in the basis  $\mathcal{U}$  is given by the following matrix-vector product,

$$[\mathbf{T}(\mathbf{x})]_u = \begin{bmatrix} 4 & -3 \\ -1 & -2 \end{bmatrix}_{uu} \begin{bmatrix} \tilde{x}_1 \\ \tilde{x}_2 \end{bmatrix}_u.$$

◁

**EXAMPLE 5.4.4:** Express the action of the differentiation transformation  $\mathbf{D} : \mathbb{P}_3 \rightarrow \mathbb{P}_2$ , given by  $\mathbf{D}(\mathbf{p})(x) = \frac{d\mathbf{p}}{dx}(x)$  as a matrix-vector product in the standard ordered bases

$$\begin{aligned} \mathcal{S} &= (\mathbf{p}_0 = 1, \mathbf{p}_1 = x, \mathbf{p}_2 = x^2, \mathbf{p}_3 = x^3) \subset \mathbb{P}_3, \\ \tilde{\mathcal{S}} &= (\mathbf{q}_0 = 1, \mathbf{q}_1 = x, \mathbf{q}_2 = x^2) \subset \mathbb{P}_2. \end{aligned}$$

**SOLUTION:** Following the Theorem 5.4.2 we only need to find the matrix  $D_{s\tilde{s}}$  associated with the transformation  $D$ ,

$$D_{s\tilde{s}} = [[D(\mathbf{p}_0)]_{\tilde{s}}, [D(\mathbf{p}_1)]_{\tilde{s}}, [D(\mathbf{p}_2)]_{\tilde{s}}, [D(\mathbf{p}_3)]_{\tilde{s}}] \Rightarrow D_{s\tilde{s}} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 3 \end{bmatrix}_{s\tilde{s}}.$$

Therefore, given any vector

$$\mathbf{p}(x) = a_0 + a_1x + a_2x^2 + a_3x^3 \Leftrightarrow \mathbf{p}_s = \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \end{bmatrix}_s$$

we obtain that  $D(\mathbf{p})(x) = a_1 + 2a_2x + 3a_3x^2$  is equivalent to

$$[D(\mathbf{p})]_{\tilde{s}} = D_{s\tilde{s}}\mathbf{p}_s \Rightarrow [D(\mathbf{p})]_{\tilde{s}} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 3 \end{bmatrix}_{s\tilde{s}} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \end{bmatrix}_s \Rightarrow [D(\mathbf{p})]_{\tilde{s}} = \begin{bmatrix} a_1 \\ 2a_2 \\ 3a_3 \end{bmatrix}_{\tilde{s}}.$$

◁

**EXAMPLE 5.4.5:** Express the action of the integration transformation  $\mathcal{S} : \mathbb{P}_2 \rightarrow \mathbb{P}_3$ , given by  $\mathcal{S}(\mathbf{q})(x) = \int_0^x \mathbf{q}(t) dt$  as a matrix-vector product in the standard ordered bases

$$\begin{aligned} \mathcal{S} &= (\mathbf{p}_0 = 1, \mathbf{p}_1 = x, \mathbf{p}_2 = x^2, \mathbf{p}_3 = x^3) \subset \mathbb{P}_3, \\ \tilde{\mathcal{S}} &= (\mathbf{q}_0 = 1, \mathbf{q}_1 = x, \mathbf{q}_2 = x^2) \subset \mathbb{P}_2. \end{aligned}$$

**SOLUTION:** Following the Theorem 5.4.2 we only need to find the matrix  $S_{\tilde{s}s}$  associated with the transformation  $\mathcal{S}$ ,

$$S_{\tilde{s}s} = [[\mathcal{S}(\mathbf{q}_0)]_s, [\mathcal{S}(\mathbf{q}_1)]_s, [\mathcal{S}(\mathbf{q}_2)]_s] \Rightarrow S_{\tilde{s}s} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & \frac{1}{2} & 0 \\ 0 & 0 & \frac{1}{3} \end{bmatrix}_{\tilde{s}s}.$$

Therefore, given any vector

$$\mathbf{q}(x) = a_0 + a_1x + a_2x^2 \Leftrightarrow \mathbf{q}_{\tilde{s}} = \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix}_{\tilde{s}}$$

we obtain that  $\mathcal{S}(\mathbf{q})(x) = a_0x + \frac{a_1}{2}x^2 + \frac{a_2}{3}x^3$  is equivalent to

$$[\mathcal{S}(\mathbf{q})]_s = S_{\tilde{s}s}\mathbf{q}_{\tilde{s}} \Rightarrow [\mathcal{S}(\mathbf{q})]_s = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & \frac{1}{2} & 0 \\ 0 & 0 & \frac{1}{3} \end{bmatrix}_{\tilde{s}s} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix}_{\tilde{s}} \Rightarrow [\mathcal{S}(\mathbf{q})]_s = \begin{bmatrix} 0 \\ a_0 \\ a_1/2 \\ a_2/3 \end{bmatrix}_s.$$

◁

**5.4.3. Composition and matrix product.** The following formula relates the composition of linear transformations with the matrix product of their associated matrices.

**Theorem 5.4.3.** *Let  $U$ ,  $V$  and  $W$  be finite-dimensional vector spaces with bases  $\mathcal{U}$ ,  $\mathcal{V}$  and  $\mathcal{W}$ , respectively. Let  $\mathbf{T} : U \rightarrow V$  and  $\mathbf{S} : V \rightarrow W$  be linear transformations with associated matrices  $\mathbf{T}_{uv}$  and  $\mathbf{S}_{vw}$ , respectively. Then, the composition  $\mathbf{S} \circ \mathbf{T} : U \rightarrow W$  given by  $(\mathbf{S} \circ \mathbf{T})(\mathbf{u}) = \mathbf{S}(\mathbf{T}(\mathbf{u}))$ , for all  $\mathbf{u} \in U$ , is a linear transformation and the associated matrix  $(\mathbf{S} \circ \mathbf{T})_{uw}$  is given by the matrix product by*

$$(\mathbf{S} \circ \mathbf{T})_{uw} = \mathbf{S}_{vw} \mathbf{T}_{uv}.$$

**Proof of Theorem 5.4.3:** First show that the composition of two linear transformations  $\mathbf{S}$  and  $\mathbf{T}$  is a linear transformation. Given any  $\mathbf{u}_1, \mathbf{u}_2 \in U$  and arbitrary scalars  $a$  and  $b$  holds

$$\begin{aligned} (\mathbf{S} \circ \mathbf{T})(a\mathbf{u}_1 + b\mathbf{u}_2) &= \mathbf{S}(\mathbf{T}(a\mathbf{u}_1 + b\mathbf{u}_2)) \\ &= \mathbf{S}(a\mathbf{T}(\mathbf{u}_1) + b\mathbf{T}(\mathbf{u}_2)) \\ &= a\mathbf{S}(\mathbf{T}(\mathbf{u}_1)) + b\mathbf{S}(\mathbf{T}(\mathbf{u}_2)) \\ &= a(\mathbf{S} \circ \mathbf{T})(\mathbf{u}_1) + b(\mathbf{S} \circ \mathbf{T})(\mathbf{u}_2). \end{aligned}$$

We now compute the matrix of the composition transformation. Denote the ordered basis in  $U$  as follows,  $\mathcal{U} = (\mathbf{u}_1, \dots, \mathbf{u}_n)$ . Then compute

$$(\mathbf{S} \circ \mathbf{T})_{uw} = [ [\mathbf{S}(\mathbf{T}(\mathbf{u}_1))]_w, \dots, [\mathbf{S}(\mathbf{T}(\mathbf{u}_n))]_w ].$$

The column  $i$ , with  $i = 1, \dots, n$ , in the matrix above has the form

$$[\mathbf{S}(\mathbf{T}(\mathbf{u}_i))]_w = \mathbf{S}_{vw} [\mathbf{T}(\mathbf{u}_i)]_v = \mathbf{S}_{vw} \mathbf{T}_{uv} \mathbf{u}_{iu}.$$

Therefore, we conclude that

$$(\mathbf{S} \circ \mathbf{T})_{uw} = \mathbf{S}_{vw} \mathbf{T}_{uv}.$$

This establishes the Theorem. □

**EXAMPLE 5.4.6:** Let  $\mathcal{S}$  be the standard ordered basis of  $\mathbb{R}^3$ , and consider the linear transformations  $\mathbf{T} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  and  $\mathbf{S} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  given by

$$\left[ \mathbf{T} \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_s \right) \right]_s = \begin{bmatrix} 2x_1 - x_2 + 3x_3 \\ -x_1 + 2x_2 - 4x_3 \\ x_2 + 3x_3 \end{bmatrix}_s, \quad \left[ \mathbf{S} \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_s \right) \right]_s = \begin{bmatrix} -x_1 \\ 2x_2 \\ 3x_3 \end{bmatrix}_s$$

- (a) Find a matrices  $\mathbf{T}_{ss}$  and  $\mathbf{S}_{ss}$ . By the way, Is  $\mathbf{T}$  injective? Is  $\mathbf{T}$  surjective?  
 (b) Find the matrix of the composition  $\mathbf{T} \circ \mathbf{S} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  in standard ordered basis  $\mathcal{S}$ .

**SOLUTION:** Since there is only one ordered basis in this problem, we represent  $\mathbf{T}_{ss}$  and  $\mathbf{S}_{ss}$  simply by  $\mathbf{T}$  and  $\mathbf{S}$ , respectively.

**Part (a):** A straightforward calculation from the definitions

$$\mathbf{T} = [[\mathbf{T}(\mathbf{e}_1)]_s, [\mathbf{T}(\mathbf{e}_2)]_s, [\mathbf{T}(\mathbf{e}_3)]_s], \quad \mathbf{S} = [[\mathbf{S}(\mathbf{e}_1)]_s, [\mathbf{S}(\mathbf{e}_2)]_s, [\mathbf{S}(\mathbf{e}_3)]_s],$$

gives us the matrices associated with  $\mathbf{T}$  and  $\mathbf{S}$  in the standard ordered basis,

$$\mathbf{T} = \begin{bmatrix} 2 & -1 & 3 \\ -1 & 2 & -4 \\ 0 & 1 & 3 \end{bmatrix}, \quad \mathbf{S} = \begin{bmatrix} -1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix}.$$

The information in matrix  $\mathbf{T}$  is useful to find out whether the linear transformation  $\mathbf{T}$  is injective and/or surjective. First, find the reduced echelon form of  $\mathbf{T}$ ,

$$\begin{bmatrix} 2 & -1 & 3 \\ -1 & 2 & -4 \\ 0 & 1 & 3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -2 & 4 \\ 2 & -1 & 3 \\ 0 & 1 & 3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -2 & 4 \\ 0 & 3 & 5 \\ 0 & 1 & 3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 10 \\ 0 & 1 & 3 \\ 0 & 0 & -4 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

We conclude that  $N(\mathbf{T}) = \{\mathbf{0}\}$ , which implies that  $\mathbf{T}$  is injective. The relation

$$\dim N(\mathbf{T}) + \dim \text{Col}(\mathbf{T}) = 3$$

together with  $\dim N(\mathbf{T}) = 0$  imply that  $\dim \text{Col}(\mathbf{T}) = 3$ , hence  $\mathbf{T}$  is surjective.

**Part (b):** The matrix of the composition  $\mathbf{T} \circ \mathbf{S}$  in the standard ordered basis  $\mathcal{S}$  is the product  $\mathbf{TS}$ , that is,

$$(\mathbf{T} \circ \mathbf{S})_{ss} = \mathbf{TS} = \begin{bmatrix} 2 & -1 & 3 \\ -1 & 2 & -4 \\ 0 & 1 & 3 \end{bmatrix} \begin{bmatrix} -1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix} \Rightarrow (\mathbf{T} \circ \mathbf{S})_{ss} = \begin{bmatrix} -2 & -2 & 9 \\ 1 & 4 & -12 \\ 0 & 2 & 9 \end{bmatrix}.$$

◁

## 5.4.4. Exercises.

5.4.1.- Consider  $\mathbf{T} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  the linear operator given by

$$\left[ \mathbf{T} \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_s \right) \right]_s = \begin{bmatrix} x_1 + x_2 \\ -2x_1 + 4x_2 \end{bmatrix}_s,$$

where  $\mathcal{S}$  denote the standard ordered basis of  $\mathbb{R}^2$ . Find the matrix  $\mathbf{T}_{uu}$  associated with the linear operator  $\mathbf{T}$  in and the ordered basis  $\mathcal{U}$  of  $\mathbb{R}^2$  given by

$$\left( [\mathbf{u}_1]_s = \begin{bmatrix} 1 \\ 1 \end{bmatrix}_s, [\mathbf{u}_2]_s = \begin{bmatrix} 1 \\ 2 \end{bmatrix}_s \right).$$

5.4.2.- Let  $\mathbf{T} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  be given by

$$\left[ \mathbf{T} \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_s \right) \right]_s = \begin{bmatrix} x_1 - x_2 \\ -x_1 + x_2 \\ x_1 - x_3 \end{bmatrix}_s,$$

$\mathcal{S}$  is the standard ordered basis of  $\mathbb{R}^3$ .

(a) Find the matrix  $\mathbf{T}_{uu}$  of the linear operator  $\mathbf{T}$  in the ordered basis

$$\mathcal{U} = \left( \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}_s, \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}_s, \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}_s \right).$$

(b) Verify that  $[\mathbf{T}(\mathbf{v})]_u = \mathbf{T}_{uu}\mathbf{v}_u$ , where

$$\mathbf{v}_s = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}_s.$$

5.4.3.- Fix an  $m \times n$  matrix  $\mathbf{A}$  and define the linear transformation  $\mathbf{T} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  as follows:  $[\mathbf{T}(\mathbf{x})]_{\tilde{\mathcal{S}}} = \mathbf{A}\mathbf{x}_s$ , where  $\mathcal{S} \subset \mathbb{R}^n$  and  $\tilde{\mathcal{S}} \subset \mathbb{R}^m$  are standard ordered bases. Show that  $\mathbf{T}_{\tilde{\mathcal{S}}\mathcal{S}} = \mathbf{A}$ .

5.4.4.- Find the matrix associated with the linear transformation  $\mathbf{T} : \mathbb{P}_3 \rightarrow \mathbb{P}_2$ ,

$$\mathbf{T}(\mathbf{p})(x) = \frac{d^2 \mathbf{p}}{dx^2}(x) - \frac{d\mathbf{p}}{dx}(x),$$

in the standard ordered bases of  $\mathbb{P}_3, \mathbb{P}_2$ .

5.4.5.- Find the matrices in the standard bases of  $\mathbb{P}_3$  and  $\mathbb{P}_2$  of the transformations

$$\mathbf{S} \circ \mathbf{D} : \mathbb{P}_3 \rightarrow \mathbb{P}_3, \quad \mathbf{D} \circ \mathbf{S} : \mathbb{P}_2 \rightarrow \mathbb{P}_2,$$

that is, the composition of the differentiation and integration transformations, as defined in Sect. 5.1.

## 5.5. CHANGE OF BASIS

In this Section we summarize the main calculations needed to find the change of vector and linear transformation components under a change of basis. We provide simple formulas to compute this change efficiently.

**5.5.1. Vector components.** We have seen in Sect. 4.4 that every vector  $\mathbf{v}$  in a finite dimensional vector space  $V$  with an ordered basis  $\mathcal{V}$ , can be expressed in a unique way as a linear combination of the basis elements, with  $\mathbf{v}_v = [\mathbf{v}]_v$  denoting the coefficients in that linear combination. These components depend on the basis chosen in  $V$ . Given two different ordered bases  $\mathcal{V}, \tilde{\mathcal{V}} \subset V$ , the components  $\mathbf{v}_v$  and  $\mathbf{v}_{\tilde{v}}$  associated with a vector  $\mathbf{v} \in V$  are in general different. We now use the matrix-vector product to find a simple formula relating  $\mathbf{v}_v$  to  $\mathbf{v}_{\tilde{v}}$ .

Let us recall the following notation. Given an  $n$ -dimensional vector space  $V$ , let  $\tilde{\mathcal{V}}$  and  $\mathcal{V}$  be two ordered bases of  $V$  given by

$$\tilde{\mathcal{V}} = (\tilde{\mathbf{v}}_1, \dots, \tilde{\mathbf{v}}_n) \quad \text{and} \quad \mathcal{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n).$$

Let  $\mathbf{I}: V \rightarrow V$  be the identity transformation, that is,  $\mathbf{I}(\mathbf{v}) = \mathbf{v}$  for all  $\mathbf{v} \in V$ , and introduce the *change of basis matrices*

$$\mathbf{l}_{\tilde{v}v} = [\tilde{\mathbf{v}}_{1v}, \dots, \tilde{\mathbf{v}}_{nv}] \quad \text{and} \quad \mathbf{l}_{v\tilde{v}} = [\mathbf{v}_{1\tilde{v}}, \dots, \mathbf{v}_{n\tilde{v}}],$$

where we denoted, as usual,  $\tilde{\mathbf{v}}_{iv} = [\tilde{\mathbf{v}}_i]_v$  and  $\mathbf{v}_{i\tilde{v}} = [\mathbf{v}_i]_{\tilde{v}}$ , for  $i = 1, \dots, n$ . Since the sets  $\mathcal{V}$  and  $\tilde{\mathcal{V}}$  are bases of  $V$ , the matrices  $\mathbf{l}_{v\tilde{v}}$  and  $\mathbf{l}_{\tilde{v}v}$  are invertible, and it is not difficult to show that  $(\mathbf{l}_{v\tilde{v}})^{-1} = \mathbf{l}_{\tilde{v}v}$ . Finally, introduce the following notation for the change of basis matrices,

$$\mathbf{P} = \mathbf{l}_{\tilde{v}v} \quad \text{and} \quad \mathbf{P}^{-1} = \mathbf{l}_{v\tilde{v}}.$$

**Theorem 5.5.1.** *Let  $V$  be a finite dimensional vector space, let  $\tilde{\mathcal{V}}$  and  $\mathcal{V}$  be two ordered bases of  $V$ , and let  $\mathbf{P} = \mathbf{l}_{\tilde{v}v}$  be the change of basis matrix. Then, the components  $\mathbf{x}_{\tilde{v}}$  and  $\mathbf{x}_v$  of any vector  $\mathbf{x} \in V$  in the ordered bases  $\tilde{\mathcal{V}}$  and  $\mathcal{V}$ , respectively, are related by the linear equation*

$$\mathbf{x}_{\tilde{v}} = \mathbf{P}^{-1}\mathbf{x}_v. \tag{5.5}$$

**REMARK:** Eq. (5.5) is equivalent to the inverse equation  $\mathbf{x}_v = \mathbf{P}\mathbf{x}_{\tilde{v}}$ .

**Proof of Theorem 5.5.1:** Let  $V$  be an  $n$ -dimensional vector space with two ordered bases  $\tilde{\mathcal{V}} = (\tilde{\mathbf{v}}_1, \dots, \tilde{\mathbf{v}}_n)$  and  $\mathcal{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$ . Given any vector  $\mathbf{x} \in V$ , then the definition of vector components  $\mathbf{x}_v = [\mathbf{x}]_v$  in the basis  $\mathcal{V}$  implies

$$\mathbf{x}_v = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}_v \quad \Leftrightarrow \quad \mathbf{x} = x_1\mathbf{v}_1 + \dots + x_n\mathbf{v}_n.$$

Express the second equation above in terms of components in the ordered basis  $\tilde{\mathcal{V}}$ ,

$$\mathbf{x}_{\tilde{v}} = x_1\mathbf{v}_{1\tilde{v}} + \dots + x_n\mathbf{v}_{n\tilde{v}} = [\mathbf{v}_{1\tilde{v}}, \dots, \mathbf{v}_{n\tilde{v}}] \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}_v \quad \Rightarrow \quad \mathbf{x}_{\tilde{v}} = \mathbf{l}_{v\tilde{v}}\mathbf{x}_v.$$

We conclude that  $\mathbf{x}_{\tilde{v}} = \mathbf{P}^{-1}\mathbf{x}_v$ . This establishes the Theorem.  $\square$

**EXAMPLE 5.5.1:** Consider the vector space  $V = \mathbb{R}^2$  with the standard ordered basis  $\mathcal{S}$  and the ordered basis  $\mathcal{U} = \left( \mathbf{u}_{1s} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}_s, \mathbf{u}_{2s} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}_s \right)$ . Given the vector with components  $\mathbf{x}_s = \begin{bmatrix} 1 \\ 3 \end{bmatrix}_s$ , find  $\mathbf{x}_u$ .

**SOLUTION:** The answer is given by Theorem 5.5.1, that says

$$\mathbf{x}_u = \mathbf{P}^{-1}\mathbf{x}_s, \quad \mathbf{P} = \mathbf{l}_{us}.$$

From the data of the problem the matrix  $\mathbf{P}$  is simple to compute, since

$$\mathbf{P} = \mathbf{l}_{us} = [\mathbf{u}_{1s}, \mathbf{u}_{2s}] = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}_{us}.$$

Computing the inverse matrix

$$\mathbf{P}^{-1} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}_{su}$$

we obtain the final result

$$\mathbf{x}_u = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}_{su} \begin{bmatrix} 1 \\ 3 \end{bmatrix}_s \Rightarrow \mathbf{x}_u = \begin{bmatrix} 2 \\ 1 \end{bmatrix}_u.$$

◁

**EXAMPLE 5.5.2:** Let  $V = \mathbb{R}^2$  with ordered bases  $\mathcal{B} = \{\mathbf{b}_1, \mathbf{b}_2\}$  and  $\mathcal{C} = \{\mathbf{c}_1, \mathbf{c}_2\}$  related by the equations

$$\mathbf{b}_1 = -\mathbf{c}_1 + 4\mathbf{c}_2, \quad \mathbf{b}_2 = 5\mathbf{c}_1 - 3\mathbf{c}_2.$$

- (a) Given  $\mathbf{x}_b = \begin{bmatrix} 5 \\ 3 \end{bmatrix}_b$ , find  $\mathbf{x}_c$ .  
 (b) Given  $\mathbf{x}_c = \begin{bmatrix} 1 \\ 1 \end{bmatrix}_c$ , find  $\mathbf{x}_b$ .

**SOLUTION:**

**Part (a):** We know that  $\mathbf{x}_c = \mathbf{P}^{-1}\mathbf{x}_b$ , where  $\mathbf{P} = \mathbf{l}_{cb}$ . From the data of the problem we know that

$$\left. \begin{array}{l} \mathbf{b}_1 = -\mathbf{c}_1 + 4\mathbf{c}_2 \Leftrightarrow \mathbf{b}_{1c} = \begin{bmatrix} -1 \\ 4 \end{bmatrix}_c, \\ \mathbf{b}_2 = 5\mathbf{c}_1 - 3\mathbf{c}_2 \Leftrightarrow \mathbf{b}_{2c} = \begin{bmatrix} 5 \\ -3 \end{bmatrix}_c, \end{array} \right\} \Rightarrow \mathbf{l}_{bc} = \begin{bmatrix} -1 & 5 \\ 4 & -3 \end{bmatrix}_{bc},$$

hence we know the change of basis matrix  $\mathbf{l}_{bc} = \mathbf{P}^{-1}$ , so we conclude that

$$\mathbf{x}_c = \begin{bmatrix} -1 & 5 \\ 4 & -3 \end{bmatrix}_{bc} \begin{bmatrix} 5 \\ 3 \end{bmatrix}_b \Rightarrow \mathbf{x}_c = \begin{bmatrix} 10 \\ 11 \end{bmatrix}_c$$

**Part (b):** keeping the definition of matrix  $\mathbf{P} = \mathbf{l}_{cb}$  as we introduced it in part (a), we know that  $\mathbf{x}_c = \mathbf{P}^{-1}\mathbf{x}_b$ . In this part (b) we need the inverse relation  $\mathbf{x}_b = \mathbf{P}\mathbf{x}_c$ . Since  $\mathbf{P}^{-1} = \begin{bmatrix} -1 & 5 \\ 4 & -3 \end{bmatrix}_{bc}$ , it is simple to obtain  $\mathbf{P} = \frac{1}{17} \begin{bmatrix} 3 & 5 \\ 4 & 1 \end{bmatrix}_{cb}$ . Using this matrix we obtain

$$\mathbf{x}_b = \frac{1}{17} \begin{bmatrix} 3 & 5 \\ 4 & 1 \end{bmatrix}_{cb} \begin{bmatrix} 1 \\ 1 \end{bmatrix}_c \Rightarrow \mathbf{x}_b = \frac{1}{17} \begin{bmatrix} 8 \\ 5 \end{bmatrix}_b.$$

◁

**5.5.2. Transformation components.** Analogously, we have seen in Sect. 5.4 that every linear transformation  $\mathbf{T}: V \rightarrow W$  between finite dimensional vector spaces  $V$  and  $W$  with ordered bases  $\mathcal{V}$  and  $\mathcal{W}$ , respectively, can be expressed in a unique way as an  $\dim W \times \dim V$  matrix  $\mathbf{T}_{vw}$ . These components depend on the bases chosen in  $V$  and  $W$ . Given two different ordered bases  $\mathcal{V}, \tilde{\mathcal{V}} \subset V$ , and given two different bases  $\mathcal{W}, \tilde{\mathcal{W}} \subset W$ , the matrices  $\mathbf{T}_{vw}$  and  $\mathbf{T}_{\tilde{v}\tilde{w}}$  associated with the linear transformation  $\mathbf{T}$  are in general different. Matrix multiplication provides a simple formula relating  $\mathbf{T}_{vw}$  and  $\mathbf{T}_{\tilde{v}\tilde{w}}$ .

Let us recall old notation and also introduce a bit of new one. Given an  $n$ -dimensional vector space  $V$ , let  $\tilde{\mathcal{V}}$  and  $\mathcal{V}$  be ordered bases of  $V$ ,

$$\tilde{\mathcal{V}} = (\tilde{\mathbf{v}}_1, \dots, \tilde{\mathbf{v}}_n) \quad \text{and} \quad \mathcal{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n);$$

and given an  $m$ -dimensional vector space  $W$ , let  $\tilde{\mathcal{W}}$  and  $\mathcal{W}$  be two ordered bases of  $W$ ,

$$\tilde{\mathcal{W}} = (\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_m) \quad \text{and} \quad \mathcal{W} = (\mathbf{w}_1, \dots, \mathbf{w}_m).$$

Let  $\mathbf{I}: V \rightarrow V$  be the identity transformation, that is,  $\mathbf{I}(\mathbf{v}) = \mathbf{v}$  for all  $\mathbf{v} \in V$ , and introduce the change of basis matrices

$$\mathbf{l}_{\tilde{v}\tilde{v}} = [\mathbf{v}_{1\tilde{v}}, \dots, \mathbf{v}_{n\tilde{v}}] \quad \text{and} \quad \mathbf{l}_{\tilde{v}v} = [\tilde{\mathbf{v}}_{1v}, \dots, \tilde{\mathbf{v}}_{nv}].$$

Let  $\mathbf{J}: W \rightarrow W$  be the identity transformation, that is,  $\mathbf{J}(\mathbf{w}) = \mathbf{w}$  for all  $\mathbf{w} \in W$ , and introduce the change of basis matrices

$$\mathbf{J}_{w\tilde{w}} = [\mathbf{w}_{1\tilde{w}}, \dots, \mathbf{w}_{m\tilde{w}}] \quad \text{and} \quad \mathbf{J}_{\tilde{w}w} = [\tilde{\mathbf{w}}_{1w}, \dots, \tilde{\mathbf{w}}_{mw}].$$

Notice that the sets  $\mathcal{V}$  and  $\tilde{\mathcal{V}}$  are bases of  $V$ , therefore the  $n \times n$  matrices  $\mathbf{l}_{v\tilde{v}}$  and  $\mathbf{l}_{\tilde{v}v}$  are invertible, and  $(\mathbf{l}_{v\tilde{v}})^{-1} = \mathbf{l}_{\tilde{v}v}$ . The similar statement is true for the  $m \times m$  matrices  $\mathbf{J}_{w\tilde{w}}$  and  $\mathbf{J}_{\tilde{w}w}$ , and  $(\mathbf{J}_{w\tilde{w}})^{-1} = \mathbf{J}_{\tilde{w}w}$ . Finally, introduce the following notation for the change of basis matrices,

$$\mathbf{P} = \mathbf{l}_{\tilde{v}v} \quad \Rightarrow \quad \mathbf{P}^{-1} = \mathbf{l}_{v\tilde{v}}, \quad \text{and} \quad \mathbf{Q} = \mathbf{J}_{\tilde{w}w} \quad \Rightarrow \quad \mathbf{Q}^{-1} = \mathbf{J}_{w\tilde{w}}.$$

**Theorem 5.5.2.** *Let  $V$  and  $W$  be finite dimensional vector spaces, let  $\tilde{\mathcal{V}}$  and  $\mathcal{V}$  be two ordered bases of  $V$ , let  $\tilde{\mathcal{W}}$  and  $\mathcal{W}$  be two ordered bases of  $W$ , and let  $\mathbf{P} = \mathbf{l}_{\tilde{v}v}$  and  $\mathbf{Q} = \mathbf{J}_{\tilde{w}w}$  be the change of basis matrices, respectively. Then, the components  $\mathbf{T}_{\tilde{v}\tilde{w}}$  and  $\mathbf{T}_{vw}$  of any linear transformation  $\mathbf{T}: V \rightarrow W$  in the bases  $\tilde{\mathcal{V}}, \tilde{\mathcal{W}}$  and  $\mathcal{V}, \mathcal{W}$ , respectively, are related by the matrix equation*

$$\mathbf{T}_{\tilde{v}\tilde{w}} = \mathbf{Q}^{-1} \mathbf{T}_{vw} \mathbf{P}. \quad (5.6)$$

**REMARK:** Eq. (5.6) is equivalent to the inverse equation  $\mathbf{T}_{vw} = \mathbf{Q} \mathbf{T}_{\tilde{v}\tilde{w}} \mathbf{P}^{-1}$ . A particular case of Theorem 5.5.2 that frequently appears in applications is when  $T$  is a linear operator. In this is the case  $\mathbf{T}: V \rightarrow V$ , so  $V = W$ . Often in applications one also has  $\mathcal{V} = \mathcal{W}$  and  $\tilde{\mathcal{V}} = \tilde{\mathcal{W}}$ , which imply  $\mathbf{P} = \mathbf{Q}$ .

**Corollary 5.5.3.** *Let  $V$  be a finite dimensional vector space, let  $\tilde{\mathcal{V}}$  and  $\mathcal{V}$  be two ordered bases of  $V$ , let  $\mathbf{P} = \mathbf{l}_{\tilde{v}v}$  be the change of basis matrix, and let  $\mathbf{T}: V \rightarrow V$  be a linear operator. Then, the components  $\tilde{\mathbf{T}} = \mathbf{T}_{\tilde{v}\tilde{v}}$  and  $\mathbf{T} = \mathbf{T}_{vv}$  of  $\mathbf{T}$  in the bases  $\tilde{\mathcal{V}}$  and  $\mathcal{V}$ , respectively, are related by the matrix equation*

$$\tilde{\mathbf{T}} = \mathbf{P}^{-1} \mathbf{T} \mathbf{P}. \quad (5.7)$$

**Proof of Theorem 5.5.2:** Let  $V$  be an  $n$ -dimensional vector space with ordered bases  $\tilde{\mathcal{V}} = (\tilde{\mathbf{v}}_1, \dots, \tilde{\mathbf{v}}_n)$  and  $\mathcal{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$ ; and let  $W$  be an  $m$ -dimensional vector space with ordered bases  $\tilde{\mathcal{W}} = (\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_m)$  and  $\mathcal{W} = (\mathbf{w}_1, \dots, \mathbf{w}_m)$ . We know that the matrix  $\mathbf{T}_{\tilde{v}\tilde{w}}$  associated with the transformation  $\mathbf{T}$  and the bases  $\tilde{\mathcal{V}}, \tilde{\mathcal{W}}$ , and the matrix  $\mathbf{T}_{vw}$  associated with the transformation  $\mathbf{T}$  and the bases  $\mathcal{V}, \mathcal{W}$  satisfy the following equations,

$$[\mathbf{T}(\mathbf{x})]_{\tilde{w}} = \mathbf{T}_{\tilde{v}\tilde{w}} \mathbf{x}_{\tilde{v}} \quad \text{and} \quad [\mathbf{T}(\mathbf{x})]_w = \mathbf{T}_{vw} \mathbf{x}_v. \quad (5.8)$$



Since  $\mathbf{T}(\mathbf{x}) \in W$ , Theorem 5.5.1 says that its components in the bases  $\tilde{\mathcal{W}}$  and  $\mathcal{W}$  are related by the equation

$$[\mathbf{T}(\mathbf{x})]_{\tilde{w}} = \mathbf{Q}^{-1}[\mathbf{T}(\mathbf{x})]_w, \quad \text{with } \mathbf{Q} = \mathbf{J}_{\tilde{w}w}.$$

Therefore, the first equation in (5.8) implies that

$$\begin{aligned} \mathbf{T}_{\tilde{v}\tilde{w}\times\tilde{v}} &= \mathbf{Q}^{-1}[\mathbf{T}(\mathbf{x})]_w \\ &= \mathbf{Q}^{-1}\mathbf{T}_{vw\times v} \\ &= \mathbf{Q}^{-1}\mathbf{T}_{vw}\mathbf{P}\times_{\tilde{v}}, \quad \text{with } \mathbf{P} = \mathbf{I}_{\tilde{v}v}, \end{aligned}$$

where in the second line above we used the second equation in (5.8), and in the third line above we used the inverse form of Theorem 5.5.1. Since the equation above holds for all  $\mathbf{x} \in V$  we conclude that

$$\mathbf{T}_{\tilde{v}\tilde{w}} = \mathbf{Q}^{-1}\mathbf{T}_{vw}\mathbf{P}.$$

This establishes the Theorem.  $\square$

**EXAMPLE 5.5.3:** Let  $\mathcal{S}$  be the standard ordered basis of  $\mathbb{R}^2$ , and let  $\mathbf{T} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be the linear operator given by

$$\left[ \mathbf{T} \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) \right]_{\mathcal{S}} = \begin{bmatrix} x_2 \\ x_1 \end{bmatrix}_{\mathcal{S}},$$

that is, a reflection along the line  $x_2 = x_1$ . Find the matrix  $\mathbf{T}_{uu}$ , where the ordered basis  $\mathcal{U}$  is given by  $\mathcal{U} = \left( \mathbf{u}_{1s} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}_{\mathcal{S}}, \mathbf{u}_{2s} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}_{\mathcal{S}} \right)$ .

**SOLUTION:** From the definition of  $\mathbf{T}$  is straightforward to obtain  $\mathbf{T}_{ss}$ , as follows,

$$\mathbf{T}_{ss} = \left[ \left[ \mathbf{T} \left( \begin{bmatrix} 1 \\ 0 \end{bmatrix} \right) \right]_{\mathcal{S}}, \left[ \mathbf{T} \left( \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right) \right]_{\mathcal{S}} \right] \Rightarrow \mathbf{T}_{ss} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}_{ss}.$$

Since  $\mathbf{T}$  is a linear operator and we want to compute the matrix  $\mathbf{T}_{uu}$  from matrix  $\mathbf{T}_{ss}$ , we can use Corollary 5.5.3, which says that these matrices are related by the similarity transformation

$$\mathbf{T}_{uu} = \mathbf{P}^{-1}\mathbf{T}_{ss}\mathbf{P}, \quad \text{where } \mathbf{P} = \mathbf{I}_{us}.$$

From the data of the problem we know that

$$\mathbf{I}_{us} = [\mathbf{u}_{1s}, \mathbf{u}_{2s}]_{us} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}_{us},$$

therefore,

$$\mathbf{P} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}_{us}, \quad \mathbf{P}^{-1} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}_{su}.$$

We then conclude that

$$\mathbf{T}_{uu} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}_{su} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}_{ss} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}_{us} \Rightarrow \mathbf{T}_{uu} = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}_{uu}.$$

**REMARK:** We can see in this Example that the matrix associated with the reflection transformation  $T$  is diagonal in the basis  $\mathcal{U}$ . For this particular transformation we have that

$$\mathbf{T}(\mathbf{u}_1) = \mathbf{u}_1, \quad \mathbf{T}(\mathbf{u}_2) = -\mathbf{u}_2.$$

Non-zero vectors  $\mathbf{v}$  with this property, that  $\mathbf{T}(\mathbf{v}) = \lambda\mathbf{v}$ , are called eigenvectors of the operator  $\mathbf{T}$  and the scalar  $\lambda$  is called eigenvalue of  $\mathbf{T}$ . In this example the elements in the basis  $\mathcal{U}$  are eigenvectors of the reflection operator, and the matrix of  $\mathbf{T}$  in this special basis is diagonal. Basis formed with eigenvectors of a given linear operator will be studied in a later on.  $\triangleleft$

**EXAMPLE 5.5.4:** Let  $\mathcal{S}$  be the standard ordered basis of  $\mathbb{R}^2$ , and let  $\mathbf{T} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be the linear operator given by

$$\left[ \mathbf{T} \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_s \right) \right]_s = \frac{(x_1 + x_2)}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix}_s,$$

that is, a projection along the line  $x_2 = x_1$ . Find the matrix  $\mathsf{T}_{uu}$ , where the ordered basis  $\mathcal{U} = \left( \mathbf{u}_{1s} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}_s, \mathbf{u}_{2s} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}_s \right)$ .

**SOLUTION:** From the definition of  $\mathbf{T}$  is straightforward to obtain  $\mathsf{T}_{ss}$ , as follows,

$$\mathsf{T}_{ss} = \left[ \left[ \mathbf{T} \left( \begin{bmatrix} 1 \\ 0 \end{bmatrix}_s \right) \right]_s, \left[ \mathbf{T} \left( \begin{bmatrix} 0 \\ 1 \end{bmatrix}_s \right) \right]_s \right] \Rightarrow \mathsf{T}_{ss} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}_{ss}.$$

Since  $\mathbf{T}$  is a linear operator and we want to compute the matrix  $\mathsf{T}_{uu}$  from matrix  $\mathsf{T}_{ss}$ , we again use Corollary 5.5.3, which says that these matrices are related by the similarity transformation

$$\mathsf{T}_{uu} = \mathsf{P}^{-1} \mathsf{T}_{ss} \mathsf{P}, \quad \text{where } \mathsf{P} = \mathbf{l}_{us}.$$

In Example 5.5.3 we have already computed

$$\mathsf{P} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}_{us}, \quad \mathsf{P}^{-1} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}_{su}.$$

We then conclude that

$$\mathsf{T}_{uu} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}_{su} \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}_{ss} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}_{us} \Rightarrow \mathsf{T}_{uu} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}_{uu}.$$

**REMARK:** Again in this Example we can see that the matrix associated with the reflection transformation  $\mathbf{T}$  is diagonal in the basis  $\mathcal{U}$ , with diagonal elements equal to one and zero. For this particular transformation we have that the basis vectors in  $\mathcal{U}$  are eigenvectors of  $\mathbf{T}$  with eigenvalues 1 and 0, that is,  $\mathbf{T}(\mathbf{u}_1) = \mathbf{u}_1$  and  $\mathbf{T}(\mathbf{u}_2) = \mathbf{0}$ .  $\triangleleft$

**EXAMPLE 5.5.5:** Let  $\mathcal{U}$  and  $\mathcal{V}$  be standard ordered bases of  $\mathbb{R}^3$  and  $\mathbb{R}^2$ , respectively, let  $\mathbf{T} : \mathbb{R}^3 \rightarrow \mathbb{R}^2$  be the linear transformation

$$\left[ \mathbf{T} \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_u \right) \right]_v = \begin{bmatrix} x_1 - x_2 + x_3 \\ x_2 - x_3 \end{bmatrix}_v,$$

and introduce the ordered bases

$$\begin{aligned} \tilde{\mathcal{U}} &= \left( \tilde{\mathbf{u}}_{1u} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}_u, \tilde{\mathbf{u}}_{2u} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}_u, \tilde{\mathbf{u}}_{3u} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}_u \right) \subset \mathbb{R}^3, \\ \tilde{\mathcal{V}} &= \left( \tilde{\mathbf{v}}_{1v} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}_v, \tilde{\mathbf{v}}_{2v} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}_v \right) \subset \mathbb{R}^2. \end{aligned}$$

Find the matrices  $\mathsf{T}_{uv}$  and  $\mathsf{T}_{\tilde{u}\tilde{v}}$ .

**SOLUTION:** We start finding the matrix  $\mathsf{T}_{uv}$ , which by definition is given by

$$\mathsf{T}_{uv} = \left[ \left[ \mathbf{T} \left( \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}_u \right) \right]_v, \left[ \mathbf{T} \left( \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}_u \right) \right]_v, \left[ \mathbf{T} \left( \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}_u \right) \right]_v \right]$$

hence we obtain  $T_{uv} = \begin{bmatrix} 1 & -1 & 1 \\ 0 & 1 & -1 \end{bmatrix}_{uv}$ . Theorem 5.5.2 says that the matrices  $T_{\tilde{u}\tilde{v}}$  and  $T_{uv}$  are related by the equation

$$T_{\tilde{u}\tilde{v}} = Q^{-1}T_{uv}P, \quad \text{where } Q = J_{\tilde{v}v} \quad \text{and} \quad P = I_{\tilde{u}u}.$$

From the data of the problem we know that

$$I_{\tilde{u}u} = [\tilde{u}_{1u}, \tilde{u}_{2u}, \tilde{u}_{3u}]_{\tilde{u}u} = \begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}_{\tilde{u}u} \Rightarrow P = \begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}_{\tilde{u}u},$$

$$J_{\tilde{v}v} = [\tilde{v}_{1v}, \tilde{v}_{2v}]_{\tilde{v}v} = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}_{\tilde{v}v} \Rightarrow Q = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}_{\tilde{v}v}, \quad Q^{-1} = \frac{1}{3} \begin{bmatrix} -1 & 2 \\ 2 & -1 \end{bmatrix}_{\tilde{v}v}.$$

Therefore, we need to compute the matrix product

$$T_{\tilde{u}\tilde{v}} = \frac{1}{3} \begin{bmatrix} -1 & 2 \\ 2 & -1 \end{bmatrix}_{\tilde{v}v} \begin{bmatrix} 1 & -1 & 1 \\ 0 & 1 & -1 \end{bmatrix}_{uv} \begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}_{\tilde{u}u}$$

and the result is  $T_{\tilde{u}\tilde{v}} = \frac{1}{3} \begin{bmatrix} 2 & 0 & -4 \\ -1 & 0 & 5 \end{bmatrix}_{\tilde{u}\tilde{v}}$ . ◁

**5.5.3. Determinant and trace of linear operators.** The type of matrix transformation given by Eq. (5.7) will be important later on, so we give such transformations a name.

**Definition 5.5.4.** The  $n \times n$  matrices  $A$  and  $B$  are related by a *similarity transformation* iff there exists an invertible  $n \times n$  matrix  $P$  such that

$$B = P^{-1}AP.$$

Therefore, similarity transformations are transformations among matrices. They appear in two different contexts. The first, or *active* context, is when matrix  $A$  and matrix  $B$  above are written in the same basis. In this case a similarity transformation changes matrix  $A$  into matrix  $B$ . The second, or *passive* context, is when matrix  $A$  and matrix  $B$  are two matrices corresponding to the same linear transformation. In this case matrices  $A$  and  $B$  are written in different basis, and the similarity transformation is the change of basis equation for the linear transformation matrices.

The following results says that, no matter in which context one uses a similarity transformation, the determinant and trace of a matrix are invariant under similarity transformations.

**Theorem 5.5.5.** If the square matrices  $A$ ,  $B$  are related by the similarity transformation  $B = P^{-1}AP$ , then holds

$$\det(B) = \det(A), \quad \text{tr}(B) = \text{tr}(A).$$

**Proof of Theorem: 5.5.5:** The determinant is invariant under similarity transformations, since

$$\det(B) = \det(P^{-1}AP) = \det(P^{-1})\det(A)\det(P) = \det(A).$$

The trace is also invariant under similarity transformations, since

$$\text{tr}(B) = \text{tr}(P^{-1}AP) = \text{tr}(PP^{-1}A) = \text{tr}(A).$$

This establishes the Theorem. □

We now use this result in the passive context for similarity transformations. This means that the determinant and trace are independent of the vector basis one uses to write down the matrix. Therefore, a determinant or a trace is an operation on a linear transformation, not on a particular matrix representation of a linear transformation. This observation suggests the following definition.

**Definition 5.5.6.** *Let  $V$  be a finite dimensional vector space, let  $\mathbf{T} \in L(V)$  be a linear operator, and let  $\mathbf{T}_{\mathcal{V}\mathcal{V}}$  be the matrix of the linear operator in an arbitrary ordered basis  $\mathcal{V}$  of  $V$ . The determinant and trace of a linear operator  $\mathbf{T} \in L(V)$  are respectively given by,*

$$\det(\mathbf{T}) = \det(\mathbf{T}_{\mathcal{V}\mathcal{V}}), \quad \text{tr}(\mathbf{T}) = \text{tr}(\mathbf{T}_{\mathcal{V}\mathcal{V}}).$$

We repeat that the determinant and trace of a linear operator are well-defined, since given any other ordered basis  $\mathcal{V}$  of  $V$ , we know that the matrices  $\mathbf{T}_{\bar{\mathcal{V}}\bar{\mathcal{V}}}$  and  $\mathbf{T}_{\mathcal{V}\mathcal{V}}$  are related by a similarity transformation. However, the determinant and trace are invariant under similarity transformations. So no matter what vector basis we use to compute determinant and trace of a linear operator, we always get the same result.

## 5.5.4. Exercises.

5.5.1.- Let  $\mathcal{U} = (\mathbf{u}_1, \mathbf{u}_2)$  be an ordered basis of  $\mathbb{R}^2$  given by

$$\mathbf{u}_1 = 2\mathbf{e}_1 - 9\mathbf{e}_2, \quad \mathbf{u}_2 = \mathbf{e}_1 + 8\mathbf{e}_2,$$

where  $\mathcal{S} = (\mathbf{e}_1, \mathbf{e}_2)$  is the standard ordered basis of  $\mathbb{R}^2$ .

(a) Find both change of basis matrices  $\mathbf{l}_{us}$  and  $\mathbf{l}_{su}$ .

(b) Given the vector  $\mathbf{x} = 2\mathbf{u}_1 + \mathbf{u}_2$ , find both  $\mathbf{x}_s$  and  $\mathbf{x}_u$ .

5.5.2.- Consider the ordered bases of  $\mathbb{R}^3$ ,  $\mathcal{B} = (\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3)$  and  $\mathcal{C} = (\mathbf{c}_1, \mathbf{c}_2, \mathbf{c}_3)$ , where

$$\mathbf{c}_1 = \mathbf{b}_1 - 2\mathbf{b}_2 + \mathbf{b}_3,$$

$$\mathbf{c}_2 = -\mathbf{b}_2 + 3\mathbf{b}_3,$$

$$\mathbf{c}_3 = -2\mathbf{b}_1 + \mathbf{b}_3.$$

(a) Find both the change of basis matrices  $\mathbf{l}_{bc}$  and  $\mathbf{l}_{cb}$ .

(b) Let  $\mathbf{x} = \mathbf{c}_1 - 2\mathbf{c}_2 + 2\mathbf{c}_3$ . Find both  $\mathbf{x}_b$  and  $\mathbf{x}_c$ .

5.5.3.- Consider the ordered bases  $\mathcal{U}$  and  $\mathcal{S}$  of  $\mathbb{R}^2$  given by

$$\mathcal{U} = \left( \mathbf{u}_{1s} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}_s, \mathbf{u}_{2s} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}_s \right),$$

$$\mathcal{S} = \left( \mathbf{e}_{1s} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}_s, \mathbf{e}_{2s} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}_s \right).$$

(a) Given  $\mathbf{x}_u = \begin{bmatrix} 3 \\ 2 \end{bmatrix}_u$  find  $\mathbf{x}_s$ .

(b) Find  $\mathbf{e}_{1u}$  and  $\mathbf{e}_{2u}$ .

5.5.4.- Consider the ordered bases of  $\mathbb{R}^2$

$$\mathcal{B} = \left( \mathbf{b}_{1s} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}_s, \mathbf{b}_{2s} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}_s \right)$$

$$\mathcal{C} = \left( \mathbf{c}_{1s} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}_s, \mathbf{c}_{2s} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}_s \right),$$

where  $\mathcal{S}$  is the standard ordered basis.

(a) Given  $\mathbf{x}_c = \begin{bmatrix} 2 \\ 3 \end{bmatrix}_c$ , find  $\mathbf{x}_s$ .

(b) For the same  $\mathbf{x}$  above, find  $\mathbf{x}_b$ .

5.5.5.- Let  $\mathcal{S} = (\mathbf{e}_1, \mathbf{e}_2)$  be the standard ordered basis of  $\mathbb{R}^2$  and  $\mathcal{B} = (\mathbf{b}_1, \mathbf{b}_2)$  be another ordered basis. Let  $\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix}$

be the matrix that transforms the components of a vector  $\mathbf{x} \in \mathbb{R}^2$  from the basis  $\mathcal{S}$  into the basis  $\mathcal{B}$ , that is,  $\mathbf{x}_b = \mathbf{A}\mathbf{x}_s$ . Find the components of the basis vectors  $\mathbf{b}_1, \mathbf{b}_2$  in the standard basis, that is, find  $\mathbf{b}_{1s}$  and  $\mathbf{b}_{2s}$ .

5.5.6.- Show that similarity is a transitive property, that is, if matrix  $\mathbf{A}$  is similar to matrix  $\mathbf{B}$ , and  $\mathbf{B}$  is similar to matrix  $\mathbf{C}$ , then  $\mathbf{A}$  is similar to  $\mathbf{C}$ .

5.5.7.- Consider  $\mathbb{R}^3$  with the standard ordered basis  $\mathcal{S}$  and the ordered basis

$$\mathcal{U} = \left( \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}_s, \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}_s, \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}_s \right).$$

Let  $\mathbf{T}: \mathbb{R}^3 \rightarrow \mathbb{R}^3$  be the linear operator

$$\left[ \mathbf{T} \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_s \right) \right]_s = \begin{bmatrix} x_1 + 2x_2 - x_3 \\ -x_2 \\ x_1 + 7x_3 \end{bmatrix}_s.$$

Find both matrices  $\mathbf{T}_{ss}$  and  $\mathbf{T}_{uu}$ .

5.5.8.- Consider  $\mathbb{R}^2$  with ordered bases

$$\mathcal{S} = \left( \mathbf{e}_{1s} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}_s, \mathbf{e}_{2s} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}_s \right),$$

$$\mathcal{U} = \left( \mathbf{u}_{1s} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}_s, \mathbf{u}_{2s} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}_s \right).$$

Let  $\mathbf{T}: \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be a linear transformation given by

$$[\mathbf{T}(\mathbf{u}_1)]_s = \begin{bmatrix} 1 \\ 3 \end{bmatrix}_s, \quad [\mathbf{T}(\mathbf{u}_2)]_s = \begin{bmatrix} 3 \\ 1 \end{bmatrix}_s.$$

Find the matrix  $\mathbf{T}_{us}$ , then the matrices  $\mathbf{T}_{ss}$ ,  $\mathbf{T}_{uu}$ , and finally  $\mathbf{T}_{su}$ .

## CHAPTER 6. INNER PRODUCT SPACES

An inner product space is a vector space with an additional structure called inner product. This additional structure is an operation that associates each pair of vectors in the vector space with a scalar. An inner product extends to any vector space the main concepts included in the dot product, which is defined on  $\mathbb{R}^n$ . These main concepts include the length of a vector, the notion of perpendicular vectors, and distance between vectors. When these ideas are introduced in function vector spaces, they allow to define the notion of convergence of an infinite sum of vectors. This, in turns, provides a way to evaluate the accuracy of approximate solutions to differential equations.

## 6.1. DOT PRODUCT

**6.1.1. Dot product in  $\mathbb{R}^2$ .** We review the definition of the dot product between vectors in  $\mathbb{R}^2$ , and we describe its main properties, including the Cauchy-Schwarz inequality. We then use the dot product to introduce the notion of length of a vector, distance and angle between vectors, including the special case of perpendicular vectors. We then review that all these notions can be generalized in a straightforward way from  $\mathbb{R}^2$  to  $\mathbb{F}^n$ ,  $n \geq 1$ .

**Definition 6.1.1.** Given any vectors  $x, y \in \mathbb{R}^2$  with components  $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ ,  $y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$  in the standard ordered basis  $\mathcal{S}$ . The **dot product** on  $\mathbb{R}^2$  with is the function  $\cdot : \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ ,

$$x \cdot y = x_1 y_1 + x_2 y_2.$$

The **dot product norm** of a vector  $x \in \mathbb{R}^2$  is the value of the function  $\| \cdot \| : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,

$$\|x\| = \sqrt{x \cdot x}.$$

The **norm distance** between  $x, y \in \mathbb{R}^2$  is the value of the function  $d : \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ ,

$$d(x, y) = \|x - y\|.$$

The dot product can be expressed using the transpose of a vector components in the standard basis, as follows,

$$x^T y = [x_1, x_2] \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = x_1 y_1 + x_2 y_2 = x \cdot y.$$

The dot product norm and the norm distance can be expressed in term of vector components in the standard ordered basis  $\mathcal{S}$  as follows,

$$\|x\| = \sqrt{(x_1)^2 + (x_2)^2}, \quad d(x, y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}.$$

The geometrical meaning of the norm and distance is clear from this expression in components, as is shown in Fig. 40. The norm of a vector is the Euclidean length from the origin point to the head point of the vector, while the distance between two vectors in the Euclidean distance between the head points of the two vectors.

It is important that we summarize the main properties of the dot product in  $\mathbb{R}^2$ , since they are the main guide to construct the generalizations of the dot product to other vector spaces.

**Theorem 6.1.2.** The dot product on  $\mathbb{R}^2$  satisfies, for every vector  $x, y, z \in \mathbb{R}^2$  and every scalar  $a, b \in \mathbb{R}$ , the following properties:

- (a)  $x \cdot y = y \cdot x$ , (Symmetry);
- (b)  $x \cdot (ay + bz) = a(x \cdot y) + b(x \cdot z)$ , (Linearity on the second argument);
- (c)  $x \cdot x \geq 0$ , and  $x \cdot x = 0$  iff  $x = 0$ , (Positive definiteness).

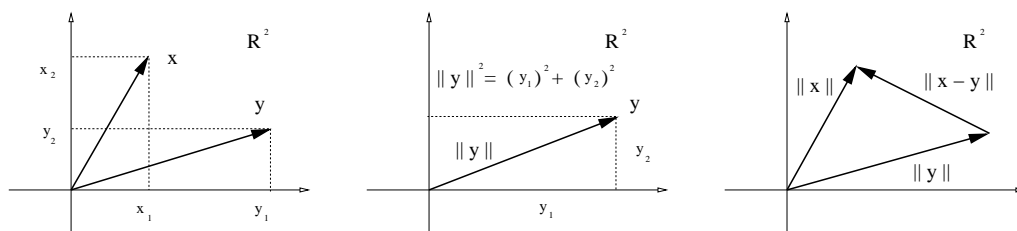


FIGURE 40. Example of the Euclidean notions of vector length and distance between vectors in  $\mathbb{R}^2$ .

**Proof of Theorem 6.1.2:** These properties are simple to obtain from the definition of the dot product.

**Part (a):** It is simple to see that

$$\mathbf{x} \cdot \mathbf{y} = x_1 y_1 + x_2 y_2 = y_1 x_1 + y_2 x_2 = \mathbf{y} \cdot \mathbf{x}.$$

**Part (b):** It is also simple to see that

$$\begin{aligned} \mathbf{x} \cdot (a\mathbf{y} + b\mathbf{z}) &= x_1(a y_1 + b z_1) + x_2(a y_2 + b z_2) \\ &= a(x_1 y_1 + x_2 y_2) + b(x_1 z_1 + x_2 z_2) \\ &= a(\mathbf{x} \cdot \mathbf{y}) + b(\mathbf{x} \cdot \mathbf{z}). \end{aligned}$$

**Part (c):** This follows from

$$\mathbf{x} \cdot \mathbf{x} = (x_1)^2 + (x_2)^2 \geq 0;$$

furthermore, in the case  $\mathbf{x} \cdot \mathbf{x} = 0$  we obtain that

$$(x_1)^2 + (x_2)^2 = 0 \quad \Leftrightarrow \quad x_1 = x_2 = 0.$$

This establishes the Theorem.  $\square$

These simple properties are crucial to establish the following result, known as Cauchy-Schwarz inequality for the dot product in  $\mathbb{R}^2$ . This inequality allows to express the dot product of two vectors in  $\mathbb{R}^2$  in terms of the angle between the vectors.

**Theorem 6.1.3 (Cauchy-Schwarz).** *The properties (a)-(c) in Theorem 6.1.2 imply that for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^2$  holds*

$$|\mathbf{x} \cdot \mathbf{y}| \leq \|\mathbf{x}\| \|\mathbf{y}\|.$$

**Proof of Theorem 6.1.3:** From the positive definiteness property we know that the following inequality holds for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^2$  and for all  $a \in \mathbb{R}$ ,

$$0 \leq \|a\mathbf{x} - \mathbf{y}\|^2 = (a\mathbf{x} - \mathbf{y}) \cdot (a\mathbf{x} - \mathbf{y}).$$

The symmetry and the linearity on the second argument imply

$$0 \leq (a\mathbf{x} - \mathbf{y}) \cdot (a\mathbf{x} - \mathbf{y}) = a^2 \|\mathbf{x}\|^2 - 2a(\mathbf{x} \cdot \mathbf{y}) + \|\mathbf{y}\|^2. \quad (6.1)$$

Since the inequality above holds for all  $a \in \mathbb{R}$ , let us choose a particular value of  $a$ , the solution of the equation

$$a \|\mathbf{x}\|^2 - (\mathbf{x} \cdot \mathbf{y}) = 0 \quad \Rightarrow \quad a = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\|^2}.$$

Introduce this particular value of  $a$  into Eq. (6.1),

$$0 \leq -\left(\frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\|^2}\right)(\mathbf{x} \cdot \mathbf{y}) + \|\mathbf{y}\|^2 \quad \Rightarrow \quad |\mathbf{x} \cdot \mathbf{y}|^2 \leq \|\mathbf{x}\|^2 \|\mathbf{y}\|^2.$$

This establishes the Theorem.  $\square$

The Cauchy-Schwarz inequality implies that we can express the dot product of two vectors in an alternative and more geometrical way, in terms of an angle related with the two vectors. The Cauchy-Schwarz inequality says

$$-1 \leq \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} \leq 1,$$

which suggests that the number  $(\mathbf{x} \cdot \mathbf{y})/(\|\mathbf{x}\| \|\mathbf{y}\|)$  can be expressed as a sine or a cosine of an appropriate angle.

**Theorem 6.1.4.** *The **angle** between vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^2$  is the number  $\theta \in [0, \pi]$  given by*

$$\cos(\theta) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}.$$

**Proof of Theorem 6.1.4:** It is not difficult to see that given any vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^2$ , the vectors  $\mathbf{x}/\|\mathbf{x}\|$  and  $\mathbf{y}/\|\mathbf{y}\|$  have unit norm. Indeed,

$$\left\| \frac{\mathbf{x}}{\|\mathbf{x}\|} \right\|^2 = \frac{(x_1)^2}{\|\mathbf{x}\|^2} + \frac{(x_2)^2}{\|\mathbf{x}\|^2} = \frac{1}{\|\mathbf{x}\|^2} [(x_1)^2 + (x_2)^2] = 1.$$

The same holds for the vector  $\mathbf{y}/\|\mathbf{y}\|$ . The expression

$$\frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} = \frac{\mathbf{x}}{\|\mathbf{x}\|} \cdot \frac{\mathbf{y}}{\|\mathbf{y}\|},$$

shows that the number  $(\mathbf{x} \cdot \mathbf{y})/(\|\mathbf{x}\| \|\mathbf{y}\|)$  is the inner product of two vectors in the unit circle, as shown in Fig. 41.

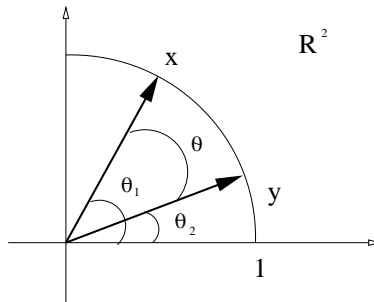


FIGURE 41. The dot product of two vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^2$  can be expressed in terms of the angle  $\theta = \theta_1 - \theta_2$  between the vectors.

Therefore, we know that

$$\frac{\mathbf{x}}{\|\mathbf{x}\|} = \begin{bmatrix} \cos(\theta_1) \\ \sin(\theta_1) \end{bmatrix}, \quad \frac{\mathbf{y}}{\|\mathbf{y}\|} = \begin{bmatrix} \cos(\theta_2) \\ \sin(\theta_2) \end{bmatrix}.$$

Their dot product is given by

$$\frac{\mathbf{x}}{\|\mathbf{x}\|} \cdot \frac{\mathbf{y}}{\|\mathbf{y}\|} = [\cos(\theta_1), \sin(\theta_1)] \begin{bmatrix} \cos(\theta_2) \\ \sin(\theta_2) \end{bmatrix} = \cos(\theta_1) \cos(\theta_2) + \sin(\theta_1) \sin(\theta_2).$$

Using the formula  $\cos(\theta_1) \cos(\theta_2) + \sin(\theta_1) \sin(\theta_2) = \cos(\theta_1 - \theta_2)$ , and denoting the angle between the vectors by  $\theta = \theta_1 - \theta_2$ , we conclude that

$$\frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} = \cos(\theta).$$

This establishes the Theorem. □



Recall the notion of perpendicular vectors.

**Definition 6.1.5.** The vectors  $x, y \in \mathbb{R}^2$  are **orthogonal**, denoted as  $x \perp y$ , iff the angle  $\theta \in [0, \pi]$  between the vectors is  $\theta = \pi/2$ .

The notion of orthogonal vectors in Def. 6.1.5 can be expressed in terms of the dot product, and it is equivalent to Pythagoras Theorem on right triangles.

**Theorem 6.1.6.** Let  $x, y \in \mathbb{R}^2$  be non-zero vectors, then the following statement holds,

$$x \perp y \Leftrightarrow x \cdot y = 0 \Leftrightarrow \|x - y\|^2 = \|x\|^2 + \|y\|^2.$$

**Proof of Theorem 6.1.6:** The non-zero vectors  $x$  and  $y \in \mathbb{R}^2$  are orthogonal iff  $\theta = \pi/2$ , which is equivalent to

$$\frac{x \cdot y}{\|x\| \|y\|} = 0 \Leftrightarrow x \cdot y = 0.$$

The last part of the Proposition comes from the following calculation,

$$\begin{aligned} \|x - y\|^2 &= (x_1 - y_1)^2 + (x_2 - y_2)^2 \\ &= (x_1)^2 + (x_2)^2 + (y_1)^2 + (y_2)^2 - 2(x_1y_1 + x_2y_2) \\ &= \|x\|^2 + \|y\|^2 - 2x \cdot y. \end{aligned}$$

Hence,  $x \perp y$  iff  $x \cdot y = 0$  iff Pythagoras Theorem holds for the triangle with sides given by  $x, y$  and hypotenuse  $x - y$ . This establishes the Theorem.  $\square$

**EXAMPLE 6.1.1:** Find the length of the vectors  $x = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$  and  $y = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$ , the angle between them, and then find a non-zero vector  $z$  orthogonal to  $x$ .

**SOLUTION:** We first find the length, that is, the norms of  $x$  and  $y$ ,

$$\begin{aligned} \|x\|^2 &= x \cdot x = x^T x = [1 \ 2] \cdot \begin{bmatrix} 1 \\ 2 \end{bmatrix} = 1 + 4 \Rightarrow \|x\| = \sqrt{5}, \\ \|y\|^2 &= y \cdot y = y^T y = [3 \ 1] \cdot \begin{bmatrix} 3 \\ 1 \end{bmatrix} = 9 + 1 \Rightarrow \|y\| = \sqrt{10}. \end{aligned}$$

We now find the angle between  $x$  and  $y$ ,

$$\cos(\theta) = \frac{x \cdot y}{\|x\| \|y\|} = \frac{[1 \ 2] \cdot \begin{bmatrix} 3 \\ 1 \end{bmatrix}}{\sqrt{5} \sqrt{10}} = \frac{5}{5\sqrt{2}} = \frac{1}{\sqrt{2}} \Rightarrow \theta = \frac{\pi}{4}.$$

We now find  $z$  such that  $z \perp x$ , that is,

$$0 = [z_1 \ z_2] \cdot \begin{bmatrix} 1 \\ 2 \end{bmatrix} = z_1 + 2z_2 \Rightarrow \begin{cases} z_1 = -2z_2 \\ z_2 \text{ free variable} \end{cases} \Rightarrow z = \begin{bmatrix} -2 \\ 1 \end{bmatrix} z_2.$$

$\triangleleft$

**6.1.2. Dot product in  $\mathbb{F}^n$ .** The notion of dot product reviewed above can be generalized in a straightforward way from  $\mathbb{R}^2$  to  $\mathbb{F}^n$ ,  $n \geq 1$ , where  $\mathbb{F} \in \{\mathbb{R}, \mathbb{C}\}$ .

**Definition 6.1.7.** The **dot product** on the vector space  $\mathbb{F}^n$ , with  $n \geq 1$ , is the function  $\cdot : \mathbb{F}^n \times \mathbb{F}^n \rightarrow \mathbb{F}$  given by

$$x \cdot y = x^* y,$$

where  $x, y$  denote components in the standard basis of  $\mathbb{F}^n$ . The **dot product norm** of a vector  $x \in \mathbb{F}^n$  is the value of the function  $\|\cdot\| : \mathbb{F}^n \rightarrow \mathbb{R}$ ,

$$\|x\| = \sqrt{x \cdot x}.$$

The **norm distance** between  $x, y \in \mathbb{F}^n$  is the value of the function  $d : \mathbb{F}^n \times \mathbb{F}^n \rightarrow \mathbb{R}$ ,

$$d(x, y) = \|x - y\|.$$

The vectors  $x, y \in \mathbb{F}^n$  are **orthogonal**, denoted as  $x \perp y$ , iff holds  $x \cdot y = 0$ .

Notice that we defined two vectors to be orthogonal by the condition that their dot product vanishes. This is the appropriate generalization to  $\mathbb{F}^n$  of the ideas we saw in  $\mathbb{R}^2$ . The concept of angle is more difficult to study. In the case that  $\mathbb{F} = \mathbb{C}$  is not clear what the angle between vectors mean. In the case  $\mathbb{F} = \mathbb{R}$  and  $n > 3$  we have to define angle by the number  $(x \cdot y)/(\|x\| \|y\|)$ . This will be done after we prove the Cauchy-Schwarz inequality, which then is used to show that the number  $(x \cdot y)/(\|x\| \|y\|) \in [-1, 1]$ . The formulas above for the dot product, norm and distance can be expressed in terms of the vector components in the standard basis as follows,

$$\begin{aligned} x \cdot y &= \bar{x}_1 y_1 + \cdots + \bar{x}_n y_n, \\ \|x\| &= \sqrt{|x_1|^2 + \cdots + |x_n|^2}, \\ d(x, y) &= \sqrt{|x_1 - y_1|^2 + \cdots + |x_n - y_n|^2}, \end{aligned}$$

where we used the standard notation  $x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$ ,  $y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$ , and  $|x_i|^2 = \bar{x}_i x_i$ , for  $i = 1, \dots, n$ .

In the particular case that  $\mathbb{F} = \mathbb{R}$  all the vector components are real numbers, so  $\bar{x}_i = x_i$ .

**EXAMPLE 6.1.2:** Find whether  $x$  is orthogonal to  $y$  and/or  $z$ , where

$$x = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}, \quad y = \begin{bmatrix} -5 \\ 4 \\ -3 \\ 2 \end{bmatrix}, \quad z = \begin{bmatrix} -4 \\ -3 \\ 2 \\ 1 \end{bmatrix}.$$

**SOLUTION:** We need to compute the dot products  $x \cdot y$  and  $x \cdot z$ . We obtain

$$\begin{aligned} x^T y &= [1 \quad 2 \quad 3 \quad 4] \begin{bmatrix} -5 \\ 4 \\ -3 \\ 2 \end{bmatrix} = -5 + 8 - 9 + 8 \Rightarrow x \cdot y = 2 \Rightarrow x \not\perp y, \\ x^T z &= [1 \quad 2 \quad 3 \quad 4] \begin{bmatrix} -4 \\ -3 \\ 2 \\ 1 \end{bmatrix} = -4 - 6 + 6 + 4 \Rightarrow x \cdot z = 0 \Rightarrow x \perp z. \end{aligned}$$

◁

**EXAMPLE 6.1.3:** Find  $x \cdot y$ , where  $x = \begin{bmatrix} 2 + 3i \\ i \\ 1 - i \end{bmatrix}$  and  $y = \begin{bmatrix} 2i \\ 1 \\ 1 + 3i \end{bmatrix}$ .

**SOLUTION:** The first product  $x \cdot y$  is given by

$$x^* y = [2 - 3i \quad -i \quad 1 + i] \begin{bmatrix} 2i \\ 1 \\ 1 + 3i \end{bmatrix} = (2 - 3i)(2i) - i + (1 + i)(1 + 3i),$$

so  $x \cdot y = 4i + 6 - i + 1 - 3 + i + 3i$ , that is,  $x \cdot y = 4 + 7i$ .

◁

The dot product satisfies the following properties.

**Theorem 6.1.8.** *The dot product on  $\mathbb{F}^n$ , with  $n \geq 1$ , satisfies for every vector  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{F}^n$  and every scalar  $a, b \in \mathbb{F}$ , the following properties:*

- (a1)  $\mathbf{x} \cdot \mathbf{y} = \mathbf{y} \cdot \mathbf{x}$ , (Symmetry  $\mathbb{F} = \mathbb{R}$ );  
 (a2)  $\mathbf{x} \cdot \mathbf{y} = \overline{\mathbf{y} \cdot \mathbf{x}}$ , (Conjugate symmetry, for  $\mathbb{F} = \mathbb{C}$ );  
 (b)  $\mathbf{x} \cdot (a\mathbf{y} + b\mathbf{z}) = a(\mathbf{x} \cdot \mathbf{y}) + b(\mathbf{x} \cdot \mathbf{z})$ , (Linearity on the second argument);  
 (c)  $\mathbf{x} \cdot \mathbf{x} \geq 0$ , and  $\mathbf{x} \cdot \mathbf{x} = 0$  iff  $\mathbf{x} = \mathbf{0}$ , (Positive definiteness).

**Proof of Theorem 6.1.8:** Use the expression of the dot product in terms of the vector components. The property in (a1) can be established as follows,

$$\mathbf{x} \cdot \mathbf{y} = x_1y_1 + \cdots + x_ny_n = y_1x_1 + \cdots + y_nx_n = \mathbf{y} \cdot \mathbf{x}.$$

The property in (a2) can be established as follows,

$$\mathbf{x} \cdot \mathbf{y} = \overline{\overline{x_1}y_1 + \cdots + \overline{x_n}y_n} = \overline{\overline{y_1}x_1 + \cdots + \overline{y_n}x_n} = \overline{\mathbf{y} \cdot \mathbf{x}}.$$

The property in (b) is shown in a similar way,

$$\mathbf{x} \cdot (a\mathbf{y} + b\mathbf{z}) = \mathbf{x}^*(a\mathbf{y} + b\mathbf{z}) = a\mathbf{x}^*\mathbf{y} + b\mathbf{x}^*\mathbf{z} = a(\mathbf{x} \cdot \mathbf{y}) + b(\mathbf{x} \cdot \mathbf{z}).$$

The property in (c) follows from

$$\mathbf{x} \cdot \mathbf{x} = \mathbf{x}^*\mathbf{x} = |x_1|^2 + \cdots + |x_n|^2 \geq 0;$$

furthermore, in the case that  $\mathbf{x} \cdot \mathbf{x} = 0$  we obtain that

$$|x_1|^2 + \cdots + |x_n|^2 = 0 \Leftrightarrow x_1 = \cdots = x_n = 0 \Leftrightarrow \mathbf{x} = \mathbf{0}.$$

This establishes the Theorem.  $\square$

The positive definiteness property (c) above shows that the dot product norm is indeed a real-valued and not a complex-valued function, since  $\mathbf{x} \cdot \mathbf{x} \geq 0$  implies that  $\|\mathbf{x}\| = \sqrt{\mathbf{x} \cdot \mathbf{x}} \in \mathbb{R}$ . In the case of  $\mathbb{F} = \mathbb{R}$ , the symmetry property and the linearity in the second argument property imply that the dot product on  $\mathbb{R}^n$  is also linear in the first argument. This is a reason to call the dot product on  $\mathbb{R}^n$  a bilinear form. Finally, notice that in the case  $\mathbb{F} = \mathbb{C}$ , the conjugate symmetry property and the linearity in the second argument imply that *the dot product on  $\mathbb{C}^n$  is conjugate linear on the first argument*. The proof is the following:

$$(a\mathbf{y} + b\mathbf{z}) \cdot \mathbf{x} = \overline{\mathbf{x} \cdot (a\mathbf{y} + b\mathbf{z})} = \overline{a(\mathbf{x} \cdot \mathbf{y}) + b(\mathbf{x} \cdot \mathbf{z})} = \overline{a} \overline{(\mathbf{x} \cdot \mathbf{y})} + \overline{b} \overline{(\mathbf{x} \cdot \mathbf{z})},$$

that is, for all  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{C}^n$  and all  $a, b \in \mathbb{C}$  holds

$$(a\mathbf{y} + b\mathbf{z}) \cdot \mathbf{x} = \overline{a}(\mathbf{y} \cdot \mathbf{x}) + \overline{b}(\mathbf{z} \cdot \mathbf{x}).$$

Hence we say that the dot product on  $\mathbb{C}^n$  is conjugate linear in the first argument.

**EXAMPLE 6.1.4:** Compute the dot product of  $\mathbf{x} = \begin{bmatrix} 2 + 3i \\ 6i - 9 \end{bmatrix}$  with  $\mathbf{y} = \begin{bmatrix} 3i \\ 2 \end{bmatrix}$ .

**SOLUTION:** This is a straightforward computation

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{x}^*\mathbf{y} = \begin{bmatrix} 2 - 3i & -6i - 9 \end{bmatrix} \begin{bmatrix} 3i \\ 2 \end{bmatrix} = 6i + 9 - 12i - 18 \Rightarrow \mathbf{x} \cdot \mathbf{y} = -9 - 6i.$$

Notice that  $\mathbf{x} = (2 + 3i)\hat{\mathbf{x}}$ , with  $\hat{\mathbf{x}} = \begin{bmatrix} 1 \\ 3i \end{bmatrix}$ , so we could use the conjugate linearity in the first argument to compute

$$\mathbf{x} \cdot \mathbf{y} = ((2 + 3i)\hat{\mathbf{x}}) \cdot \mathbf{y} = (2 - 3i)(\hat{\mathbf{x}} \cdot \mathbf{y}) = (2 - 3i) \begin{bmatrix} 1 & -3i \end{bmatrix} \begin{bmatrix} 3i \\ 2 \end{bmatrix} = (2 - 3i)(3i - 6i),$$

and we obtain the same result,  $\mathbf{x} \cdot \mathbf{y} = -9 - 6i$ . Finally, notice that  $\mathbf{y} \cdot \mathbf{x} = -9 + 6i$ .  $\triangleleft$

An important result is that the dot product in  $\mathbb{F}^n$  satisfies the Cauchy-Schwarz inequality.

**Theorem 6.1.9 (Cauchy-Schwarz).** *The properties (a1)-(c) in Theorem 6.1.8 imply that for all  $x, y \in \mathbb{F}^n$  holds*

$$|x \cdot y| \leq \|x\| \|y\|.$$

**REMARK:** The proof of the Cauchy-Schwarz inequality only uses the three properties of the dot product presented in Theorem 6.1.8. Any other function  $f : \mathbb{F}^n \times \mathbb{F}^n \rightarrow \mathbb{F}$  having these three properties also satisfies the Cauchy-Schwarz inequality.

**Proof of Theorem 6.1.9:** From the positive definiteness property we know that the following inequality holds for all  $x, y \in \mathbb{F}^n$  and for all  $a \in \mathbb{F}$ ,

$$0 \leq \|ax - y\|^2 = (ax - y) \cdot (ax - y).$$

The symmetry and the linearity on the second argument imply

$$0 \leq (ax - y) \cdot (ax - y) = a \bar{a} \|x\|^2 - \bar{a} (x \cdot y) - a (y \cdot x) + \|y\|^2. \quad (6.2)$$

Since the inequality above holds for all  $a \in \mathbb{F}$ , let us choose a particular value of  $a$ , the solution of the equation

$$a \bar{a} \|x\|^2 - \bar{a} (x \cdot y) = 0 \quad \Rightarrow \quad a = \frac{x \cdot y}{\|x\|^2}.$$

Introduce this particular value of  $a$  into Eq. (6.2),

$$0 \leq -\left(\frac{x \cdot y}{\|x\|^2}\right)(\overline{x \cdot y}) + \|y\|^2 \quad \Rightarrow \quad |x \cdot y|^2 \leq \|x\|^2 \|y\|^2.$$

This establishes the Theorem.  $\square$

In the case  $\mathbb{F} = \mathbb{R}$ , the Cauchy-Schwarz inequality in  $\mathbb{R}^n$  implies that the number  $(x \cdot y)/(\|x\| \|y\|) \in [-1, 1]$ , which is a necessary and sufficient condition for the following definition of angle between two vectors in  $\mathbb{R}^n$ .

**Definition 6.1.10.** *The **angle** between vectors  $x, y \in \mathbb{R}^n$  is the number  $\theta \in [0, \pi]$  given by*

$$\cos(\theta) = \frac{x \cdot y}{\|x\| \|y\|}.$$

The dot product norm function in Definition 6.1.7 satisfies the following properties.

**Theorem 6.1.11.** *The dot product norm function on  $\mathbb{F}^n$ , with  $n \geq 1$ , satisfies for every vector  $x, y \in \mathbb{F}^n$  and every scalar  $a \in \mathbb{F}$  the following properties:*

- (a)  $\|x\| \geq 0$ , and  $\|x\| = 0$  iff  $x = 0$ , (*Positive definiteness*);
- (b)  $\|ax\| = |a| \|x\|$ , (*Scaling*);
- (c)  $\|x + y\| \leq \|x\| + \|y\|$ , (*Triangle inequality*).

**Proof of Theorem 6.1.11:** Properties (a) and (b) are straightforward to show from the definition of dot product, and their proof is left as an exercise. We show here how to obtain the triangle inequality, property (c). The proof uses the Cauchy-Schwarz inequality presented in Theorem 6.1.9. Given any vectors  $x, y \in \mathbb{F}^n$  holds

$$\begin{aligned} \|x + y\|^2 &= (x + y) \cdot (x + y) \\ &= \|x\|^2 + (x \cdot y) + (y \cdot x) + \|y\|^2 \\ &\leq \|x\|^2 + 2|x \cdot y| + \|y\|^2 \\ &\leq \|x\|^2 + 2\|x\| \|y\| + \|y\|^2 = (\|x\| + \|y\|)^2, \end{aligned}$$

We conclude that  $\|x + y\| \leq (\|x\| + \|y\|)$ . This establishes the Theorem.  $\square$

A vector  $v \in \mathbb{F}^n$  is called *normal or unit vector* iff  $\|v\| = 1$ . Examples of unit vectors are the standard basis vectors. Unit vectors parallel to a given vector are simple to find.

**Theorem 6.1.12.** *If  $\mathbf{v} \in \mathbb{F}^n$  is non-zero, then  $\frac{\mathbf{v}}{\|\mathbf{v}\|}$  is a unit vector parallel to  $\mathbf{v}$ .*

**Proof of Theorem 6.1.12:** Notice that  $\mathbf{u} = \frac{\mathbf{v}}{\|\mathbf{v}\|}$  is parallel to  $\mathbf{v}$ , and it is straightforward to check that  $\mathbf{u}$  is a unit vector, since

$$\|\mathbf{u}\| = \left\| \frac{\mathbf{v}}{\|\mathbf{v}\|} \right\| = \frac{1}{\|\mathbf{v}\|} \|\mathbf{v}\| = 1.$$

This establishes the Theorem. □

**EXAMPLE 6.1.5:** Find a unit vector parallel to  $\mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$ .

**SOLUTION:** First compute the norm of  $\mathbf{x}$ ,

$$\|\mathbf{x}\| = \sqrt{1 + 4 + 9} = \sqrt{14},$$

therefore  $\mathbf{u} = \frac{1}{\sqrt{14}} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$  is a unit vector parallel to  $\mathbf{v}$ . ◁

## 6.1.3. Exercises.

**6.1.1.-** Consider the vector space  $\mathbb{R}^4$  with standard basis  $\mathcal{S}$  and dot product. Find the norm of  $u$  and  $v$ , their distance and the angle between them, where

$$u = \begin{bmatrix} 2 \\ 1 \\ -4 \\ -2 \end{bmatrix}, \quad v = \begin{bmatrix} 1 \\ -1 \\ 1 \\ -1 \end{bmatrix}.$$

**6.1.2.-** Use the dot product on  $\mathbb{R}^2$  to find two unit vectors orthogonal to

$$x = \begin{bmatrix} 3 \\ 2 \end{bmatrix}.$$

**6.1.3.-** Use the dot product on  $\mathbb{C}^2$  to find a unit vector parallel to

$$x = \begin{bmatrix} 1 + 2i \\ 2 - i \end{bmatrix}.$$

**6.1.4.-** Consider the vector space  $\mathbb{R}^2$  with the dot product.

- Give an example of a linearly independent set  $\{x, y\}$  with  $x \not\perp y$ .
- Give an example of a linearly dependent set  $\{x, y\}$  with  $x \perp y$ .

**6.1.5.-** Consider the vector space  $\mathbb{F}^n$  with the dot product, and let  $\operatorname{Re}$  denote the real part of a complex number. Show that for all  $x, y \in \mathbb{F}^n$  holds

$$\|x - y\|^2 = \|x\|^2 + \|y\|^2 - 2\operatorname{Re}(x \cdot y).$$

**6.1.6.-** Use the result in Exercise **6.1.5** above to prove the following generalizations of the Pythagoras Theorem to  $\mathbb{F}^n$  with the dot product.

(a) For  $x, y \in \mathbb{R}^n$  holds

$$x \perp y \Leftrightarrow \|x + y\|^2 = \|x\|^2 + \|y\|^2.$$

(b) For  $x, y \in \mathbb{C}^n$  holds

$$x \perp y \Rightarrow \|x + y\|^2 = \|x\|^2 + \|y\|^2.$$

**6.1.7.-** Prove that the parallelogram law holds for the dot product norm in  $\mathbb{F}^n$ , that is, show that for all  $x, y \in \mathbb{F}^n$  holds

$$\|x + y\|^2 + \|x - y\|^2 = 2\|x\|^2 + 2\|y\|^2.$$

This law states that the sum of the squares of the lengths of the four sides of a parallelogram formed by  $x$  and  $y$  equals the sum of the square of the lengths of the two diagonals.

## 6.2. INNER PRODUCT

**6.2.1. Inner product.** An inner product on a vector space is a generalization of the dot product on  $\mathbb{R}^n$  or  $\mathbb{C}^n$  introduced in Sect. 6.1. The inner product is not defined with a particular formula, or requiring a particular basis in the vector space. Instead, the inner product is defined by a list of properties that must satisfy. We did something similar when we introduced the concept of a vector space. In that case we defined a vector space as a set of any kind of elements where linear combinations are possible, instead of defining the set by explicitly giving its elements.

**Definition 6.2.1.** Let  $V$  be a vector space over the scalar field  $\mathbb{F} \in \{\mathbb{R}, \mathbb{C}\}$ . A function  $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{F}$  is called an **inner product** iff for every  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in V$  and every  $a, b \in \mathbb{F}$  the function  $\langle \cdot, \cdot \rangle$  satisfies:

- (a1)  $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$ , (Symmetry, for  $\mathbb{F} = \mathbb{R}$ );  
 (a2)  $\langle \mathbf{x}, \mathbf{y} \rangle = \overline{\langle \mathbf{y}, \mathbf{x} \rangle}$ , (Conjugate symmetry, for  $\mathbb{F} = \mathbb{C}$ );  
 (b)  $\langle \mathbf{x}, (a\mathbf{y} + b\mathbf{z}) \rangle = a\langle \mathbf{x}, \mathbf{y} \rangle + b\langle \mathbf{x}, \mathbf{z} \rangle$ , (Linearity on the second argument);  
 (c)  $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$ , and  $\langle \mathbf{x}, \mathbf{x} \rangle = 0$  iff  $\mathbf{x} = \mathbf{0}$ , (Positive definiteness).

An **inner product space** is a pair  $(V, \langle \cdot, \cdot \rangle)$  of a vector space with an inner product.

Different inner products can be defined on a given vector space. The dot product is an inner product in  $\mathbb{F}^n$ . A different inner product can be defined in  $\mathbb{F}^n$ , as can be seen in the following example.

**EXAMPLE 6.2.1:** We show that  $\mathbb{R}^n$  can have different inner products.

- (a) The dot product on  $\mathbb{R}^n$  is an inner product, since the expression

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}_s^T \mathbf{y}_s = [x_1 \ \cdots \ x_n]_s \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_s = x_1 y_1 + \cdots + x_n y_n, \quad (6.3)$$

satisfies all the properties in Definition 6.2.1, with  $\mathcal{S}$  the standard ordered basis in  $\mathbb{R}^n$ .

- (b) A different inner product in  $\mathbb{R}^n$  can be introduced by a formula similar to the one in Eq. (6.3) by choosing a different ordered basis. If  $\mathcal{U}$  is any ordered basis of  $\mathbb{R}^n$ , then

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}_u^T \mathbf{y}_u.$$

defines an inner product on  $\mathbb{R}^n$ . The inner product defined using the basis  $\mathcal{U}$  is not equal to the inner product defined using the standard basis  $\mathcal{S}$ . Let  $\mathbf{P} = \mathbf{I}_{us}$  be the change of basis matrix, then we know that  $\mathbf{x}_u = \mathbf{P}^{-1} \mathbf{x}_s$ . The inner product above can be expressed in terms of the  $\mathcal{S}$  basis as follows,

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}_s^T \mathbf{M} \mathbf{y}_s, \quad \mathbf{M} = (\mathbf{P}^{-1})^T (\mathbf{P}^{-1}),$$

and in general,  $\mathbf{M} \neq \mathbf{I}_n$ . Therefore, the inner product above is not equal to the dot product. Also, see Example 6.2.2. ◁

**EXAMPLE 6.2.2:** Let  $\mathcal{S}$  be the standard ordered basis in  $\mathbb{R}^2$ , and introduce the ordered basis  $\mathcal{U}$  as the following rescaling of  $\mathcal{S}$ ,

$$\mathcal{U} = (\mathbf{u}_1 = \frac{1}{2} \mathbf{e}_1, \mathbf{u}_2 = \frac{1}{3} \mathbf{e}_2).$$

Express the inner product  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}_u^T \mathbf{y}_u$  in terms of  $\mathbf{x}_s$  and  $\mathbf{y}_s$ . Is this inner product the same as the dot product?

**SOLUTION:** The definition of the inner product says that  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}_u^T \mathbf{y}_u$ . Introducing the notation  $\mathbf{x}_u = \begin{bmatrix} \tilde{x}_1 \\ \tilde{x}_2 \end{bmatrix}_u$  and  $\mathbf{y}_u = \begin{bmatrix} \tilde{y}_1 \\ \tilde{y}_2 \end{bmatrix}_u$ , we obtain the usual expression  $\langle \mathbf{x}, \mathbf{y} \rangle = \tilde{x}_1 \tilde{y}_1 + \tilde{x}_2 \tilde{y}_2$ . The components  $\mathbf{x}_u$  and  $\mathbf{x}_s$  are related by the change of basis formula

$$\mathbf{x}_u = \mathbf{P}^{-1} \mathbf{x}_s, \quad \mathbf{P} = \mathbf{I}_{us} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/3 \end{bmatrix}_{us} \quad \Rightarrow \quad \mathbf{P}^{-1} = \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}_{su} = (\mathbf{P}^{-1})^T.$$

Therefore,

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}_u^T \mathbf{y}_u = \mathbf{x}_s^T (\mathbf{P}^{-1})^T \mathbf{P}^{-1} \mathbf{y}_s = [x_1, \quad x_2]_s \begin{bmatrix} 4 & 0 \\ 0 & 9 \end{bmatrix}_{su} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}_s$$

where we used the standard notation  $\mathbf{x}_s = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_s$  and  $\mathbf{y}_s = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}_s$ . We conclude that

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}_s^T (\mathbf{P}^{-1})^2 \mathbf{y}_s \quad \Leftrightarrow \quad \langle \mathbf{x}, \mathbf{y} \rangle = 4x_1y_1 + 9x_2y_2.$$

The inner product  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}_u^T \mathbf{y}_u$  is **different** from the dot product  $\mathbf{x} \cdot \mathbf{y} = \mathbf{x}_s^T \mathbf{y}_s$ .  $\triangleleft$

**EXAMPLE 6.2.3:** Determine whether the function  $\langle \cdot, \cdot \rangle : \mathbb{R}^3 \times \mathbb{R}^3 \rightarrow \mathbb{R}$  below is an inner product in  $\mathbb{R}^3$ , where

$$\langle \mathbf{x}, \mathbf{y} \rangle = x_1y_1 + x_2y_2 + x_3y_3 + 3x_1y_2 + 3x_2y_1.$$

**SOLUTION:** The function  $\langle \cdot, \cdot \rangle$  seems to be symmetric and linear. It is not so clear whether this function is positive, because of the presence of crossed terms. So, before spending time to prove the symmetry and linearity properties, we first concentrate on the property that might fail, positivity. If positivity fails, we don't need to prove the remaining properties. The crossed terms  $3(x_1y_2 + x_2y_1)$  in the definition of the inner product suggest that the product is not positive, since the factor 3 makes them too important compared with the other terms. Let us try to find an example, that is a vector  $\mathbf{x} \neq \mathbf{0}$  such that  $\langle \mathbf{x}, \mathbf{x} \rangle \leq 0$ . Let us try with

$$\mathbf{x} = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} \quad \Rightarrow \quad \langle \mathbf{x}, \mathbf{x} \rangle = 1^2 + (-1)^2 + 0^2 + 3[(1)(-1) + (-1)(1)] = 2 - 6 = -4.$$

Since  $\langle \mathbf{x}, \mathbf{x} \rangle = -4$ , the function  $\langle \cdot, \cdot \rangle$  is not positive, hence **it is not an inner product**.  $\triangleleft$

**EXAMPLE 6.2.4:** Consider the vector space  $\mathbb{F}^{m,n}$  of all  $m \times n$  matrices. Show that an inner product on that space is the function  $\langle \cdot, \cdot \rangle_F : \mathbb{F}^{m,n} \times \mathbb{F}^{m,n} \rightarrow \mathbb{F}$

$$\langle \mathbf{A}, \mathbf{B} \rangle_F = \text{tr}(\mathbf{A}^* \mathbf{B}).$$

The inner product is called *Frobenius inner product*.

**SOLUTION:** We show that the Frobenius function  $\langle \cdot, \cdot \rangle_F$  above satisfies the three properties in Def. 6.2.1. We use the component notation  $\mathbf{A} = [A_{ij}]$ ,  $\mathbf{B} = [B_{kl}]$ , with  $i, k = 1, \dots, m$  and  $j, l = 1, \dots, n$ , so

$$(\mathbf{A}^* \mathbf{B})_{jl} = \sum_{i=1}^m (\overline{A}^T)_{ji} B_{il} = \sum_{i=1}^m \overline{A}_{ij} B_{il} \quad \Rightarrow \quad \langle \mathbf{A}, \mathbf{B} \rangle_F = \sum_{j=1}^n \sum_{i=1}^n \overline{A}_{ij} B_{ij}.$$

The first property is satisfied, since

$$\langle \mathbf{A}, \mathbf{A} \rangle_F = \text{tr}(\mathbf{A}^* \mathbf{A}) = \sum_{j=1}^n \sum_{i=1}^m |A_{ij}|^2 \geq 0,$$



and  $\langle \mathbf{A}, \mathbf{A} \rangle_F = 0$  iff  $A_{ij} = 0$  for every indices  $i, j$ , which is equivalently to  $\mathbf{A} = \mathbf{0}$ . The second property is satisfied, since

$$\langle \mathbf{A}, \mathbf{B} \rangle_F = \sum_{j=1}^n \sum_{i=1}^n \overline{A_{ij}} B_{ij} = \overline{\sum_{j=1}^n \sum_{i=1}^n B_{ij} A_{ij}} = \overline{\langle \mathbf{B}, \mathbf{A} \rangle_F}.$$

The same proof can be expressed in index-free notation using the properties of the trace,

$$\operatorname{tr}(\mathbf{A}^* \mathbf{B}) = \operatorname{tr}(\overline{\mathbf{A}}^T \mathbf{B}) = \operatorname{tr}[(\overline{\mathbf{A}}^T \mathbf{B})^T] = \operatorname{tr}(\mathbf{B}^T \overline{\mathbf{A}}) = \overline{\operatorname{tr}(\mathbf{B}^* \mathbf{A})},$$

that is,  $\langle \mathbf{A}, \mathbf{B} \rangle_F = \overline{\langle \mathbf{B}, \mathbf{A} \rangle_F}$ . The third property comes from the distributive property of the matrix product, that is,

$$\langle \mathbf{A}, (a\mathbf{B} + b\mathbf{C}) \rangle_F = \operatorname{tr}(\mathbf{A}^* (a\mathbf{B} + b\mathbf{C})) = a \operatorname{tr}(\mathbf{A}^* \mathbf{B}) + b \operatorname{tr}(\mathbf{A}^* \mathbf{C}) = a \langle \mathbf{A}, \mathbf{B} \rangle_F + b \langle \mathbf{A}, \mathbf{C} \rangle_F.$$

This establishes that  $\langle \cdot, \cdot \rangle_F$  is an inner product.  $\triangleleft$

**EXAMPLE 6.2.5:** Compute the Frobenius inner product  $\langle \mathbf{A}, \mathbf{B} \rangle_F$  where

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 1 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 3 & 2 & 1 \\ 2 & 1 & 2 \end{bmatrix} \in \mathbb{R}^{2,3}.$$

**SOLUTION:** Since the matrices have real coefficients, the Frobenius inner product has the form  $\langle \mathbf{A}, \mathbf{B} \rangle_F = \operatorname{tr}(\mathbf{A}^T \mathbf{B})$ . So, we need to compute the diagonal elements in the product

$$\mathbf{A}^T \mathbf{B} = \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} 3 & 2 & 1 \\ 2 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 7 & * & * \\ * & 8 & * \\ * & * & 5 \end{bmatrix} \Rightarrow \langle \mathbf{A}, \mathbf{B} \rangle_F = 7 + 8 + 5 \Rightarrow \langle \mathbf{A}, \mathbf{B} \rangle_F = 20.$$

$\triangleleft$

**EXAMPLE 6.2.6:** Consider the vector space  $\mathbb{P}_n([-1, 1])$  of polynomials with real coefficients having degree less or equal  $n \geq 1$  and being defined on the interval  $[-1, 1]$ . Show that an inner product in this space is the following:

$$\langle \mathbf{p}, \mathbf{q} \rangle = \int_{-1}^1 \mathbf{p}(x) \mathbf{q}(x) dx. \quad \mathbf{p}, \mathbf{q} \in \mathbb{P}_n.$$

**SOLUTION:** We need to verify the three properties in the Definition 6.2.1. The positive definiteness property is satisfied, since

$$\langle \mathbf{p}, \mathbf{p} \rangle = \int_{-1}^1 [\mathbf{p}(x)]^2 dx \geq 0,$$

and in the case  $\langle \mathbf{p}, \mathbf{p} \rangle = 0$  this implies that the integrand must vanish, that is,  $[\mathbf{p}(x)]^2 = 0$ , which is equivalent to  $\mathbf{p} = 0$ . The symmetry property is satisfied, since  $\mathbf{p}(x) \mathbf{q}(x) = \mathbf{q}(x) \mathbf{p}(x)$ , which implies that  $\langle \mathbf{p}, \mathbf{q} \rangle = \langle \mathbf{q}, \mathbf{p} \rangle$ . The linearity property on the second argument is also satisfied, since

$$\begin{aligned} \langle \mathbf{p}, (a\mathbf{q} + b\mathbf{r}) \rangle &= \int_{-1}^1 \mathbf{p}(x) [a\mathbf{q}(x) + b\mathbf{r}(x)] dx \\ &= a \int_{-1}^1 \mathbf{p}(x) \mathbf{q}(x) dx + b \int_{-1}^1 \mathbf{p}(x) \mathbf{r}(x) dx \\ &= a \langle \mathbf{p}, \mathbf{q} \rangle + b \langle \mathbf{p}, \mathbf{r} \rangle. \end{aligned}$$

This establishes that  $\langle \cdot, \cdot \rangle$  is an inner product.  $\triangleleft$

**EXAMPLE 6.2.7:** Consider the vector space  $C^k([a, b], \mathbb{R})$ , with  $k \geq 0$  and  $a < b$ , of  $k$ -times continuously differentiable real-valued functions  $\mathbf{f} : [a, b] \rightarrow \mathbb{R}$ . An inner product in this vector space is given by

$$\langle \mathbf{f}, \mathbf{g} \rangle = \int_a^b \mathbf{f}(x)\mathbf{g}(x) dx.$$

Any positive function  $\mu \in C^0([a, b], \mathbb{R})$  determines an inner product in  $C^k([a, b], \mathbb{R})$  as follows

$$\langle \mathbf{f}, \mathbf{g} \rangle_\mu = \int_a^b \mu(x)\mathbf{f}(x)\mathbf{g}(x) dx.$$

The function  $\mu$  is called a weigh function. An inner product in the vector space  $C^k([a, b], \mathbb{C})$  of  $k$ -times continuously differentiable complex-valued functions  $\mathbf{f} : [a, b] \subset \mathbb{R} \rightarrow \mathbb{C}$  is the following,

$$\langle \mathbf{f}, \mathbf{g} \rangle = \int_a^b \overline{\mathbf{f}(x)}\mathbf{g}(x) dx.$$

◁

An inner product satisfies the following inequality.

**Theorem 6.2.2 (Cauchy-Schwarz).** *If  $(V, \langle \cdot, \cdot \rangle)$  is an inner product space over  $\mathbb{F}$ , then for every  $\mathbf{x}, \mathbf{y} \in V$  holds*

$$|\langle \mathbf{x}, \mathbf{y} \rangle|^2 \leq \langle \mathbf{x}, \mathbf{x} \rangle \langle \mathbf{y}, \mathbf{y} \rangle.$$

Furthermore, equality holds iff  $\mathbf{y} = a\mathbf{x}$ , with  $a = \langle \mathbf{x}, \mathbf{y} \rangle / \langle \mathbf{x}, \mathbf{x} \rangle$ .

**Proof of Theorem 6.2.2:** From the positive definiteness property we know that for every  $\mathbf{x}, \mathbf{y} \in V$  and every scalar  $a \in \mathbb{F}$  holds  $0 \leq \langle (a\mathbf{x} - \mathbf{y}), (a\mathbf{x} - \mathbf{y}) \rangle$ . The symmetry and the linearity on the second argument imply

$$0 \leq \langle (a\mathbf{x} - \mathbf{y}), (a\mathbf{x} - \mathbf{y}) \rangle = a\bar{a}\langle \mathbf{x}, \mathbf{x} \rangle - \bar{a}\langle \mathbf{x}, \mathbf{y} \rangle - a\langle \mathbf{y}, \mathbf{x} \rangle + \langle \mathbf{y}, \mathbf{y} \rangle. \quad (6.4)$$

Since the inequality above holds for all  $a \in \mathbb{F}$ , let us choose a particular value of  $a$ , the solution of the equation

$$a\bar{a}\langle \mathbf{x}, \mathbf{x} \rangle - \bar{a}\langle \mathbf{x}, \mathbf{y} \rangle = 0 \quad \Rightarrow \quad a = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle}.$$

Introduce this particular value of  $a$  into Eq. (6.4),

$$0 \leq -\left(\frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle}\right)\overline{\langle \mathbf{x}, \mathbf{y} \rangle} + \langle \mathbf{y}, \mathbf{y} \rangle \quad \Rightarrow \quad |\langle \mathbf{x}, \mathbf{y} \rangle|^2 \leq \langle \mathbf{x}, \mathbf{x} \rangle \langle \mathbf{y}, \mathbf{y} \rangle.$$

Finally, notice that equality holds iff  $a\mathbf{x} = \mathbf{y}$ , and in this case, computing the inner product with  $\mathbf{x}$  we obtain  $a\langle \mathbf{x}, \mathbf{x} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle$ . This establishes the Theorem.  $\square$

**6.2.2. Inner product norm.** The inner product on a vector space determines a particular notion of length, or norm, of a vector, and we call it the inner product norm. After we introduce this norm we show its main properties. In Chapter 8 later on we use these properties to define a more general notion of norm as any function on the vector space satisfying these properties. The inner product norm is just a particular case of this broader notion of length. A normed space is a vector space with any norm.

**Definition 6.2.3.** *The **inner product norm** determined in an inner product space  $(V, \langle \cdot, \cdot \rangle)$  is the function  $\| \cdot \| : V \rightarrow \mathbb{R}$  given by*

$$\| \mathbf{x} \| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}.$$

The Cauchy-Schwarz inequality is often expressed using the inner product norm as follows: For every  $\mathbf{x}, \mathbf{y} \in V$  holds

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|.$$

A vector  $\mathbf{x} \in V$  is a *normal or unit vector* iff  $\|\mathbf{x}\| = 1$ .

**Theorem 6.2.4.** *If  $\mathbf{v} \neq \mathbf{0}$  belongs to  $(V, \langle \cdot, \cdot \rangle)$ , then  $\frac{\mathbf{v}}{\|\mathbf{v}\|}$  is a unit vector parallel to  $\mathbf{v}$ .*

The proof is the same of Theorem 6.1.12.

**EXAMPLE 6.2.8:** Consider the inner product space  $(\mathbb{F}^{m,n}, \langle \cdot, \cdot \rangle_F)$ , where  $\mathbb{F}^{m,n}$  is the vector space of all  $m \times n$  matrices and  $\langle \cdot, \cdot \rangle_F$  is the Frobenius inner product defined in Example 6.2.4. The associated inner product norm is called the *Frobenius norm* and is given by

$$\|\mathbf{A}\|_F = \sqrt{\langle \mathbf{A}, \mathbf{A} \rangle_F} = \sqrt{\text{tr}(\mathbf{A}^* \mathbf{A})}.$$

If  $\mathbf{A} = [A_{ij}]$ , with  $i = 1, \dots, m$  and  $j = 1, \dots, n$ , then

$$\|\mathbf{A}\|_F = \left( \sum_{i=1}^m \sum_{j=1}^n |A_{ij}|^2 \right)^{1/2}.$$

◁

**EXAMPLE 6.2.9:** Find an explicit expression for the Frobenius norm of any element  $\mathbf{A} \in \mathbb{F}^{2,2}$ .

**SOLUTION:** The Frobenius norm of an arbitrary matrix  $\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \in \mathbb{F}^{2,2}$  is given by

$$\|\mathbf{A}\|_F^2 = \text{tr} \left( \begin{bmatrix} \overline{A_{11}} & \overline{A_{21}} \\ \overline{A_{12}} & \overline{A_{22}} \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \right).$$

Since we are only interested in the diagonal elements of the matrix product in the equation above, we obtain

$$\|\mathbf{A}\|_F^2 = \text{tr} \begin{bmatrix} |A_{11}|^2 + |A_{21}|^2 & * \\ * & |A_{12}|^2 + |A_{22}|^2 \end{bmatrix}$$

which gives the formula

$$\|\mathbf{A}\|_F^2 = |A_{11}|^2 + |A_{12}|^2 + |A_{21}|^2 + |A_{22}|^2.$$

This is the explicit expression of the sum  $\|\mathbf{A}\|_F = \left( \sum_{i=1}^2 \sum_{j=1}^2 |A_{ij}|^2 \right)^{1/2}$ .

◁

The inner product norm function has the following properties.

**Theorem 6.2.5.** *The inner product norm introduced in Definition 6.2.3 satisfies that for every  $\mathbf{x}, \mathbf{y} \in V$  and every  $a \in \mathbb{F}$  holds,*

- (a)  $\|\mathbf{x}\| \geq 0$ , and  $\|\mathbf{x}\| = 0$  iff  $\mathbf{x} = \mathbf{0}$ , (*Positive definiteness*);
- (b)  $\|a\mathbf{x}\| = |a| \|\mathbf{x}\|$ , (*Scaling*);
- (c)  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ , (*Triangle inequality*).

**Proof of Theorem 6.2.5:** Properties (a) and (b) are straightforward to show from the definition of inner product, and their proof is left as an exercise. We show here how to

obtain the triangle inequality, property (c). Given any vectors  $\mathbf{x}, \mathbf{y} \in V$  holds

$$\begin{aligned} \|\mathbf{x} + \mathbf{y}\|^2 &= \langle (\mathbf{x} + \mathbf{y}), (\mathbf{x} + \mathbf{y}) \rangle \\ &= \|\mathbf{x}\|^2 + \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{y}, \mathbf{x} \rangle + \|\mathbf{y}\|^2 \\ &\leq \|\mathbf{x}\|^2 + 2|\langle \mathbf{x}, \mathbf{y} \rangle| + \|\mathbf{y}\|^2 \\ &\leq \|\mathbf{x}\|^2 + 2\|\mathbf{x}\|\|\mathbf{y}\| + \|\mathbf{y}\|^2 = (\|\mathbf{x}\| + \|\mathbf{y}\|)^2, \end{aligned}$$

where the last inequality comes from the Cauchy-Schwarz inequality. We then conclude that  $\|\mathbf{x} + \mathbf{y}\|^2 \leq (\|\mathbf{x}\| + \|\mathbf{y}\|)^2$ . This establishes the Theorem.  $\square$

**6.2.3. Norm distance.** The norm on an inner product space determines a particular notion of distance between vectors. After we introduce this norm we show its main properties.

**Definition 6.2.6.** The *norm distance* between two vectors in a vector space  $V$  with a norm function  $\|\cdot\| : V \rightarrow \mathbb{R}$  is the value of the function  $d : V \times V \rightarrow \mathbb{R}$  given by

$$d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|.$$

**Theorem 6.2.7.** The norm distance in Definition 6.2.6 satisfies for every  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in V$  that

- (a)  $d(\mathbf{x}, \mathbf{y}) \geq 0$ , and  $d(\mathbf{x}, \mathbf{y}) = 0$  iff  $\mathbf{x} = \mathbf{y}$ , (*Positive definiteness*);
- (b)  $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$ , (*Symmetry*);
- (c)  $d(\mathbf{x}, \mathbf{y}) \leq d(\mathbf{x}, \mathbf{z}) + d(\mathbf{z}, \mathbf{y})$ , (*Triangle inequality*).

**Proof of Theorem 6.2.7:** Properties (a) and (b) are straightforward from properties (a) and (b), and their proof are left as an exercise. We show how the triangle inequality for the distance comes from the triangle inequality for the norm. Indeed

$$d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \|(\mathbf{x} - \mathbf{z}) - (\mathbf{y} - \mathbf{z})\| \leq \|\mathbf{x} - \mathbf{z}\| + \|\mathbf{y} - \mathbf{z}\| = d(\mathbf{x}, \mathbf{z}) + d(\mathbf{z}, \mathbf{y}),$$

where we used the symmetry of the distance function on the last term above. This establishes the Theorem.  $\square$

The presence of an inner product, and hence a norm and a distance, on a vector space permits to introduce the notion of convergence of an infinite sequence of vectors. We say that the sequence  $\{\mathbf{x}_n\}_{n=0}^{\infty} \subset V$  *converges* to  $\mathbf{x} \in V$  iff

$$\lim_{n \rightarrow \infty} d(\mathbf{x}_n, \mathbf{x}) = 0.$$

Some of the most important concepts related to convergence are closeness of a subspace, completeness of the vector space, and the continuity of linear operators and linear transformations. In the case of finite dimensional vector spaces the situation is straightforward. All subspaces are closed, all inner product spaces are complete and all linear operators and linear transformations are continuous. However, in the case of infinite dimensional vector spaces, things are not so simple.

## 6.2.4. Exercises.

**6.2.1.-** Determine which of the following functions  $\langle \cdot, \cdot \rangle : \mathbb{R}^3 \times \mathbb{R}^3 \rightarrow \mathbb{R}$  defines an inner product on  $\mathbb{R}^3$ . Justify your answers.

- (a)  $\langle x, y \rangle = x_1y_1 + x_3y_3$ ;
- (b)  $\langle x, y \rangle = x_1y_1 - x_2y_2 + x_3y_3$ ;
- (c)  $\langle x, y \rangle = 2x_1y_1 + x_2y_2 + 4x_3y_3$ ;
- (d)  $\langle x, y \rangle = x_1^2y_1^2 + x_2^2y_2^2 + x_3^2y_3^2$ .

We used the standard notation

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}.$$

**6.2.2.-** Prove that an inner product function  $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{F}$  satisfies the following properties:

- (a)  $\langle x, y \rangle = 0$  for all  $x \in V$ , then  $y = 0$ .
- (b)  $\langle ax, y \rangle = \bar{a} \langle x, y \rangle$  for all  $x, y \in V$ .

**6.2.3.-** Given a matrix  $M \in \mathbb{R}^{2,2}$  introduce the function  $\langle \cdot, \cdot \rangle_M : \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ ,

$$\langle y, x \rangle_M = y^T M x.$$

For each of the matrices  $M$  below determine whether  $\langle \cdot, \cdot \rangle_M$  defines an inner product or not. Justify your answers.

- (a)  $M = \begin{bmatrix} 4 & 1 \\ 1 & 9 \end{bmatrix}$ ;
- (b)  $M = \begin{bmatrix} 4 & -3 \\ 3 & 9 \end{bmatrix}$ ;
- (c)  $M = \begin{bmatrix} 4 & 1 \\ 0 & 9 \end{bmatrix}$ .

**6.2.4.-** Fix any  $A \in \mathbb{R}^{n,n}$  with  $N(A) = \{0\}$  and introduce  $M = A^T A$ . Prove that  $\langle \cdot, \cdot \rangle_M : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ , given by

$$\langle y, x \rangle = y^T M x.$$

is an inner product in  $\mathbb{R}^n$ .

**6.2.5.-** Find  $k \in \mathbb{R}$  such that the matrices  $A, B \in \mathbb{R}^{2,2}$  are perpendicular in the Frobenius inner product,

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 2 \\ k & 1 \end{bmatrix}.$$

**6.2.6.-** Evaluate the Frobenius norm for the matrices

$$A = \begin{bmatrix} 1 & -2 \\ -1 & 2 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}.$$

**6.2.7.-** Prove that  $\|A\|_F = \|A^*\|_F$  for all  $A \in \mathbb{F}^{m,n}$ .

**6.2.8.-** Consider the vector space  $\mathbb{P}_2([0, 1])$  with inner product

$$\langle p, q \rangle = \int_0^1 p(x)q(x) dx.$$

Find a unit vector parallel to

$$p(x) = 3 - 5x^2.$$

## 6.3. ORTHOGONAL VECTORS

**6.3.1. Definition and examples.** In the previous Section we introduced the notion of inner product in a vector space. This structure provides a notion of vector norm and distance between vectors. In this Section we explore another concept provided by an inner product; the notion of perpendicular vectors and the notion of angle between vectors in a real vector space. We start defining perpendicular vectors on any inner product space.

**Definition 6.3.1.** Two vectors  $\mathbf{x}, \mathbf{y}$  in an inner product space  $(V, \langle \cdot, \cdot \rangle)$  are **orthogonal** or **perpendicular**, denoted as  $\mathbf{x} \perp \mathbf{y}$ , iff holds  $\langle \mathbf{x}, \mathbf{y} \rangle = 0$ .

The Pythagoras Theorem holds on any inner product space.

**Theorem 6.3.2.** Let  $(V, \langle \cdot, \cdot \rangle)$  be an inner product space over the field  $\mathbb{F}$ .

(a) If  $\mathbb{F} = \mathbb{R}$ , then  $\mathbf{x} \perp \mathbf{y} \Leftrightarrow \|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2$ ;

(b) If  $\mathbb{F} = \mathbb{C}$ , then  $\mathbf{x} \perp \mathbf{y} \Rightarrow \|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2$ .

**Proof of Theorem 6.3.2:** Both statements derive from the following equation:

$$\begin{aligned} \|\mathbf{x} - \mathbf{y}\|^2 &= \langle (\mathbf{x} - \mathbf{y}), (\mathbf{x} - \mathbf{y}) \rangle \\ &= \langle \mathbf{x}, \mathbf{x} \rangle + \langle \mathbf{y}, \mathbf{y} \rangle - \langle \mathbf{x}, \mathbf{y} \rangle - \langle \mathbf{y}, \mathbf{x} \rangle \\ &= \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - 2 \operatorname{Re}(\langle \mathbf{x}, \mathbf{y} \rangle). \end{aligned} \quad (6.5)$$

In the case  $\mathbb{F} = \mathbb{R}$  holds  $\operatorname{Re}(\langle \mathbf{x}, \mathbf{y} \rangle) = \langle \mathbf{x}, \mathbf{y} \rangle$ , so Part (a) follows. If  $\mathbb{F} = \mathbb{C}$ , then  $\langle \mathbf{x}, \mathbf{y} \rangle$  implies  $\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2$ , so Part (b) follows. (Notice that the converse statement is not true in the case  $\mathbb{F} = \mathbb{C}$ , since Eq. (6.5) together with the hypothesis  $\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2$  do not fix  $\operatorname{Im}(\langle \mathbf{x}, \mathbf{y} \rangle)$ .) This establishes the Theorem.  $\square$

In the case of real vector space the Cauchy-Schwarz inequality stated in Theorem 6.2.2 allows us to define the angle between vectors.

**Definition 6.3.3.** The **angle** between two vectors  $\mathbf{x}, \mathbf{y}$  in a real vector inner product space  $(V, \langle \cdot, \cdot \rangle)$  is the number  $\theta \in [0, \pi]$  solution of

$$\cos(\theta) = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|}.$$

**EXAMPLE 6.3.1:** Consider the inner product space  $(\mathbb{R}^{2,2}, \langle \cdot, \cdot \rangle_F)$  and show that the following matrices are orthogonal,

$$\mathbf{A} = \begin{bmatrix} 1 & 3 \\ -1 & 4 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} -5 & 2 \\ 5 & 1 \end{bmatrix}.$$

**SOLUTION:** Since we need to compute the Frobenius inner product  $\langle \mathbf{A}, \mathbf{B} \rangle_F$ , we first compute the matrix

$$\mathbf{A}^T \mathbf{B} = \begin{bmatrix} 1 & -1 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} -5 & 2 \\ 5 & 1 \end{bmatrix} = \begin{bmatrix} -10 & 1 \\ 5 & 10 \end{bmatrix}.$$

Therefore  $\langle \mathbf{A}, \mathbf{B} \rangle_F = \operatorname{tr}(\mathbf{A}^T \mathbf{B}) = 0$ , so we conclude that  $\mathbf{A} \perp \mathbf{B}$ .  $\triangleleft$

**EXAMPLE 6.3.2:** Consider the vector space  $V = C^\infty([-\ell, \ell], \mathbb{R})$  with the inner product

$$\langle \mathbf{f}, \mathbf{g} \rangle = \int_{-\ell}^{\ell} \mathbf{f}(x) \mathbf{g}(x) dx.$$

Consider the functions  $\mathbf{u}_n(x) = \cos\left(\frac{n\pi x}{\ell}\right)$  and  $\mathbf{v}_m(x) = \sin\left(\frac{m\pi x}{\ell}\right)$ , where  $n, m$  are integers.

- (a) Show that  $\mathbf{u}_n \perp \mathbf{v}_m$  for all  $n, m$ .  
 (b) Show that  $\mathbf{u}_n \perp \mathbf{u}_m$  for all  $n \neq m$ .

(c) Show that  $\mathbf{v}_n \perp \mathbf{v}_m$  for all  $n \neq m$ .

**SOLUTION:** Recall the identities

$$\sin(\theta) \cos(\phi) = \frac{1}{2} [\sin(\theta - \phi) + \sin(\theta + \phi)], \quad (6.6)$$

$$\cos(\theta) \cos(\phi) = \frac{1}{2} [\cos(\theta - \phi) + \cos(\theta + \phi)], \quad (6.7)$$

$$\sin(\theta) \sin(\phi) = \frac{1}{2} [\cos(\theta - \phi) - \cos(\theta + \phi)]. \quad (6.8)$$

**Part (a):** Using identity in Eq. (6.6) is simple to show that

$$\begin{aligned} \langle \mathbf{u}_n, \mathbf{v}_m \rangle &= \int_{-\ell}^{\ell} \cos\left(\frac{n\pi x}{\ell}\right) \sin\left(\frac{m\pi x}{\ell}\right) dx \\ &= \frac{1}{2} \int_{-\ell}^{\ell} \left[ \sin\left(\frac{(n-m)\pi x}{\ell}\right) + \sin\left(\frac{(n+m)\pi x}{\ell}\right) \right]. \end{aligned} \quad (6.9)$$

First, assume that both  $n - m$  and  $n + m$  are non-zero,

$$\langle \mathbf{u}_n, \mathbf{v}_m \rangle = -\frac{1}{2} \left[ \frac{\ell}{(n-m)\pi} \cos\left(\frac{(n-m)\pi x}{\ell}\right) \Big|_{-\ell}^{\ell} + \frac{\ell}{(n+m)\pi} \cos\left(\frac{(n+m)\pi x}{\ell}\right) \Big|_{-\ell}^{\ell} \right]. \quad (6.10)$$

Since  $\cos((n \pm m)\pi) = \cos(-(n \pm m)\pi)$ , we conclude that both terms above vanish.

Second, in the case that  $n - m = 0$  the first term in Eq. (6.9) vanishes identically and we need to compute the term with  $(n + m)$ , which also vanishes by the second term in Eq. (6.10). Analogously, in the case of  $(n + m) = 0$  the second term in Eq. (6.9) vanishes identically and we need to compute the term with  $(n - m)$  which also vanishes by the first term in Eq. (6.10). Therefore,  $\langle \mathbf{u}_n, \mathbf{v}_m \rangle = 0$  for all  $n, m$  integers, and so  $\mathbf{u}_n \perp \mathbf{v}_m$  in this case.

**Part (b):** Using identity in Eq. (6.7) is simple to show that

$$\begin{aligned} \langle \mathbf{u}_n, \mathbf{u}_m \rangle &= \int_{-\ell}^{\ell} \cos\left(\frac{n\pi x}{\ell}\right) \cos\left(\frac{m\pi x}{\ell}\right) dx \\ &= \frac{1}{2} \int_{-\ell}^{\ell} \left[ \cos\left(\frac{(n-m)\pi x}{\ell}\right) + \cos\left(\frac{(n+m)\pi x}{\ell}\right) \right]. \end{aligned} \quad (6.11)$$

We know that  $n - m$  is non-zero. Now, assume that  $n + m$  is non-zero, then

$$\begin{aligned} \langle \mathbf{u}_n, \mathbf{u}_m \rangle &= \frac{1}{2} \left[ \frac{\ell}{(n-m)\pi} \sin\left(\frac{(n-m)\pi x}{\ell}\right) \Big|_{-\ell}^{\ell} + \frac{\ell}{(n+m)\pi} \sin\left(\frac{(n+m)\pi x}{\ell}\right) \Big|_{-\ell}^{\ell} \right] \\ &= \frac{\ell}{(n-m)\pi} \sin((n-m)\pi) + \frac{\ell}{(n+m)\pi} \sin((n+m)\pi). \end{aligned} \quad (6.12)$$

Since  $\sin((n \pm m)\pi) = 0$  for  $(n \pm m)$  integer, we conclude that both terms above vanish.

In the case of  $(n + m) = 0$  the second term in Eq. (6.11) vanishes identically and we need to compute the term with  $(n - m)$  which also vanishes by the first term in Eq. (6.12). Therefore,  $\langle \mathbf{u}_n, \mathbf{u}_m \rangle = 0$  for all  $n \neq m$  integers, and so  $\mathbf{u}_n \perp \mathbf{u}_m$  in this case.

**Part (c):** Using identity in Eq. (6.8) is simple to show that

$$\begin{aligned} \langle \mathbf{v}_n, \mathbf{v}_m \rangle &= \int_{-\ell}^{\ell} \sin\left(\frac{n\pi x}{\ell}\right) \sin\left(\frac{m\pi x}{\ell}\right) dx \\ &= \frac{1}{2} \int_{-\ell}^{\ell} \left[ \cos\left(\frac{(n-m)\pi x}{\ell}\right) - \cos\left(\frac{(n+m)\pi x}{\ell}\right) \right]. \end{aligned} \quad (6.13)$$

Since the only difference between Eq. (6.13) and (6.11) is the sign of the second term, repeating the argument done in case (b) we conclude that  $\langle \mathbf{v}_n, \mathbf{v}_m \rangle = 0$  for all  $n \neq m$  integers, and so  $\mathbf{v}_n \perp \mathbf{v}_m$  in this case.  $\triangleleft$

**6.3.2. Orthonormal basis.** We saw that an important property of a basis is that every vector in a vector space can be decomposed in a unique way in terms of the basis vectors. This decomposition is particularly simple to find in an inner product space when the basis is an orthonormal basis. Before we introduce such basis we define an orthonormal set.

**Definition 6.3.4.** The set  $U = \{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ ,  $p \geq 1$ , in an inner product space  $(V, \langle \cdot, \cdot \rangle)$  is called an **orthonormal set** iff for all  $i, j = 1, \dots, p$  holds

$$\langle \mathbf{u}_i, \mathbf{u}_j \rangle = \begin{cases} 0 & \text{if } i \neq j, \\ 1 & \text{if } i = j. \end{cases}$$

The set  $U$  is called an **orthogonal set** iff holds that  $\langle \mathbf{u}_i, \mathbf{u}_j \rangle = 0$  if  $i \neq j$  and  $\langle \mathbf{u}_i, \mathbf{u}_i \rangle \neq 0$ .

**EXAMPLE 6.3.3:** Consider the vector space  $V = C^\infty([-l, l], \mathbb{R})$  with the inner product

$$\langle \mathbf{f}, \mathbf{g} \rangle = \int_{-l}^l \mathbf{f}(x)\mathbf{g}(x) dx.$$

Show that the set

$$U = \left\{ \mathbf{u}_0 = \frac{1}{\sqrt{2l}}, \mathbf{u}_n(x) = \frac{1}{\sqrt{l}} \cos\left(\frac{n\pi x}{l}\right), \mathbf{v}_m(x) = \frac{1}{\sqrt{l}} \sin\left(\frac{m\pi x}{l}\right) \right\}_{n=1}^{\infty}$$

is an orthonormal set.

**SOLUTION:** We have shown in Example 6.3.2 that  $U$  is an orthogonal set. We only need to compute the norm of the vectors  $\mathbf{u}_0$ ,  $\mathbf{u}_n$  and  $\mathbf{v}_n$ , for  $n = 1, 2, \dots$ . The norm of the first vector is simple to compute,

$$\langle \mathbf{u}_0, \mathbf{u}_0 \rangle = \int_{-l}^l \frac{1}{2l} dx = 1.$$

The norm of the cosine functions is computed as follows,

$$\begin{aligned} \langle \mathbf{u}_n, \mathbf{u}_n \rangle &= \frac{1}{l} \int_{-l}^l \cos^2\left(\frac{n\pi x}{l}\right) dx \\ &= \frac{1}{2l} \int_{-l}^l \left[1 + \cos\left(\frac{2n\pi x}{l}\right)\right] dx \\ &= 1 + \frac{l}{2n\pi} \left[ \sin\left(\frac{2n\pi x}{l}\right) \Big|_{-l}^l \right] \Rightarrow \langle \mathbf{u}_n, \mathbf{u}_n \rangle = 1. \end{aligned}$$

A similar calculation for the sine functions gives the result

$$\begin{aligned} \langle \mathbf{v}_n, \mathbf{v}_n \rangle &= \frac{1}{l} \int_{-l}^l \sin^2\left(\frac{n\pi x}{l}\right) dx \\ &= \frac{1}{2l} \int_{-l}^l \left[1 - \cos\left(\frac{2n\pi x}{l}\right)\right] dx \\ &= 1 - \frac{l}{2n\pi} \left[ \sin\left(\frac{2n\pi x}{l}\right) \Big|_{-l}^l \right] \Rightarrow \langle \mathbf{v}_n, \mathbf{v}_n \rangle = 1. \end{aligned}$$

Therefore,  $U$  is an orthonormal set.  $\triangleleft$

A straightforward result is the following:

**Theorem 6.3.5.** An orthogonal set in an inner product space is linearly independent.



**Proof of Theorem 6.3.5:** Let  $U = \{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ ,  $p \geq 1$ , be an orthogonal set. The zero vector is not included since  $\langle \mathbf{u}_i, \mathbf{u}_i \rangle \neq 0$  for all  $i = 1, \dots, p$ . Let  $c_1, \dots, c_p \in \mathbb{F}$  be scalars such that

$$c_1 \mathbf{u}_1 + \dots + c_p \mathbf{u}_p = \mathbf{0}.$$

Then, for any  $\mathbf{u}_i \in U$  holds

$$0 = \langle \mathbf{u}_i, (c_1 \mathbf{u}_1 + \dots + c_p \mathbf{u}_p) \rangle = c_1 \langle \mathbf{u}_i, \mathbf{u}_1 \rangle + \dots + c_p \langle \mathbf{u}_i, \mathbf{u}_p \rangle.$$

Since the set  $U$  is orthogonal,  $\langle \mathbf{u}_i, \mathbf{u}_j \rangle = 0$  for  $i \neq j$  and  $\langle \mathbf{u}_i, \mathbf{u}_i \rangle \neq 0$ , so we conclude that

$$c_i \langle \mathbf{u}_i, \mathbf{u}_i \rangle = 0 \quad \Rightarrow \quad c_i = 0, \quad i = 1, \dots, p.$$

Therefore,  $U$  is a linearly independent set. This establishes the Theorem. □

**Definition 6.3.6.** A basis  $\mathcal{U}$  of an inner product space is called an **orthonormal basis** (**orthogonal basis**) iff the basis  $\mathcal{U}$  is an orthonormal (orthogonal) set.

**EXAMPLE 6.3.4:** Consider the inner product space  $(\mathbb{R}^2, \cdot)$ . Determine whether the following bases are orthonormal, orthogonal or neither:

$$\mathcal{S} = \left\{ \mathbf{e}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \mathbf{e}_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\}, \quad \mathcal{U} = \left\{ \mathbf{u}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \mathbf{u}_2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right\}, \quad \mathcal{V} = \left\{ \mathbf{v}_1 = \begin{bmatrix} 1 \\ 3 \end{bmatrix}, \mathbf{v}_2 = \begin{bmatrix} 3 \\ 1 \end{bmatrix} \right\}.$$

**SOLUTION:** The basis  $\mathcal{S}$  is **orthonormal**, since  $\mathbf{e}_1 \cdot \mathbf{e}_2 = 0$  and  $\mathbf{e}_1 \cdot \mathbf{e}_1 = \mathbf{e}_2 \cdot \mathbf{e}_2 = 1$ . The basis  $\mathcal{U}$  is **orthogonal** since  $\mathbf{u}_1 \cdot \mathbf{u}_2 = 0$ , but it is not orthonormal. Finally, the basis  $\mathcal{V}$  is **neither orthonormal nor orthogonal**. ◁

**Theorem 6.3.7.** Given the set  $U = \{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ ,  $p \geq 1$ , in the inner product space  $(\mathbb{F}^n, \cdot)$ , introduce the matrix  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_p]$ . Then the following statements hold:

- (a)  $U$  is an orthonormal set iff matrix  $\mathbf{U}$  satisfies  $\mathbf{U}^* \mathbf{U} = \mathbf{I}_p$ .
- (b)  $U$  is an orthonormal basis of  $\mathbb{F}^n$  iff matrix  $\mathbf{U}$  satisfies  $\mathbf{U}^{-1} = \mathbf{U}^*$ .

Matrices satisfying the property mentioned in part (b) of Theorem 6.3.7 appear frequently in Quantum Mechanics, so they are given a name.

**Definition 6.3.8.** A matrix  $\mathbf{U} \in \mathbb{F}^{n,n}$  is called **unitary** iff holds  $\mathbf{U}^{-1} = \mathbf{U}^*$ .

These matrices are called unitary because they do not change the norm of vectors in  $\mathbb{F}^n$  equipped with the dot product. Indeed, given any  $\mathbf{x} \in \mathbb{F}^n$  with norm  $\|\mathbf{x}\|$ , the vector  $\mathbf{U}\mathbf{x} \in \mathbb{F}^n$  has the same norm, since

$$\|\mathbf{U}\mathbf{x}\|^2 = (\mathbf{U}\mathbf{x})^* (\mathbf{U}\mathbf{x}) = \mathbf{x}^* \mathbf{U}^* \mathbf{U} \mathbf{x} = \mathbf{x}^* \mathbf{U}^{-1} \mathbf{U} \mathbf{x} = \mathbf{x}^* \mathbf{x} = \|\mathbf{x}\|^2.$$

**Proof of Theorem 6.3.7:**

**Part (a):** This is proved by a straightforward computation,

$$\mathbf{U}^* \mathbf{U} = \begin{bmatrix} \mathbf{u}_1^* \\ \vdots \\ \mathbf{u}_p^* \end{bmatrix} [\mathbf{u}_1, \dots, \mathbf{u}_p] = \begin{bmatrix} \mathbf{u}_1^* \mathbf{u}_1 & \dots & \mathbf{u}_1^* \mathbf{u}_p \\ \vdots & & \vdots \\ \mathbf{u}_p^* \mathbf{u}_1 & \dots & \mathbf{u}_p^* \mathbf{u}_p \end{bmatrix} = \begin{bmatrix} \mathbf{u}_1 \cdot \mathbf{u}_1 & \dots & \mathbf{u}_1 \cdot \mathbf{u}_p \\ \vdots & & \vdots \\ \mathbf{u}_p \cdot \mathbf{u}_1 & \dots & \mathbf{u}_p \cdot \mathbf{u}_p \end{bmatrix} = \mathbf{I}_p.$$

**Part (b):** It follows from part (a): If  $U$  is a basis of  $\mathbb{F}^n$ , then  $p = n$ ; since  $U$  is an orthonormal set, part (a) implies  $\mathbf{U}^* \mathbf{U} = \mathbf{I}_n$ . Since  $\mathbf{U}$  is an  $n \times n$  matrix, it follows that  $\mathbf{U}^* = \mathbf{U}^{-1}$ . This establishes the Theorem. □

**EXAMPLE 6.3.5:** Consider  $\mathbf{v}_1 = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}$ ,  $\mathbf{v}_2 = \begin{bmatrix} 2 \\ 0 \\ -1 \end{bmatrix}$ , in the inner product space  $(\mathbb{R}^3, \cdot)$ .

- (a) Show that  $\mathbf{v}_1 \perp \mathbf{v}_2$ ;

- (b) Find  $\mathbf{x} \in \mathbb{R}^3$  such that  $\mathbf{x} \perp \mathbf{v}_1$  and  $\mathbf{x} \perp \mathbf{v}_2$ .  
 (c) Rescale the elements of  $\{\mathbf{v}_1, \mathbf{v}_2, \mathbf{x}\}$  so that the new set is an orthonormal set.

**SOLUTION:**

Part (a):

$$\begin{bmatrix} 1 & 1 & 2 \end{bmatrix} \begin{bmatrix} 2 \\ 0 \\ -1 \end{bmatrix} = 2 + 0 - 2 = 0 \quad \Rightarrow \quad \mathbf{v}_1 \perp \mathbf{v}_2.$$

Part (b): We need to find  $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$  such that

$$\mathbf{v}_1 \cdot \mathbf{x} = \begin{bmatrix} 1 & 1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 0, \quad \mathbf{v}_2 \cdot \mathbf{x} = \begin{bmatrix} 2 & 0 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 0,$$

The equations above can be written in matrix notation as

$$\mathbf{A}^T \mathbf{x} = \mathbf{0}, \quad \text{where } \mathbf{A} = [\mathbf{v}_1, \mathbf{v}_2] = \begin{bmatrix} 1 & 2 \\ 1 & 0 \\ 2 & -1 \end{bmatrix}.$$

Gauss elimination operation on  $\mathbf{A}^T$  imply

$$\begin{bmatrix} 1 & 1 & 2 \\ 2 & 0 & -1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & -1/2 \\ 0 & 1 & 5/2 \end{bmatrix} \Rightarrow \begin{cases} x_1 = \frac{1}{2}x_3, \\ x_2 = -\frac{5}{2}x_3, \\ x_3 \text{ free.} \end{cases}$$

There is a solution for any choice of  $x_3 \neq 0$ , so we choose  $x_3 = 2$ , that is,  $\mathbf{x} = \begin{bmatrix} 1 \\ -5 \\ 2 \end{bmatrix}$ .

Part (c): The vectors  $\mathbf{v}_1$ ,  $\mathbf{v}_2$  and  $\mathbf{x}$  are mutually orthogonal. Their norms are:

$$\|\mathbf{v}_1\| = \sqrt{6}, \quad \|\mathbf{v}_2\| = \sqrt{5}, \quad \|\mathbf{x}\| = \sqrt{30}.$$

Therefore, the orthonormal set is

$$\left\{ \mathbf{u}_1 = \frac{1}{\sqrt{6}} \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}, \mathbf{u}_2 = \frac{1}{\sqrt{5}} \begin{bmatrix} 2 \\ 0 \\ -1 \end{bmatrix}, \mathbf{u}_3 = \frac{1}{\sqrt{30}} \begin{bmatrix} 1 \\ -5 \\ 2 \end{bmatrix} \right\}.$$

Finally, notice that the inverse of the matrix

$$\mathbf{U} = \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{30}} \\ \frac{1}{\sqrt{6}} & 0 & -\frac{5}{\sqrt{30}} \\ \frac{2}{\sqrt{6}} & -\frac{1}{\sqrt{5}} & \frac{2}{\sqrt{30}} \end{bmatrix} \quad \text{is} \quad \mathbf{U}^{-1} = \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{6}} & \frac{2}{\sqrt{6}} \\ \frac{2}{\sqrt{5}} & 0 & -\frac{1}{\sqrt{5}} \\ \frac{1}{\sqrt{30}} & -\frac{5}{\sqrt{30}} & \frac{2}{\sqrt{30}} \end{bmatrix} = \mathbf{U}^T.$$

◁

**6.3.3. Vector components.** Given a basis in a finite dimensional vector space, we know that every vector in the vector space can be decomposed in a unique way as a linear combination of the basis vectors. What we do not know is a general formula to compute the vector components in such a basis. However, in the particular case that the vector space admits an inner product and the basis is an orthonormal basis, we have such formula for the vector components. The formula is very simple and is the main result in the following statement.

**Theorem 6.3.9.** If  $(V, \langle \cdot, \cdot \rangle)$  is an  $n$ -dimensional inner product space with an orthonormal basis  $\mathcal{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ , then every vector  $\mathbf{x} \in V$  can be decomposed as

$$\mathbf{x} = \langle \mathbf{u}_1, \mathbf{x} \rangle \mathbf{u}_1 + \dots + \langle \mathbf{u}_n, \mathbf{x} \rangle \mathbf{u}_n. \quad (6.14)$$

**Proof of Theorem 6.3.9:** Since  $\mathcal{U}$  is a basis, we know that for all  $\mathbf{x} \in V$  there exist scalars  $c_1, \dots, c_n$  such that

$$\mathbf{x} = c_1 \mathbf{u}_1 + \dots + c_n \mathbf{u}_n.$$

Therefore, the inner product  $\langle \mathbf{u}_i, \mathbf{x} \rangle$  for any  $i = 1, \dots, n$  is given by

$$\langle \mathbf{u}_i, \mathbf{x} \rangle = c_1 \langle \mathbf{u}_i, \mathbf{u}_1 \rangle + \dots + c_n \langle \mathbf{u}_i, \mathbf{u}_n \rangle.$$

Since  $\mathcal{U}$  is an orthonormal set,  $\langle \mathbf{u}_i, \mathbf{x} \rangle = c_i$ . This establishes the Theorem.  $\square$

The result in Theorem 6.3.9 provides a remarkable simple formula for vector components in an orthonormal basis. We will get back to this subject in some depth in Section 7.1. In that Section we will name the coefficients  $\langle \mathbf{u}_i, \mathbf{x} \rangle$  in Eq. (6.14) as *Fourier coefficients* of the vector  $\mathbf{x}$  in the orthonormal set  $\mathcal{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ . When the set  $\mathcal{U}$  is an orthonormal ordered basis of a vector space  $V$ , the coordinate map  $[\ ]_{\mathcal{U}} : V \rightarrow \mathbb{F}^n$  is expressed in terms of the Fourier coefficients as follows,

$$[\mathbf{x}]_{\mathcal{U}} = \begin{bmatrix} \langle \mathbf{u}_1, \mathbf{x} \rangle \\ \vdots \\ \langle \mathbf{u}_n, \mathbf{x} \rangle \end{bmatrix}.$$

So, the coordinate map has a simple expression when it is defined by an orthonormal basis.

**EXAMPLE 6.3.6:** Consider the inner product space  $(\mathbb{R}^3, \cdot)$  with the standard ordered basis

$\mathcal{S}$ , and find the vector components of  $\mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$  in the orthonormal ordered basis

$$\mathcal{U} = \left( \mathbf{u}_1 = \frac{1}{\sqrt{6}} \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}, \mathbf{u}_2 = \frac{1}{\sqrt{5}} \begin{bmatrix} 2 \\ 0 \\ -1 \end{bmatrix}, \mathbf{u}_3 = \frac{1}{\sqrt{30}} \begin{bmatrix} 1 \\ -5 \\ 2 \end{bmatrix} \right).$$

**SOLUTION:** The vector components of  $\mathbf{x}$  in the orthonormal basis  $\mathcal{U}$  are given by

$$\mathbf{x}_{\mathcal{U}} = \begin{bmatrix} \langle \mathbf{u}_1, \mathbf{x} \rangle \\ \langle \mathbf{u}_2, \mathbf{x} \rangle \\ \langle \mathbf{u}_3, \mathbf{x} \rangle \end{bmatrix} \Rightarrow \mathbf{u} = \begin{bmatrix} \frac{9}{\sqrt{6}} \\ -\frac{1}{\sqrt{5}} \\ -\frac{3}{\sqrt{30}} \end{bmatrix}.$$

**REMARK:** We have done a change of basis, from the standard basis  $\mathcal{S}$  to the  $\mathcal{U}$  basis. In fact, we can express the calculation above as follows

$$\mathbf{x}_{\mathcal{U}} = \mathbf{P}^{-1} \mathbf{x}_{\mathcal{S}}, \quad \text{where } \mathbf{P} = \mathbf{I}_{\mathcal{U}\mathcal{S}} = \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{30}} \\ \frac{1}{\sqrt{6}} & 0 & -\frac{5}{\sqrt{30}} \\ \frac{2}{\sqrt{6}} & -\frac{1}{\sqrt{5}} & \frac{2}{\sqrt{30}} \end{bmatrix} = \mathbf{U}.$$

Since  $\mathcal{U}$  is an orthonormal basis,  $\mathbf{U}^{-1} = \mathbf{U}^T$ , so we conclude that

$$\mathbf{x}_{\mathcal{U}} = \mathbf{U}^T \mathbf{x}_{\mathcal{S}} = \begin{bmatrix} \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{6}} & \frac{2}{\sqrt{6}} \\ \frac{2}{\sqrt{5}} & 0 & -\frac{1}{\sqrt{5}} \\ \frac{1}{\sqrt{30}} & -\frac{5}{\sqrt{30}} & \frac{2}{\sqrt{30}} \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} \frac{9}{\sqrt{6}} \\ -\frac{1}{\sqrt{5}} \\ -\frac{3}{\sqrt{30}} \end{bmatrix}.$$

$\triangleleft$

## 6.3.4. Exercises.

**6.3.1.-** Prove that the following form of the Pythagoras Theorem holds on complex vector spaces: Two vectors  $\mathbf{x}$ ,  $\mathbf{y}$  in an inner product space  $(V, \langle \cdot, \cdot \rangle)$  over  $\mathbb{C}$  are orthogonal iff for all  $a, b \in \mathbb{C}$  holds

$$\|a\mathbf{x} + b\mathbf{y}\|^2 = \|a\mathbf{x}\|^2 + \|b\mathbf{y}\|^2.$$

**6.3.2.-** Consider the vector space  $\mathbb{R}^3$  with the dot product. Find all vectors  $\mathbf{x} \in \mathbb{R}^3$  which are orthogonal to the vector

$$\mathbf{v} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}.$$

**6.3.3.-** Consider the vector space  $\mathbb{R}^3$  with the dot product. Find all vectors  $\mathbf{x} \in \mathbb{R}^3$  which are orthogonal to the vectors

$$\mathbf{v}_1 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}.$$

**6.3.4.-** Let  $\mathbb{P}_3([-1, 1])$  be the space of polynomials up to degree three defined on the interval  $[-1, 1] \subset \mathbb{R}$  with the inner product

$$\langle \mathbf{p}, \mathbf{q} \rangle = \int_{-1}^1 \mathbf{p}(x)\mathbf{q}(x) dx.$$

Show that the set  $(\mathbf{p}_0, \mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3)$  is an orthogonal basis of  $\mathbb{P}_3$ , where

$$\begin{aligned} \mathbf{p}_0(x) &= 1, \\ \mathbf{p}_1(x) &= x, \\ \mathbf{p}_2(x) &= \frac{1}{2}(3x^2 - 1), \\ \mathbf{p}_3(x) &= \frac{1}{2}(5x^3 - 3x). \end{aligned}$$

(These polynomials are the first four of the Legendre polynomials.)

**6.3.5.-** Consider the vector space  $\mathbb{R}^3$  with the dot product.

(a) Show that the following ordered basis  $\mathcal{U}$  is orthonormal,

$$\left( \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix}, \frac{1}{\sqrt{3}} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \frac{1}{\sqrt{6}} \begin{bmatrix} -1 \\ -1 \\ 2 \end{bmatrix} \right).$$

(b) Use part (a) to find the components in the ordered basis  $\mathcal{U}$  of the vector

$$\mathbf{x} = \begin{bmatrix} 1 \\ 0 \\ -2 \end{bmatrix}.$$

**6.3.6.-** Consider the vector space  $\mathbb{R}^{2,2}$  with the Frobenius inner product.

(a) Show that the ordered basis given by  $\mathcal{U} = (\mathbf{E}_1, \mathbf{E}_2, \mathbf{E}_3, \mathbf{E}_4)$  is orthonormal, where

$$\mathbf{E}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathbf{E}_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$

$$\mathbf{E}_3 = \frac{1}{2} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{E}_4 = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}.$$

(b) Use part (a) to find the components in the ordered basis  $\mathcal{U}$  of the matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}.$$

**6.3.7.-** Consider the inner product space  $(\mathbb{R}^{2,2}, \langle \cdot, \cdot \rangle_F)$ , and find the cosine of the angle between the matrices

$$\mathbf{A} = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 2 & -2 \\ 2 & 0 \end{bmatrix}.$$

**6.3.8.-** Find the third column in matrix  $\mathbf{U}$  below such that  $\mathbf{U}^T = \mathbf{U}^{-1}$ , where

$$\mathbf{U} = \begin{bmatrix} 1/\sqrt{3} & 1/\sqrt{14} & U_{13} \\ 1/\sqrt{3} & 2/\sqrt{14} & U_{23} \\ 1/\sqrt{3} & -3/\sqrt{14} & U_{33} \end{bmatrix}.$$

## 6.4. ORTHOGONAL PROJECTIONS

**6.4.1. Orthogonal projection onto subspaces.** Given any subspace of an inner product space, every vector in the vector space can be decomposed as a sum of two vectors; a vector in the subspace and a vector perpendicular to the subspace. The picture one often has in mind is a plane in  $\mathbb{R}^3$  equipped with the dot product, as in Fig. 43. Any vector in  $\mathbb{R}^3$  can be decomposed as a vector on the plane plus a vector perpendicular to the plane. In this Section we provide expressions to compute this type of decompositions. We start splitting a vector in orthogonal components with respect to a one dimensional subspace. The study of this simple case describes the main ideas and the main notation used in orthogonal decompositions. Later on we present the decomposition of a vector onto an  $n$ -dimensional subspace.

**Theorem 6.4.1.** Fix a vector  $\mathbf{u} \neq \mathbf{0}$  in an inner product space  $(V, \langle \cdot, \cdot \rangle)$ . Given any vector  $\mathbf{x} \in V$  decompose it as  $\mathbf{x} = \mathbf{x}_\parallel + \mathbf{x}_\perp$  where  $\mathbf{x}_\parallel \in \text{Span}(\{\mathbf{u}\})$ . Then,  $\mathbf{x}_\perp \perp \mathbf{u}$  iff holds

$$\mathbf{x}_\parallel = \frac{\langle \mathbf{u}, \mathbf{x} \rangle}{\|\mathbf{u}\|^2} \mathbf{u}. \quad (6.15)$$

Furthermore, in the case that  $\mathbf{u}$  is a unit vector holds  $\mathbf{x}_\parallel = \langle \mathbf{u}, \mathbf{x} \rangle \mathbf{u}$ .

The main idea of this decomposition can be understood in the inner product space  $(\mathbb{R}^2, \cdot)$  and it is sketched in Fig. 42. It is obvious that every vector  $\mathbf{x}$  can be expressed in many different ways as the sum of two vectors. What is special of the decomposition in Eq. 6.15 is that  $\mathbf{x}_\parallel$  has the precise length such that  $\mathbf{x}_\perp$  is orthogonal to  $\mathbf{x}_\parallel$  (see Fig. 42).

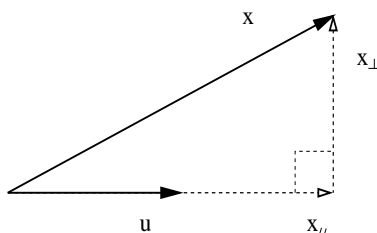


FIGURE 42. Orthogonal decomposition of the vector  $\mathbf{x} \in \mathbb{R}^2$  onto the subspace spanned by vector  $\mathbf{u}$ .

**Proof of Theorem 6.4.1:** Since  $\mathbf{x}_\parallel \in \text{Span}(\{\mathbf{u}\})$ , there exists a scalar  $a$  such that  $\mathbf{x}_\parallel = a\mathbf{u}$ . Therefore  $\mathbf{u} \perp \mathbf{x}_\perp$  iff holds that  $\langle \mathbf{u}, \mathbf{x}_\perp \rangle = 0$ . A straightforward computation shows,

$$0 = \langle \mathbf{u}, \mathbf{x}_\perp \rangle = \langle \mathbf{u}, \mathbf{x} \rangle - \langle \mathbf{u}, \mathbf{x}_\parallel \rangle = \langle \mathbf{u}, \mathbf{x} \rangle - a \langle \mathbf{u}, \mathbf{u} \rangle \quad \Leftrightarrow \quad a = \frac{\langle \mathbf{u}, \mathbf{x} \rangle}{\|\mathbf{u}\|^2}.$$

We conclude that the decomposition  $\mathbf{x} = \mathbf{x}_\parallel + \mathbf{x}_\perp$  satisfies

$$\mathbf{x}_\perp \perp \mathbf{u} \quad \Leftrightarrow \quad \mathbf{x}_\parallel = \frac{\langle \mathbf{u}, \mathbf{x} \rangle}{\|\mathbf{u}\|^2} \mathbf{u}.$$

In the case that  $\mathbf{u}$  is a unit vector holds  $\|\mathbf{u}\| = 1$ , so  $\mathbf{x}_\parallel$  is given by  $\mathbf{x}_\parallel = \langle \mathbf{u}, \mathbf{x} \rangle \mathbf{u}$ . This establishes the Theorem.  $\square$

**EXAMPLE 6.4.1:** Consider the inner product space  $(\mathbb{R}^3, \cdot)$  and decompose the vector  $\mathbf{x}$  in orthogonal components with respect to the vector  $\mathbf{u}$ , where

$$\mathbf{x} = \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}.$$

**SOLUTION:** We first compute  $\mathbf{x}_\parallel = \frac{\mathbf{u} \cdot \mathbf{x}}{\|\mathbf{u}\|^2} \mathbf{u}$ . Since

$$\mathbf{u} \cdot \mathbf{x} = \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix} = 3 + 4 + 3 = 10, \quad \|\mathbf{u}\|^2 = \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = 1 + 4 + 9 = 14,$$

we obtain  $\mathbf{x}_\parallel = \frac{5}{7} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$ . We now compute  $\mathbf{x}_\perp$  as follows,

$$\mathbf{x}_\perp = \mathbf{x} - \mathbf{x}_\parallel = \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix} - \frac{5}{7} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = \frac{1}{7} \begin{bmatrix} 21 \\ 14 \\ 7 \end{bmatrix} - \frac{1}{7} \begin{bmatrix} 5 \\ 10 \\ 15 \end{bmatrix} = \frac{1}{7} \begin{bmatrix} 16 \\ 4 \\ -8 \end{bmatrix} \Rightarrow \mathbf{x}_\perp = \frac{4}{7} \begin{bmatrix} 4 \\ 1 \\ -2 \end{bmatrix}.$$

Therefore,  $\mathbf{x}$  can be decomposed as

$$\mathbf{x} = \frac{5}{7} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + \frac{4}{7} \begin{bmatrix} 4 \\ 1 \\ -2 \end{bmatrix}.$$

**REMARK:** We can verify that this decomposition is orthogonal with respect to  $\mathbf{u}$ , since

$$\mathbf{u} \cdot \mathbf{x}_\perp = \frac{4}{7} \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 4 \\ 1 \\ -2 \end{bmatrix} = \frac{4}{7} (4 + 2 - 6) = 0.$$

◁

We now decompose a vector into orthogonal components with respect to a  $p$ -dimensional subspace with  $p \geq 1$ .

**Theorem 6.4.2.** Fix an orthogonal set  $\mathcal{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ , with  $p \geq 1$ , in an inner product space  $(V, \langle \cdot, \cdot \rangle)$ . Given any vector  $\mathbf{x} \in V$ , decompose it as  $\mathbf{x} = \mathbf{x}_\parallel + \mathbf{x}_\perp$ , where  $\mathbf{x}_\parallel \in \text{Span}(\mathcal{U})$ . Then,  $\mathbf{x}_\perp \perp \mathbf{u}_i$ , for  $i = 1, \dots, p$  iff holds

$$\mathbf{x}_\parallel = \frac{\langle \mathbf{u}_1, \mathbf{x} \rangle}{\|\mathbf{u}_1\|^2} \mathbf{u}_1 + \dots + \frac{\langle \mathbf{u}_p, \mathbf{x} \rangle}{\|\mathbf{u}_p\|^2} \mathbf{u}_p. \quad (6.16)$$

Furthermore, in the case that  $\mathcal{U}$  is an orthonormal basis holds

$$\mathbf{x}_\parallel = \langle \mathbf{u}_1, \mathbf{x} \rangle \mathbf{u}_1 + \dots + \langle \mathbf{u}_p, \mathbf{x} \rangle \mathbf{u}_p. \quad (6.17)$$

**REMARK:** The set  $\mathcal{U}$  in Theorem 6.4.2 must be an orthogonal set. If the vectors  $\mathbf{u}_1, \dots, \mathbf{u}_p$  are not mutually orthogonal, then the vector  $\mathbf{x}_\parallel$  computed in Eq. (6.17) is not the orthogonal projection of vector  $\mathbf{x}$ , that is,  $(\mathbf{x} - \mathbf{x}_\parallel) \not\perp \mathbf{u}_i$  for  $i = 1, \dots, p$ . Therefore, before using Eq. (6.17) in a particular application one should verify that the set  $\mathcal{U}$  one is working with is in fact an orthogonal set. A particular case of the orthogonal projection of a vector in the inner product space  $(\mathbb{R}^3, \cdot)$  onto a plane is sketched in Fig. 43.

**Proof of Theorem 6.4.2:** Since  $\mathbf{x}_\parallel \in \text{Span}(\mathcal{U})$ , there exist scalars  $a_i$ , for  $i = 1, \dots, p$  such that  $\mathbf{x}_\parallel = a_1 \mathbf{u}_1 + \dots + a_p \mathbf{u}_p$ . The vector  $\mathbf{x}_\perp \perp \mathbf{u}_i$  iff holds that  $\langle \mathbf{u}_i, \mathbf{x}_\perp \rangle = 0$ . A straightforward computation shows that, for  $i = 1, \dots, p$  holds

$$\begin{aligned} 0 = \langle \mathbf{u}_i, \mathbf{x}_\perp \rangle &= \langle \mathbf{u}_i, \mathbf{x} \rangle - \langle \mathbf{u}_i, \mathbf{x}_\parallel \rangle \\ &= \langle \mathbf{u}_i, \mathbf{x} \rangle - a_1 \langle \mathbf{u}_i, \mathbf{u}_1 \rangle - \dots - a_p \langle \mathbf{u}_i, \mathbf{u}_p \rangle \\ &= \langle \mathbf{u}_i, \mathbf{x} \rangle - a_i \langle \mathbf{u}_i, \mathbf{u}_i \rangle \Leftrightarrow a_i = \frac{\langle \mathbf{u}_i, \mathbf{x} \rangle}{\|\mathbf{u}_i\|^2}. \end{aligned}$$

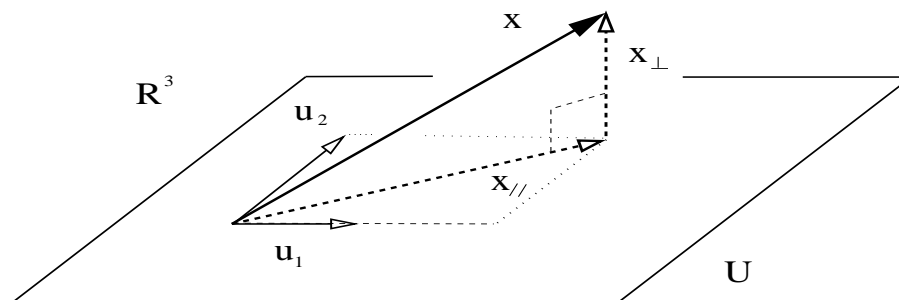


FIGURE 43. Orthogonal decomposition of the vector  $\mathbf{x} \in \mathbb{R}^3$  onto the subspace  $U$  spanned by the vectors  $\mathbf{u}_1$  and  $\mathbf{u}_2$ .

We conclude that the decomposition  $\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_\perp$  satisfies  $\mathbf{x}_\perp \perp \mathbf{u}_i$  for  $i = 1, \dots, p$  iff holds

$$\mathbf{x}_1 = \frac{\langle \mathbf{u}_1, \mathbf{x} \rangle}{\|\mathbf{u}_1\|^2} \mathbf{u}_1 + \dots + \frac{\langle \mathbf{u}_p, \mathbf{x} \rangle}{\|\mathbf{u}_p\|^2} \mathbf{u}_p.$$

In the case that  $\mathcal{U}$  is an orthonormal set holds  $\|\mathbf{u}_i\| = 1$  for  $i = 1, \dots, p$ , so  $\mathbf{x}_1$  is given by Eq. (6.17). This establishes the Theorem.  $\square$

**EXAMPLE 6.4.2:** Consider the inner product space  $(\mathbb{R}^3, \cdot)$  and decompose the vector  $\mathbf{x}$  in orthogonal components with respect to the subspace  $U$ , where

$$\mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \quad U = \text{Span}\left(\left\{ \mathbf{u}_1 = \begin{bmatrix} 2 \\ 5 \\ -1 \end{bmatrix}, \mathbf{u}_2 = \begin{bmatrix} -2 \\ 1 \\ 1 \end{bmatrix} \right\}\right).$$

**SOLUTION:** In order to use Eq. (6.16) we need an orthogonal basis of  $U$ . So, need need to verify whether  $\mathbf{u}_1$  is orthogonal to  $\mathbf{u}_2$ . This is indeed the case, since

$$\mathbf{u}_1 \cdot \mathbf{u}_2 = [2 \quad 5 \quad -1] \begin{bmatrix} -2 \\ 1 \\ 1 \end{bmatrix} = -4 + 5 - 1 = 0.$$

So now we use  $\mathbf{u}_1$  and  $\mathbf{u}_2$  to compute  $\mathbf{x}_1$  using Eq. (6.17). We need the quantities

$$\begin{aligned} \mathbf{u}_1 \cdot \mathbf{x} &= [2 \quad 5 \quad -1] \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = 9, & \mathbf{u}_1 \cdot \mathbf{u}_1 &= [2 \quad 5 \quad -1] \begin{bmatrix} 2 \\ 5 \\ -1 \end{bmatrix} = 30, \\ \mathbf{u}_2 \cdot \mathbf{x} &= [-2 \quad 1 \quad 1] \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = 3, & \mathbf{u}_2 \cdot \mathbf{u}_2 &= [-2 \quad 1 \quad 1] \begin{bmatrix} -2 \\ 1 \\ 1 \end{bmatrix} = 6. \end{aligned}$$

Now is simple to compute  $\mathbf{x}_1$ , since

$$\mathbf{x}_1 = \frac{9}{30} \begin{bmatrix} 2 \\ 5 \\ -1 \end{bmatrix} + \frac{3}{6} \begin{bmatrix} -2 \\ 1 \\ 1 \end{bmatrix} = \frac{3}{10} \begin{bmatrix} 2 \\ 5 \\ -1 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} -2 \\ 1 \\ 1 \end{bmatrix} = \frac{1}{10} \begin{bmatrix} 6 \\ 15 \\ -3 \end{bmatrix} + \frac{1}{10} \begin{bmatrix} -10 \\ 5 \\ 5 \end{bmatrix},$$

therefore,

$$\mathbf{x}_1 = \frac{1}{10} \begin{bmatrix} -4 \\ 20 \\ 2 \end{bmatrix} \Rightarrow \mathbf{x}_1 = \frac{1}{5} \begin{bmatrix} -2 \\ 10 \\ 1 \end{bmatrix}.$$

The vector  $\mathbf{x}_\perp$  is obtained as follows,

$$\mathbf{x}_\perp = \mathbf{x} - \mathbf{x}_{||} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} - \frac{1}{5} \begin{bmatrix} -2 \\ 10 \\ 1 \end{bmatrix} = \frac{1}{5} \begin{bmatrix} 5 \\ 10 \\ 15 \end{bmatrix} - \frac{1}{5} \begin{bmatrix} -2 \\ 10 \\ 1 \end{bmatrix} = \frac{1}{5} \begin{bmatrix} 7 \\ 0 \\ 14 \end{bmatrix} \Rightarrow \mathbf{x}_\perp = \frac{7}{5} \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix}.$$

We conclude that

$$\mathbf{x} = \frac{1}{5} \begin{bmatrix} -2 \\ 10 \\ 1 \end{bmatrix} + \frac{7}{5} \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix}.$$

**REMARK:** We can verify that  $\mathbf{x}_\perp \perp U$ , since

$$\begin{aligned} \mathbf{u}_1 \cdot \mathbf{x}_\perp &= \frac{7}{5} [2 \quad 5 \quad -1] \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix} = \frac{7}{5} (2 + 0 - 2) = 0, \\ \mathbf{u}_2 \cdot \mathbf{x}_\perp &= \frac{7}{5} [-2 \quad 1 \quad 1] \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix} = \frac{7}{5} (-2 + 0 + 2) = 0. \end{aligned}$$

◁

**6.4.2. Orthogonal complement.** We have seen that any vector in an inner product space can be decomposed as a sum of orthogonal vectors, that is,  $\mathbf{x} = \mathbf{x}_{||} + \mathbf{x}_\perp$  with  $\langle \mathbf{x}_{||}, \mathbf{x}_\perp \rangle = 0$ . Inner product spaces can be decomposed in a somehow analogous way as a direct sum of two mutually orthogonal subspaces. In order to understand such decomposition, we need to introduce the notion of the orthogonal complement of a subspace. Then we will show that a finite dimensional inner product space is the direct sum of a subspace and its orthogonal complement. It is precisely this result that motivates the word “complement” in the name orthogonal complement of a subspace.

**Definition 6.4.3.** The *orthogonal complement* of a subspace  $W$  in an inner product space  $(V, \langle \cdot, \cdot \rangle)$ , denoted as  $W^\perp$ , is the set  $W^\perp = \{\mathbf{u} \in V : \langle \mathbf{u}, \mathbf{w} \rangle = 0 \quad \forall \mathbf{w} \in W\}$ .

**EXAMPLE 6.4.3:** In the inner product space  $(\mathbb{R}^3, \cdot)$ , the orthogonal complement to a line is a plane, and the orthogonal complement to a plane is a line, as it is shown in Fig. 44. ◁

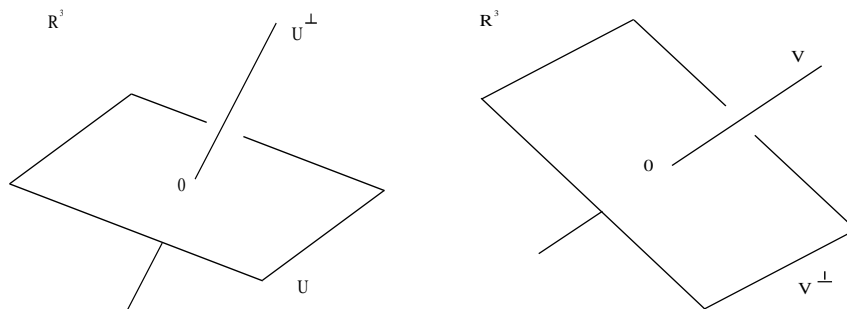


FIGURE 44. The orthogonal complement to the plane  $U$  is the line  $U^\perp$ , and the orthogonal complement to the line  $V$  is the plane  $V^\perp$ .

As the sketch in Fig. 44 suggests, the orthogonal complement of a subspace is a subspace.

**Theorem 6.4.4.** The orthogonal complement  $W^\perp$  of a subspace  $W$  in an inner product space is also a subspace.



**Proof of Theorem 6.4.4:** Let  $\mathbf{u}_1, \mathbf{u}_2 \in W^\perp$ , that is,  $\langle \mathbf{u}_i, \mathbf{w} \rangle = 0$  for all  $\mathbf{w} \in W$  and  $i = 1, 2$ . Then, any linear combination  $a\mathbf{u}_1 + b\mathbf{u}_2$  also belongs to  $W^\perp$ , since

$$\langle (a\mathbf{u}_1 + b\mathbf{u}_2), \mathbf{w} \rangle = \bar{a} \langle \mathbf{u}_1, \mathbf{w} \rangle + \bar{b} \langle \mathbf{u}_2, \mathbf{w} \rangle = 0 + 0 \quad \forall \mathbf{w} \in W.$$

This establishes the Theorem.  $\square$

**EXAMPLE 6.4.4:** Find  $W^\perp$  for the subspace  $W = \text{Span}\left(\left\{\mathbf{w}_1 = \begin{bmatrix} -1 \\ 2 \\ 3 \end{bmatrix}\right\}\right)$  in  $(\mathbb{R}^3, \cdot)$ .

**SOLUTION:** We need to find the set of all  $\mathbf{x} \in \mathbb{R}^3$  such that  $\mathbf{x} \cdot \mathbf{w}_1 = 0$ . That is,

$$\begin{bmatrix} -1 & 2 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 0 \quad \Rightarrow \quad x_1 = 2x_2 + 3x_3.$$

The solution is

$$\mathbf{x} = \begin{bmatrix} 2x_2 + 3x_3 \\ x_2 \\ x_3 \end{bmatrix} \quad \Rightarrow \quad \mathbf{x} = \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} x_2 + \begin{bmatrix} 3 \\ 0 \\ 1 \end{bmatrix} x_3,$$

hence we obtain

$$W^\perp = \text{Span}\left(\left\{\begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 3 \\ 0 \\ 1 \end{bmatrix}\right\}\right).$$

The orthogonal complement of a line is a plane, as sketched in the second picture in Fig. 44.

**REMARK:** We can verify that the result is correct, since

$$\begin{bmatrix} -1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} = -2 + 2 + 0 = 0, \quad \begin{bmatrix} -1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 3 \\ 0 \\ 1 \end{bmatrix} = -3 + 0 + 3 = 0.$$

$\triangleleft$

**EXAMPLE 6.4.5:** Find  $W^\perp$  for the subspace  $W = \text{Span}\left(\left\{\mathbf{w}_1 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \mathbf{w}_2 = \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix}\right\}\right)$  in  $(\mathbb{R}^3, \cdot)$ .

**SOLUTION:** We need to find the set of all  $\mathbf{x} \in \mathbb{R}^3$  such that  $\mathbf{x} \cdot \mathbf{w}_1 = 0$  and  $\mathbf{x} \cdot \mathbf{w}_2 = 0$ . That is,

$$\begin{bmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 0$$

We can use Gauss elimination to find the solution,

$$\begin{bmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & 2 \end{bmatrix} \quad \Rightarrow \quad \begin{cases} x_1 = x_3, \\ x_2 = -2x_3, \\ x_3 \text{ free.} \end{cases} \quad \Rightarrow \quad \mathbf{x} = \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} x_3,$$

hence we obtain

$$W^\perp = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}\right\}\right).$$

So, the orthogonal complement of a plane is a line, as sketched in the first picture in Fig. 44.

**REMARK:** We can verify that the result is correct, since

$$\begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} = 1 - 4 + 3 = 0, \quad \begin{bmatrix} 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} = 3 - 4 + 1 = 0.$$

◁

As we mentioned above, the reason for the word “complement” in the name of an orthogonal complement is that the vector space can be split into the sum of two subspaces with zero intersection. We summarize this below.

**Theorem 6.4.5.** *If  $W$  is a subspace in a finite dimensional inner product space  $V$ , then*

$$V = W \oplus W^\perp.$$

**Proof of Theorem 6.4.5:** We first show that  $V = W + W^\perp$ . In order to do that, first choose an orthonormal basis  $\mathcal{W} = \{\mathbf{w}_1, \dots, \mathbf{w}_p\}$  for  $W$ , we here we have assumed  $\dim W = p \leq n = \dim V$ . Then, for every vector  $\mathbf{x} \in V$  holds that it can be decomposed as

$$\mathbf{x} = \mathbf{x}_\parallel + \mathbf{x}_\perp, \quad \mathbf{x}_\parallel = \langle \mathbf{w}_1, \mathbf{x} \rangle \mathbf{w}_1 + \dots + \langle \mathbf{w}_p, \mathbf{x} \rangle \mathbf{w}_p, \quad \mathbf{x}_\perp = \mathbf{x} - \mathbf{x}_\parallel.$$

Theorem 6.4.2 says that  $\mathbf{x}_\perp \perp \mathbf{w}_i$  for  $i = 1, \dots, p$ . This implies that  $\mathbf{x}_\perp \in W^\perp$  and, since  $\mathbf{x}$  is an arbitrary vector in  $V$ , we have established that  $V = W + W^\perp$ . We now show that  $W \cap W^\perp = \{\mathbf{0}\}$ . Indeed, if  $\mathbf{u} \in W^\perp$ , then  $\langle \mathbf{u}, \mathbf{w} \rangle = 0$  for all  $\mathbf{w} \in W$ . If  $\mathbf{u} \in W$ , then choosing  $\mathbf{w} = \mathbf{u}$  in the equation above we get  $\langle \mathbf{u}, \mathbf{u} \rangle = 0$ , which implies  $\mathbf{u} = \mathbf{0}$ . Therefore,  $W \cap W^\perp = \{\mathbf{0}\}$  and we conclude that  $V = W \oplus W^\perp$ . This establishes the Theorem.  $\square$

Since the orthogonal complement  $W^\perp$  of a subspace  $W$  is itself a subspace, we can compute  $(W^\perp)^\perp$ . The following statement says that the result is the original subspace  $W$ .

**Theorem 6.4.6.** *If  $W$  is a subspace in an finite dimensional inner product space, then*

$$(W^\perp)^\perp = W.$$

**Proof of Theorem 6.4.6:** First notice that  $W \subset (W^\perp)^\perp$ . Indeed, given any fixed vector  $\mathbf{w} \in W$ , the definition of  $W^\perp$  implies that every vector  $\mathbf{u} \in W^\perp$  satisfies  $\langle \mathbf{w}, \mathbf{u} \rangle = 0$ . This condition says that  $\mathbf{w} \in (W^\perp)^\perp$ , hence we conclude  $W \subset (W^\perp)^\perp$ . Second, Theorem 6.4.2 says that

$$V = W^\perp \oplus (W^\perp)^\perp.$$

In particular, this decomposition implies that  $\dim V = \dim W^\perp + \dim (W^\perp)^\perp$ . Again Theorem 6.4.2 also says that

$$V = W \oplus W^\perp,$$

which in particular implies that  $\dim V = \dim W + \dim W^\perp$ . These last two results put together imply  $\dim W = \dim (W^\perp)^\perp$ . It is from this last result together with our previous result,  $W \subset (W^\perp)^\perp$ , that we obtain  $W = (W^\perp)^\perp$ . This establishes the Theorem.  $\square$

## 6.4.3. Exercises.

**6.4.1.-** Consider the inner product space  $(\mathbb{R}^3, \cdot)$  and use Theorem 6.4.1 to find the orthogonal decomposition of vector  $\mathbf{x}$  along vector  $\mathbf{u}$ , where

$$\mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \quad \mathbf{u} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}.$$

**6.4.2.-** Consider the subspace  $W$  given by

$$\text{Span}\left(\left\{\mathbf{w}_1 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}, \mathbf{w}_2 = \begin{bmatrix} 2 \\ -1 \\ 3 \end{bmatrix}\right\}\right)$$

in the inner product space  $(\mathbb{R}^3, \cdot)$ .

- Find an orthogonal decomposition of the vector  $\mathbf{w}_2$  with respect to the vector  $\mathbf{w}_1$ . Using this decomposition, find an orthogonal basis for the space  $W$ .
- Find the decomposition of the vector  $\mathbf{x}$  below in orthogonal components with respect to the subspace  $W$ , where

$$\mathbf{x} = \begin{bmatrix} 4 \\ 3 \\ 0 \end{bmatrix}.$$

**6.4.3.-** Consider the subspace  $W$  given by

$$\text{Span}\left(\left\{\mathbf{w}_1 = \begin{bmatrix} 2 \\ 0 \\ -2 \end{bmatrix}, \mathbf{w}_2 = \begin{bmatrix} 2 \\ -1 \\ 0 \end{bmatrix}\right\}\right)$$

in the inner product space  $(\mathbb{R}^3, \cdot)$ . Decompose the vector  $\mathbf{x}$  below into orthogonal components with respect to  $W$ , where

$$\mathbf{x} = \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix}.$$

(Notice that  $\mathbf{w}_1 \not\perp \mathbf{w}_2$ .)

**6.4.4.-** Given the matrix  $A$  below, find a basis for the space  $R(A)^\perp$ , where

$$A = \begin{bmatrix} 1 & 2 \\ 1 & 1 \\ 2 & 0 \end{bmatrix}.$$

**6.4.5.-** Consider the subspace

$$W = \text{Span}\left\{\begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix}\right\}$$

in the inner product space  $(\mathbb{R}^3, \cdot)$ .

- Find a basis for the space  $W^\perp$ , that is, find a basis for the orthogonal complement of the space  $W$ .
- Use Theorem 6.4.1 to transform the basis of  $W^\perp$  found in part (a) into an orthogonal basis.

**6.4.6.-** Consider the inner product space  $(\mathbb{R}^4, \cdot)$ , and find a basis for the orthogonal complement of the subspace  $W$  given by

$$W = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 2 \\ 0 \\ 3 \end{bmatrix}, \begin{bmatrix} 2 \\ 4 \\ 1 \\ 6 \end{bmatrix}\right\}\right).$$

**6.4.7.-** Let  $X$  and  $Y$  be subspaces of a finite dimensional inner product space  $(V, \langle \cdot, \cdot \rangle)$ . Prove the following:

- $X \subset Y \Rightarrow Y^\perp \subset X^\perp$ ;
- $(X + Y)^\perp = X^\perp \cap Y^\perp$ ;
- Use part (b) to show that

$$(X \cap Y)^\perp = X^\perp + Y^\perp.$$

## 6.5. GRAM-SCHMIDT METHOD

We now describe the Gram-Schmidt orthogonalization method, which is a method to transform a linearly independent set of vectors into an orthonormal set. The method is based on projecting the  $i$ -th vector in the set onto the subspace spanned by the previous  $(i - 1)$  vectors.

**Theorem 6.5.1 (Gram-Schmidt).** *Let  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_p\}$  be a linearly independent set in an inner product space  $(V, \langle \cdot, \cdot \rangle)$ . If the set  $Y = \{\mathbf{y}_1, \dots, \mathbf{y}_p\}$  is defined as follows,*

$$\begin{aligned} \mathbf{y}_1 &= \mathbf{x}_1, \\ \mathbf{y}_2 &= \mathbf{x}_2 - \frac{\langle \mathbf{y}_1, \mathbf{x}_2 \rangle}{\|\mathbf{y}_1\|^2} \mathbf{y}_1, \\ &\vdots \\ \mathbf{y}_p &= \mathbf{x}_p - \frac{\langle \mathbf{y}_1, \mathbf{x}_p \rangle}{\|\mathbf{y}_1\|^2} \mathbf{y}_1 - \dots - \frac{\langle \mathbf{y}_{(p-1)}, \mathbf{x}_p \rangle}{\|\mathbf{y}_{(p-1)}\|^2} \mathbf{y}_{(p-1)}, \end{aligned}$$

then  $Y$  is an orthogonal set with  $\text{Span}(Y) = \text{Span}(X)$ . Furthermore, the set

$$Z = \left\{ \mathbf{z}_1 = \frac{\mathbf{y}_1}{\|\mathbf{y}_1\|}, \dots, \mathbf{z}_p = \frac{\mathbf{y}_p}{\|\mathbf{y}_p\|} \right\}$$

is an orthonormal set with  $\text{Span}(Z) = \text{Span}(Y)$ .

**REMARK:** Using the notation in Theorem 6.4.2 we can write  $\mathbf{y}_2 = \mathbf{x}_{2\perp}$ , where the projection is onto the subspace  $\text{Span}(\{\mathbf{y}_1\})$ . Analogously,  $\mathbf{y}_i = \mathbf{x}_{i\perp}$ , for  $i = 2, \dots, p$ , where the projection is onto the subspace  $\text{Span}(\{\mathbf{y}_1, \dots, \mathbf{y}_{i-1}\})$ .

**Proof of Theorem 6.5.1:** We first show that  $Y$  is an orthogonal set. It is simple to see that  $\mathbf{y}_2 \in \text{Span}(\{\mathbf{y}_1\})^\perp$ , since

$$\langle \mathbf{y}_1, \mathbf{y}_2 \rangle = \langle \mathbf{y}_1, \mathbf{x}_2 \rangle - \frac{\langle \mathbf{y}_1, \mathbf{x}_2 \rangle}{\|\mathbf{y}_1\|^2} \langle \mathbf{y}_1, \mathbf{y}_1 \rangle = 0.$$

Assume that  $\mathbf{y}_i \in \text{Span}(\{\mathbf{y}_1, \dots, \mathbf{y}_{i-1}\})^\perp$ , we then show that  $\mathbf{y}_{i+1} \in \text{Span}(\{\mathbf{y}_1, \dots, \mathbf{y}_i\})^\perp$ . Indeed, for  $j = 1, \dots, i$  holds

$$\langle \mathbf{y}_j, \mathbf{y}_{i+1} \rangle = \langle \mathbf{y}_j, \mathbf{x}_{i+1} \rangle - \frac{\langle \mathbf{y}_j, \mathbf{x}_{i+1} \rangle}{\|\mathbf{y}_j\|^2} \langle \mathbf{y}_j, \mathbf{y}_j \rangle = 0,$$

where we used that  $\mathbf{y}_j \in \text{Span}(\{\mathbf{y}_1, \dots, \mathbf{y}_{i-1}\})^\perp$ , for all  $j = 1, \dots, i$ . Therefore,  $Y$  is an orthogonal set (and so, a linearly independent set).

The proof that  $\text{Span}(X) = \text{Span}(Y)$  has two steps: On the one hand, the elements in  $Y$  are linear combinations of elements in  $X$ , hence  $\text{Span}(Y) \subset \text{Span}(X)$ ; on the other hand  $\dim \text{Span}(X) = \dim \text{Span}(Y)$ , since  $X$  and  $Y$  are both linearly independent sets with the same number of elements. We conclude that  $\text{Span}(X) = \text{Span}(Y)$ . It is straightforward to see that  $Z$  is an orthonormal set, and since every element  $\mathbf{z}_i \in Z$  is proportional to every  $\mathbf{y}_i \in Y$ , then  $\text{Span}(Y) = \text{Span}(Z)$ . This establishes the Theorem.  $\square$

**EXAMPLE 6.5.1:** Use the Gram-Schmidt method to find an orthonormal basis for the inner product space  $(\mathbb{R}^3, \cdot)$  from the ordered basis

$$X = \left( \mathbf{x}_1 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \mathbf{x}_2 = \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix}, \mathbf{x}_3 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right).$$

**SOLUTION:** We first find an orthogonal basis. The first element is

$$y_1 = x_1 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \Rightarrow \|y_1\|^2 = 2.$$

The second element is

$$y_2 = x_2 - \frac{y_1 \cdot x_2}{\|y_1\|^2} y_1,$$

where

$$y_1 \cdot x_2 = [1 \ 1 \ 0] \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} = 3.$$

A simple calculation shows

$$y_2 = \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} - \frac{3}{2} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 4 \\ 2 \\ 0 \end{bmatrix} - \frac{1}{2} \begin{bmatrix} 3 \\ 3 \\ 0 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix},$$

therefore,

$$y_2 = \frac{1}{2} \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} \Rightarrow \|y_2\|^2 = \frac{1}{2}.$$

Finally, the last element is

$$y_3 = x_3 - \frac{y_1 \cdot x_3}{\|y_1\|^2} y_1 - \frac{y_2 \cdot x_3}{\|y_2\|^2} y_2,$$

where

$$y_1 \cdot x_3 = [1 \ 1 \ 0] \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = 2, \quad y_2 \cdot x_3 = \frac{1}{2} [1 \ -1 \ 0] \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = 0.$$

Another simple calculation shows

$$y_3 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - \frac{2}{2} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix},$$

therefore,

$$y_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \Rightarrow \|y_3\|^2 = 1.$$

The set  $Y = \{y_1, y_2, y_3\}$  is an orthogonal set, while an orthonormal set is given by

$$Z = \left\{ z_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, z_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix}, z_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right\}.$$

◁

**EXAMPLE 6.5.2:** Consider the vector space  $\mathbb{P}_3([-1, 1])$  with the inner product

$$\langle \mathbf{p}, \mathbf{q} \rangle = \int_{-1}^1 \mathbf{p}(x) \mathbf{q}(x) dx.$$

Given the basis  $\{\mathbf{p}_0 = 1, \mathbf{p}_1 = x, \mathbf{p}_2 = x^2, \mathbf{p}_3 = x^3\}$ , use the Gram-Schmidt method starting with the vector  $\mathbf{p}_0$  to find an orthogonal basis for  $\mathbb{P}_3([-1, 1])$ .

**SOLUTION:** The first element in the new basis is

$$\mathbf{q}_0 = \mathbf{p}_0 = 1 \quad \Rightarrow \quad \|\mathbf{q}_0\|^2 = \int_{-1}^1 dx = 2.$$

The second element is

$$\mathbf{q}_1 = \mathbf{p}_1 - \frac{\langle \mathbf{q}_0, \mathbf{p}_1 \rangle}{\|\mathbf{q}_0\|^2} \mathbf{q}_0.$$

It is simple to see that

$$\langle \mathbf{q}_0, \mathbf{p}_1 \rangle = \int_{-1}^1 x dx = \frac{1}{2} x^2 \Big|_{-1}^1 = 0.$$

So we conclude that

$$\mathbf{q}_1 = \mathbf{p}_1 = x \quad \Rightarrow \quad \|\mathbf{q}_1\|^2 = \int_{-1}^1 x^2 dx = \frac{1}{3} x^3 \Big|_{-1}^1 \quad \Rightarrow \quad \|\mathbf{q}_1\|^2 = \frac{2}{3}.$$

The third element in the basis is

$$\mathbf{q}_2 = \mathbf{p}_2 - \frac{\langle \mathbf{q}_0, \mathbf{p}_2 \rangle}{\|\mathbf{q}_0\|^2} \mathbf{q}_0 - \frac{\langle \mathbf{q}_1, \mathbf{p}_2 \rangle}{\|\mathbf{q}_1\|^2} \mathbf{q}_1.$$

It is simple to see that

$$\langle \mathbf{q}_0, \mathbf{p}_2 \rangle = \int_{-1}^1 x^2 dx = \frac{1}{3} x^3 \Big|_{-1}^1 = \frac{2}{3},$$

$$\langle \mathbf{q}_1, \mathbf{p}_2 \rangle = \int_{-1}^1 x^3 dx = \frac{1}{4} x^4 \Big|_{-1}^1 = 0.$$

Hence we obtain

$$\mathbf{q}_2 = \mathbf{p}_2 - \frac{2}{3} \frac{1}{2} \mathbf{q}_0 = x^2 - \frac{1}{3} \quad \Rightarrow \quad \mathbf{q}_2 = \frac{1}{3} (3x^2 - 1).$$

The norm square of this vector is

$$\begin{aligned} \|\mathbf{q}_2\|^2 &= \frac{1}{9} \int_{-1}^1 (3x^2 - 1)(3x^2 - 1) dx \\ &= \frac{1}{9} \int_{-1}^1 (9x^4 - 6x^2 + 1) dx \\ &= \frac{1}{9} \left( \frac{9}{5} x^5 - 2x^3 + x \right) \Big|_{-1}^1 \\ &= \frac{8}{45}. \end{aligned}$$

Finally, the fourth vector of the orthogonal basis is given by

$$\mathbf{q}_3 = \mathbf{p}_3 - \frac{\langle \mathbf{q}_0, \mathbf{p}_3 \rangle}{\|\mathbf{q}_0\|^2} \mathbf{q}_0 - \frac{\langle \mathbf{q}_1, \mathbf{p}_3 \rangle}{\|\mathbf{q}_1\|^2} \mathbf{q}_1 - \frac{\langle \mathbf{q}_2, \mathbf{p}_3 \rangle}{\|\mathbf{q}_2\|^2} \mathbf{q}_2.$$

It is simple to see that

$$\langle \mathbf{q}_0, \mathbf{p}_3 \rangle = \int_{-1}^1 x^3 dx = \frac{1}{4} x^4 \Big|_{-1}^1 = 0,$$

$$\langle \mathbf{q}_1, \mathbf{p}_3 \rangle = \int_{-1}^1 x^4 dx = \frac{1}{5} x^5 \Big|_{-1}^1 = \frac{2}{5},$$

$$\langle \mathbf{q}_2, \mathbf{p}_3 \rangle = \frac{1}{3} \int_{-1}^1 (3x^2 - 1) x^3 dx = \frac{1}{3} \left( \frac{1}{2} x^6 - \frac{1}{4} x^4 \right) \Big|_{-1}^1 = 0.$$

Hence we obtain

$$\mathbf{q}_3 = \mathbf{p}_3 - \frac{2}{5} \frac{3}{2} \mathbf{q}_1 = x^3 - \frac{3}{5}x \quad \Rightarrow \quad \mathbf{q}_3 = \frac{1}{5}(5x^3 - 3x).$$

The orthogonal basis is then given by

$$\left\{ \mathbf{q}_0 = 1, \mathbf{q}_1 = x, \mathbf{q}_2 = \frac{1}{3}(3x^2 - 1), \mathbf{q}_3 = \frac{1}{5}(5x^3 - 3x) \right\}.$$

These polynomials are proportional to the first three Legendre polynomials. The Legendre polynomials form an orthogonal set in the space  $\mathbb{P}_\infty([-1, 1])$  of polynomials of all degrees. They play an important role in physics, since Legendre polynomials are solution of a particular differential equation that often appears in physics.  $\triangleleft$

## 6.5.1. Exercises.

**6.5.1.-** Find an orthonormal basis for the subspace of  $\mathbb{R}^3$  spanned by the vectors

$$\left\{ \mathbf{u}_1 = \begin{bmatrix} -2 \\ 2 \\ -1 \end{bmatrix}, \mathbf{u}_2 = \begin{bmatrix} 1 \\ -3 \\ 1 \end{bmatrix} \right\},$$

using the Gram-Schmidt process starting with the vector  $\mathbf{u}_1$ .

**6.5.2.-** Let  $W \subset \mathbb{R}^3$  be the subspace

$$\text{Span} \left\{ \mathbf{u}_1 = \begin{bmatrix} 0 \\ 2 \\ 0 \end{bmatrix}, \mathbf{u}_2 = \begin{bmatrix} 3 \\ 1 \\ 4 \end{bmatrix} \right\}.$$

- (a) Find an orthonormal basis for  $W$  using the Gram-Schmidt method starting with the vector  $\mathbf{u}_1$ .  
 (b) Decompose the vector  $\mathbf{x}$  below as

$$\mathbf{x} = \mathbf{x}_W + \mathbf{x}_W^\perp,$$

with  $\mathbf{x}_W \in W$  and  $\mathbf{x}_W^\perp \in W^\perp$ , where

$$\mathbf{x} = \begin{bmatrix} 5 \\ 1 \\ 0 \end{bmatrix}.$$

**6.5.3.-** Knowing that the column vectors in matrix  $A$  below form a linearly independent set, use the Gram-Schmidt method to find an orthonormal basis for  $R(A)$ , where

$$A = \begin{bmatrix} 1 & 2 & 5 \\ 0 & 2 & 0 \\ 1 & 0 & -1 \end{bmatrix}.$$

**6.5.4.-** Use the Gram-Schmidt method to find an orthonormal basis for  $R(A)$ , where

$$A = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 2 & -2 \\ 1 & 0 & 3 \end{bmatrix}.$$

**6.5.5.-** Consider the vector space  $\mathbb{P}_2([0, 1])$  with inner product

$$\langle \mathbf{p}, \mathbf{q} \rangle = \int_0^1 \mathbf{p}(x)\mathbf{q}(x) dx.$$

Use the Gram-Schmidt method on the ordered basis

$$(\mathbf{p}_0 = 1, \mathbf{p}_1 = x, \mathbf{p}_2 = x^2),$$

starting with vector  $\mathbf{p}_0$ , to obtain an orthogonal basis for  $\mathbb{P}_2([0, 1])$ .



## 6.6. THE ADJOINT OPERATOR

**6.6.1. The Riesz Representation Theorem.** The Riesz Representation Theorem is a statement concerning linear functionals on an inner product space. Recall that a *linear functional* on a vector space  $V$  over a scalar field  $\mathbb{F}$  is a linear function  $f : V \rightarrow \mathbb{F}$ , that is, for all  $\mathbf{x}, \mathbf{y} \in V$  and all  $a, b \in \mathbb{F}$  holds  $f(a\mathbf{x} + b\mathbf{y}) = af(\mathbf{x}) + bf(\mathbf{y}) \in \mathbb{F}$ . An example of a linear functional on  $\mathbb{R}^3$  is the function

$$\mathbb{R}^3 \ni \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \mapsto f(\mathbf{x}) = 3x_1 + 2x_2 + x_3 \in \mathbb{R}.$$

This function can be expressed in terms of the dot product in  $\mathbb{R}^3$  as follows

$$f(\mathbf{x}) = \mathbf{u} \cdot \mathbf{x}, \quad \mathbf{u} = \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix}.$$

The Riesz Representation Theorem says that what we did in this example can be done in the general case. In an inner product space  $(V, \langle \cdot, \cdot \rangle)$  every linear functional  $f$  can be expressed in terms of the inner product.

**Theorem 6.6.1.** *Consider a finite dimensional inner product space  $(V, \langle \cdot, \cdot \rangle)$  over the scalar field  $\mathbb{F}$ . For every linear functional  $f : V \rightarrow \mathbb{F}$  there exists a unique vector  $\mathbf{u}_f \in V$  such that holds*

$$f(\mathbf{v}) = \langle \mathbf{u}_f, \mathbf{v} \rangle \quad \forall \mathbf{v} \in V.$$

**Proof of Theorem 9.4.1:** Introduce the set

$$N = \{ \mathbf{v} \in V : f(\mathbf{v}) = 0 \} \subset V.$$

This set is the analogous to linear functionals of the null space of linear operators. Since  $f$  is a linear function the set  $N$  is a subspace of  $V$ . (Proof: Given two elements  $\mathbf{v}_1, \mathbf{v}_2 \in N$  and two scalars  $a, b \in \mathbb{F}$ , holds that  $f(a\mathbf{v}_1 + b\mathbf{v}_2) = af(\mathbf{v}_1) + bf(\mathbf{v}_2) = 0 + 0$ , so  $(a\mathbf{v}_1 + b\mathbf{v}_2) \in N$ .) Introduce the orthogonal complement of  $N$ , that is,

$$N^\perp = \{ \mathbf{w} \in V : \langle \mathbf{w}, \mathbf{v} \rangle = 0 \forall \mathbf{v} \in N \},$$

which is also a subspace of  $V$ . If  $N^\perp = \{\mathbf{0}\}$ , then  $N = (N^\perp)^\perp = (\{\mathbf{0}\})^\perp = V$ . Since the null space of  $f$  is the whole vector space, the functional  $f$  is identically zero, so only for the choice  $\mathbf{u}_f = \mathbf{0}$  holds  $f(\mathbf{v}) = \langle \mathbf{0}, \mathbf{v} \rangle$  for all  $\mathbf{v} \in V$ .

In the case that  $N^\perp \neq \{\mathbf{0}\}$  we now show that this space cannot be very big, in fact it has dimension one, as the following argument shows. Choose  $\tilde{\mathbf{u}} \in N^\perp$  such that  $f(\tilde{\mathbf{u}}) = 1$ . Then notice that for every  $\mathbf{w} \in N^\perp$  the vector  $\mathbf{w} - f(\mathbf{w})\tilde{\mathbf{u}}$  is trivially in  $N^\perp$  but it is also in  $N$ , since

$$f(\mathbf{w} - f(\mathbf{w})\tilde{\mathbf{u}}) = f(\mathbf{w}) - f(\mathbf{w})f(\tilde{\mathbf{u}}) = f(\mathbf{w}) - f(\mathbf{w}) = 0.$$

A vector both in  $N$  and  $N^\perp$  must vanish, so  $\mathbf{w} = f(\mathbf{w})\tilde{\mathbf{u}}$ . Then every vector in  $N^\perp$  is proportional to  $\tilde{\mathbf{u}}$ , so  $\dim N^\perp = 1$ . This information is used to split any vector  $\mathbf{v} \in V$  as follows  $\mathbf{v} = a\tilde{\mathbf{u}} + \mathbf{x}$  where  $\mathbf{x} \in N$  and  $a \in \mathbb{F}$ . It is clear that

$$f(\mathbf{v}) = f(a\tilde{\mathbf{u}} + \mathbf{x}) = af(\tilde{\mathbf{u}}) + f(\mathbf{x}) = af(\tilde{\mathbf{u}}) = a.$$

However, the function with values  $g(\mathbf{v}) = \left\langle \frac{\tilde{\mathbf{u}}}{\|\tilde{\mathbf{u}}\|^2}, \mathbf{v} \right\rangle$  has precisely the same values as  $f$ , since for all  $\mathbf{v} \in V$  holds

$$g(\mathbf{v}) = \left\langle \frac{\tilde{\mathbf{u}}}{\|\tilde{\mathbf{u}}\|^2}, \mathbf{v} \right\rangle = \left\langle \frac{\tilde{\mathbf{u}}}{\|\tilde{\mathbf{u}}\|^2}, (a\tilde{\mathbf{u}} + \mathbf{x}) \right\rangle = \frac{a}{\|\tilde{\mathbf{u}}\|^2} \langle \tilde{\mathbf{u}}, \tilde{\mathbf{u}} \rangle + \frac{1}{\|\tilde{\mathbf{u}}\|^2} \langle \tilde{\mathbf{u}}, \mathbf{x} \rangle = a.$$

Therefore, choosing  $\mathbf{u}_f = \tilde{\mathbf{u}}/\|\tilde{\mathbf{u}}\|^2$ , holds that

$$f(\mathbf{v}) = \langle \mathbf{u}_f, \mathbf{v} \rangle \quad \forall \mathbf{v} \in V.$$

Since  $\dim N^\perp = 1$ , the choice of  $\mathbf{u}_f$  is unique. This establishes the Theorem.  $\square$

**6.6.2. The adjoint operator.** Given a linear operator defined on an inner product space, a new linear operator can be defined through an equation involving the inner product.

**Proposition 6.6.2.** *If  $\mathbf{T} \in L(V)$  is a linear operator on a finite-dimensional inner product space  $(V, \langle \cdot, \cdot \rangle)$ , then there exists a unique linear operator  $\mathbf{T}^* \in L(V)$ , called the **adjoint** of  $\mathbf{T}$ , such that for all vectors  $\mathbf{u}, \mathbf{v} \in V$  holds*

$$\langle \mathbf{v}, \mathbf{T}^*(\mathbf{u}) \rangle = \langle \mathbf{T}(\mathbf{v}), \mathbf{u} \rangle.$$

**Proof of Proposition 9.4.2:** We first establish the following statement: For every vector  $\mathbf{u} \in V$  there exists a unique vector  $\mathbf{w} \in V$  such that

$$\langle \mathbf{T}(\mathbf{v}), \mathbf{u} \rangle = \langle \mathbf{v}, \mathbf{w} \rangle \quad \forall \mathbf{v} \in V. \quad (6.18)$$

The proof starts noticing that for a fixed  $\mathbf{u} \in V$  the scalar-valued function  $f_{\mathbf{u}} : V \rightarrow \mathbb{F}$  given by  $f_{\mathbf{u}}(\mathbf{v}) = \langle \mathbf{u}, \mathbf{T}(\mathbf{v}) \rangle$  is a linear functional. Therefore, the Riesz Representation Theorem 6.6.1 implies that there exists a unique vector  $\mathbf{w} \in V$  such that  $f_{\mathbf{u}}(\mathbf{v}) = \langle \mathbf{w}, \mathbf{v} \rangle$ . This establishes that for every vector  $\mathbf{u} \in V$  there exists a unique vector  $\mathbf{w} \in V$  such that Eq. (6.18) holds. Now that this statement is proven we can define a map, that we choose to denote as  $\mathbf{T}^* : V \rightarrow V$ , given by  $\mathbf{u} \mapsto \mathbf{T}^*(\mathbf{u}) = \mathbf{w}$ . We now show that this map  $\mathbf{T}^*$  is linear. Indeed, for all  $\mathbf{u}_1, \mathbf{u}_2 \in V$  and all  $a, b \in \mathbb{F}$  holds

$$\begin{aligned} \langle \mathbf{v}, \mathbf{T}^*(a\mathbf{u}_1 + b\mathbf{u}_2) \rangle &= \langle \mathbf{T}(\mathbf{v}), (a\mathbf{u}_1 + b\mathbf{u}_2) \rangle \quad \forall \mathbf{v} \in V, \\ &= a \langle \mathbf{T}(\mathbf{v}), \mathbf{u}_1 \rangle + b \langle \mathbf{T}(\mathbf{v}), \mathbf{u}_2 \rangle \\ &= a \langle \mathbf{v}, \mathbf{T}^*(\mathbf{u}_1) \rangle + b \langle \mathbf{v}, \mathbf{T}^*(\mathbf{u}_2) \rangle \\ &= \langle \mathbf{v}, [a \mathbf{T}^*(\mathbf{u}_1) + b \mathbf{T}^*(\mathbf{u}_2)] \rangle \quad \forall \mathbf{v} \in V, \end{aligned}$$

hence  $\mathbf{T}^*(a\mathbf{u}_1 + b\mathbf{u}_2) = a \mathbf{T}^*(\mathbf{u}_1) + b \mathbf{T}^*(\mathbf{u}_2)$ . This establishes the Proposition.  $\square$

The next result relates the adjoint of a linear operator with the concept of the adjoint of a square matrix introduced in Sect. 2.2. Recall that given a basis in the vector space, every linear operator has associated a unique square matrix. Let us use the notation  $[\mathbf{T}]$  and  $[\mathbf{T}^*]$  for the matrices on a given basis of the operators  $\mathbf{T}$  and  $\mathbf{T}^*$ , respectively.

**Proposition 6.6.3.** *Let  $(V, \langle \cdot, \cdot \rangle)$  be a finite-dimensional vector space, let  $\mathcal{V}$  be an orthonormal basis of  $V$ , and let  $[\mathbf{T}]$  be the matrix of the linear operator  $\mathbf{T} \in L(V)$  in the basis  $\mathcal{V}$ . Then, the matrix of the adjoint operator  $\mathbf{T}^*$  in the basis  $\mathcal{V}$  is given by  $[\mathbf{T}^*] = [\mathbf{T}]^*$ .*

Proposition 6.6.3 says that the matrix of the adjoint operator is the adjoint of the matrix of the operator, however this is true only in the case that the basis used to compute the respective matrices is orthonormal. If the basis is not orthonormal, the relation between the matrices  $[\mathbf{T}]$  and  $[\mathbf{T}^*]$  is more involved.

**Proof of Proposition 6.6.3:** Let  $\mathcal{V} = \{\mathbf{e}_1, \dots, \mathbf{e}_n\}$  be an orthonormal basis of  $V$ , that is,

$$\langle \mathbf{e}_i, \mathbf{e}_j \rangle = \begin{cases} 0 & \text{if } i \neq j, \\ 1 & \text{if } i = j. \end{cases}$$

The components of two arbitrary vectors  $\mathbf{u}, \mathbf{v} \in V$  in the basis  $\mathcal{V}$  is denoted as follows

$$\mathbf{u} = \sum_i u_i \mathbf{e}_i, \quad \mathbf{v} = \sum_i v_i \mathbf{e}_i.$$

The action of the operator  $\mathbf{T}$  can also be decomposed in the basis  $\mathcal{V}$  as follows

$$\mathbf{T}(\mathbf{e}_j) = \sum_i [\mathbf{T}]_{ij} \mathbf{e}_i, \quad [\mathbf{T}]_{ij} = [\mathbf{T}(\mathbf{e}_j)]_i.$$

We use the same notation for the adjoint operator, that is,

$$\mathbf{T}^*(\mathbf{e}_j) = \sum_i [\mathbf{T}^*]_{ij} \mathbf{e}_i, \quad [\mathbf{T}^*]_{ij} = [\mathbf{T}^*(\mathbf{e}_j)]_i.$$

The adjoint operator is defined such that the equation  $\langle \mathbf{v}, \mathbf{T}^*(\mathbf{u}) \rangle = \langle \mathbf{T}(\mathbf{v}), \mathbf{u} \rangle$  holds for all  $\mathbf{u}, \mathbf{v} \in V$ . This equation can be expressed in terms of components in the basis  $\mathcal{V}$  as follows

$$\sum_{ijk} \langle v_i \mathbf{e}_i, u_j [\mathbf{T}^*(\mathbf{e}_j)]_k \mathbf{e}_k \rangle = \sum_{ijk} \langle v_i [\mathbf{T}(\mathbf{e}_i)]_k \mathbf{e}_k, u_j \mathbf{e}_j \rangle,$$

that is,

$$\sum_{ijk} \bar{v}_i u_j [\mathbf{T}^*]_{kj} \langle \mathbf{e}_i, \mathbf{e}_k \rangle = \sum_{ijk} \bar{v}_i [\mathbf{T}]_{ki} u_j \langle \mathbf{e}_k, \mathbf{e}_j \rangle.$$

Since the basis  $\mathcal{V}$  is orthonormal we obtain the equation

$$\sum_{ij} \bar{v}_i u_j [\mathbf{T}^*]_{ij} = \sum_{ijk} \bar{v}_i [\mathbf{T}]_{ji} u_j,$$

which holds for all vectors  $\mathbf{u}, \mathbf{v} \in V$ , so we conclude

$$[\mathbf{T}^*]_{ij} = \overline{[\mathbf{T}]_{ji}} \quad \Leftrightarrow \quad [\mathbf{T}^*] = \overline{[\mathbf{T}]^T} \quad \Leftrightarrow \quad [\mathbf{T}^*] = [\mathbf{T}]^*.$$

This establishes the Proposition.  $\square$

**EXAMPLE 6.6.1:** Consider the inner product space  $(\mathbb{C}^3, \cdot)$ . Find the adjoint of the linear operator  $\mathbf{T}$  with matrix in the standard basis of  $\mathbb{C}^3$  given by

$$[\mathbf{T}(\mathbf{x})] = \begin{bmatrix} x_1 + 2ix_2 + ix_3 \\ ix_1 - x_3 \\ x_1 - x_2 + ix_3 \end{bmatrix}, \quad [\mathbf{x}] = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}.$$

**SOLUTION:** The matrix of this operator in the standard basis of  $\mathbb{C}^3$  is given by

$$[\mathbf{T}] = \begin{bmatrix} 1 & 2i & i \\ i & 0 & -1 \\ 1 & -1 & i \end{bmatrix}.$$

Since the standard basis is an orthonormal basis with respect to the dot product, Proposition 6.6.3 implies that

$$[\mathbf{T}^*] = [\mathbf{T}]^* = \begin{bmatrix} 1 & 2i & i \\ i & 0 & -1 \\ 1 & -1 & i \end{bmatrix}^* = \begin{bmatrix} 1 & -i & 1 \\ -2i & 0 & -1 \\ -i & -1 & -i \end{bmatrix} \Rightarrow [\mathbf{T}^*(\mathbf{x})] = \begin{bmatrix} x_1 - ix_2 + x_3 \\ -2ix_1 - x_3 \\ -ix_1 - x_2 - ix_3 \end{bmatrix}.$$

$\triangleleft$

**6.6.3. Normal operators.** Recall now that the commutator of two linear operators  $\mathbf{T}, \mathbf{S} \in L(V)$  is the linear operator  $[\mathbf{T}, \mathbf{S}] \in L(V)$  given by

$$[\mathbf{T}, \mathbf{S}](\mathbf{u}) = \mathbf{T}(\mathbf{S}(\mathbf{u})) - \mathbf{S}(\mathbf{T}(\mathbf{u})) \quad \forall \mathbf{u} \in V.$$

Two operators  $\mathbf{T}, \mathbf{S} \in L(V)$  are said to commute iff their commutator vanishes, that is,  $[\mathbf{T}, \mathbf{S}] = \mathbf{0}$ . Examples of operators that commute are two rotations on the plane. Examples of operators that do not commute are two arbitrary rotations in space.

**Definition 6.6.4.** A linear operator  $\mathbf{T}$  defined on a finite-dimensional inner product space  $(V, \langle \cdot, \cdot \rangle)$  is called a **normal operator** iff holds  $[\mathbf{T}, \mathbf{T}^*] = \mathbf{0}$ , that is, the operator commutes with its adjoint.

An interesting characterization of normal operators is the following: A linear operator  $\mathbf{T}$  on an inner product space is normal iff  $\|\mathbf{T}(\mathbf{u})\| = \|\mathbf{T}^*(\mathbf{u})\|$  holds for all  $\mathbf{u} \in V$ . Normal operators are particularly important because for these operators hold the Spectral Theorem, which we study in Chapter 9.

Two particular cases of normal operators are often used in physics. A linear operator  $\mathbf{T}$  on an inner product space is called a **unitary operator** iff  $\mathbf{T}^* = \mathbf{T}^{-1}$ , that is, the adjoint is the inverse operator. Unitary operators are normal operators, since

$$\mathbf{T}^* = \mathbf{T}^{-1} \quad \Rightarrow \quad \begin{cases} \mathbf{T}\mathbf{T}^* = \mathbf{I}, \\ \mathbf{T}^*\mathbf{T} = \mathbf{I}, \end{cases} \quad \Rightarrow \quad [\mathbf{T}, \mathbf{T}^*] = \mathbf{0}.$$

Unitary operators preserve the length of a vector, since

$$\|\mathbf{v}\|^2 = \langle \mathbf{v}, \mathbf{v} \rangle = \langle \mathbf{v}, \mathbf{T}^{-1}(\mathbf{T}(\mathbf{v})) \rangle = \langle \mathbf{v}, \mathbf{T}^*(\mathbf{T}(\mathbf{v})) \rangle = \langle \mathbf{T}(\mathbf{v}), \mathbf{T}(\mathbf{v}) \rangle = \|\mathbf{T}(\mathbf{v})\|^2.$$

Unitary operators defined on a complex inner product space are particularly important in quantum mechanics. The particular case of unitary operators on a real inner product space are called **orthogonal operators**. So, orthogonal operators do not change the length of a vector. Examples of orthogonal operators are rotations in  $\mathbb{R}^3$ .

A linear operator  $\mathbf{T}$  on an inner product space is called an **Hermitian operator** iff  $\mathbf{T}^* = \mathbf{T}$ , that is, the adjoint is the original operator. This definition agrees with the definition of Hermitian matrices given in Chapter 2.

**EXAMPLE 6.6.2:** Consider the inner product space  $(\mathbb{C}^3, \cdot)$  and the linear operator  $\mathbf{T}$  with matrix in the standard basis of  $\mathbb{C}^3$  given by

$$[\mathbf{T}(\mathbf{x})] = \begin{bmatrix} x_1 - ix_2 + x_3 \\ ix_1 - x_3 \\ x_1 - x_2 + x_3 \end{bmatrix}, \quad [\mathbf{x}] = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}.$$

Show that  $\mathbf{T}$  is Hermitian.

**SOLUTION:** We need to compute the adjoint of  $\mathbf{T}$ . The matrix of this operator in the standard basis of  $\mathbb{C}^3$  is given by

$$[\mathbf{T}] = \begin{bmatrix} 1 & -i & 1 \\ i & 0 & -1 \\ 1 & -1 & 1 \end{bmatrix}.$$

Since the standard basis is an orthonormal basis with respect to the dot product, Proposition 6.6.3 implies that

$$[\mathbf{T}^*] = [\mathbf{T}]^* = \begin{bmatrix} 1 & -i & 1 \\ i & 0 & -1 \\ 1 & -1 & 1 \end{bmatrix}^* = \begin{bmatrix} 1 & i & 1 \\ -i & 0 & -1 \\ 1 & -1 & 1 \end{bmatrix} = [\mathbf{T}].$$

Therefore,  $\mathbf{T}^* = \mathbf{T}$ . ◀

#### 6.6.4. Bilinear forms.

**Definition 6.6.5.** A **bilinear form** on a vector space  $V$  over  $\mathbb{F}$  is a function  $a : V \times V \rightarrow \mathbb{F}$  linear on both arguments, that is, for all  $\mathbf{u}, \mathbf{v}_1, \mathbf{v}_2 \in V$  and all  $b_1, b_2 \in \mathbb{F}$  holds

$$\begin{aligned} a(\mathbf{u}, (b_1\mathbf{v}_1 + b_2\mathbf{v}_2)) &= b_1a(\mathbf{u}, \mathbf{v}_1) + b_2a(\mathbf{u}, \mathbf{v}_2), \\ a((b_1\mathbf{v}_1 + b_2\mathbf{v}_2), \mathbf{u}) &= b_1a(\mathbf{v}_1, \mathbf{u}) + b_2a(\mathbf{v}_2, \mathbf{u}). \end{aligned}$$

The bilinear form  $a : V \times V \rightarrow \mathbb{F}$  is called **symmetric** iff for all  $\mathbf{u}, \mathbf{v} \in V$  holds

$$a(\mathbf{u}, \mathbf{v}) = a(\mathbf{v}, \mathbf{u}).$$

An example of a symmetric bilinear form is any the inner product on a real vector space. Indeed, given a real inner product space  $(V, \langle \cdot, \cdot \rangle)$ , the function  $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$  is a symmetric bilinear form, since it is symmetric and

$$\langle \mathbf{u}, (b_1 \mathbf{v}_1 + b_2 \mathbf{v}_2) \rangle = b_1 \langle \mathbf{u}, \mathbf{v}_1 \rangle + b_2 \langle \mathbf{u}, \mathbf{v}_2 \rangle.$$

We will shortly see that another example of a bilinear form appears on the weak formulation of the boundary value problem in Eq. (7.30).

On the other hand, an inner product on a complex vector space is not a bilinear form, since it is conjugate linear on the first argument instead of linear. Such functions are called sesquilinear forms. That is, a **sesquilinear form** on a complex vector space  $V$  is a function  $a : V \times V \rightarrow \mathbb{C}$  that is conjugate linear on the first argument and linear on the second argument. Sesquilinear forms are important in the case of studying differential equations involving complex functions.

**Definition 6.6.6.** Consider a bilinear form  $a : V \times V \rightarrow \mathbb{F}$  on an inner product space  $(V, \langle \cdot, \cdot \rangle)$  over  $\mathbb{F}$ . The bilinear form  $a$  is called **positive** iff there exists a real number  $k > 0$  such that for all  $\mathbf{u} \in V$  holds

$$a(\mathbf{u}, \mathbf{u}) \geq k \|\mathbf{u}\|^2.$$

The bilinear form  $a$  is called **bounded** iff there exists a real number  $K > 0$  such that for all  $\mathbf{u}, \mathbf{v} \in V$  holds

$$a(\mathbf{u}, \mathbf{v}) \leq K \|\mathbf{u}\| \|\mathbf{v}\|.$$

An example of a positive bilinear form is any inner product on a real vector space. The Schwarz inequality implies that such inner product is also a bounded bilinear form. In fact, an inner product on a real vector space is a symmetric, positive, bounded bilinear form. We will shortly see that the bilinear form that appears on a weak formulation of the boundary value problem in Eq. (7.30) is symmetric, positive and bounded. We will see that these properties imply the existence and uniqueness of solutions to the weak formulation of the boundary value problem.

**6.6.5. Exercises.****6.6.1.-** .**6.6.2.-** .

## CHAPTER 7. APPROXIMATION METHODS

## 7.1. BEST APPROXIMATION

The first half of this Section is dedicated to show that the Fourier series approximation of a function is a particular case of the orthogonal decomposition of a vector onto a subspace in an inner product space, studied in Section 6.4. Once we realize that it is not difficult to see why such Fourier approximations are useful. We show that the orthogonal projection  $\mathbf{x}_\perp$  of a vector  $\mathbf{x}$  onto a subspace  $U$  is the vector in the subspace  $U$  closest to  $\mathbf{x}$ . This is the origin of the name “best approximation” for  $\mathbf{x}_\perp$ . What is intuitively clear in  $(\mathbb{R}^3, \cdot)$  is true in every inner product space, hence it is true for the Fourier series approximation of a function. In the second half of this Section we show a deep relation between the Null space of a matrix and the Range space of the adjoint matrix. The former is the orthogonal complement of the latter. A consequence of this relation is a simple proof to the property that a matrix and its adjoint matrix have the same rank. Another consequence is given in the next Section, where we obtain a simple equation, called the normal equation, to find a least squares solution of an inconsistent linear system.

**7.1.1. Fourier expansions.** We have seen that orthonormal bases have a practical advantage over arbitrary basis. The components  $[\mathbf{x}]_u$  of a vector  $\mathbf{x}$  in an orthonormal basis  $\mathcal{U}_n = \{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  of the inner product space  $(V, \langle \cdot, \cdot \rangle)$  are given by the simple expression

$$[\mathbf{x}]_u = \begin{bmatrix} \langle \mathbf{u}_1, \mathbf{x} \rangle \\ \vdots \\ \langle \mathbf{u}_n, \mathbf{x} \rangle \end{bmatrix} \Leftrightarrow \mathbf{x} = \langle \mathbf{u}_1, \mathbf{x} \rangle \mathbf{u}_1 + \dots + \langle \mathbf{u}_n, \mathbf{x} \rangle \mathbf{u}_n.$$

In the case that an orthonormal set  $\mathcal{U}_p$  is not a basis of  $V$ , that is,  $p < \dim V$ , one can always introduce the orthogonal projection of a vector  $\mathbf{x} \in V$  onto the subspace  $U_p = \text{Span}(\mathcal{U}_p)$ ,

$$\mathbf{x}_\perp = \langle \mathbf{u}_1, \mathbf{x} \rangle \mathbf{u}_1 + \dots + \langle \mathbf{u}_p, \mathbf{x} \rangle \mathbf{u}_p.$$

We have seen that in this case,  $\mathbf{x}_\perp \neq \mathbf{x}$ . In fact we called the difference vector  $\mathbf{x}_\perp = \mathbf{x} - \mathbf{x}_\perp$ , since this vector satisfies that  $\mathbf{x}_\perp \perp \mathbf{x}_\perp$ . We now give the projection vector  $\mathbf{x}_\perp$  a new name.

**Definition 7.1.1.** The **Fourier expansion** of a vector  $\mathbf{x} \in V$  with respect to an orthonormal set  $\mathcal{U}_p = \{\mathbf{u}_1, \dots, \mathbf{u}_p\} \subset (V, \langle \cdot, \cdot \rangle)$  is the unique vector  $\mathbf{x}_\perp \in \text{Span}(\mathcal{U}_p)$  given by

$$\mathbf{x}_\perp = \langle \mathbf{u}_1, \mathbf{x} \rangle \mathbf{u}_1 + \dots + \langle \mathbf{u}_p, \mathbf{x} \rangle \mathbf{u}_p. \quad (7.1)$$

The scalars  $\langle \mathbf{u}_i, \mathbf{x} \rangle$  are the **Fourier coefficient** of the vector  $\mathbf{x}$  with respect to the set  $\mathcal{U}_p$ .

A reason to the name “Fourier expansion” to the orthogonal projection of a vector onto a subspace is given in the following example.

**EXAMPLE 7.1.1:** Given the vector space of continuous functions  $V = C([-l, l], \mathbb{R})$  with inner product given by  $\langle \mathbf{f}, \mathbf{g} \rangle = \int_{-l}^l \mathbf{f}(x)\mathbf{g}(x) dx$ , find the Fourier expansion of an arbitrary function  $\mathbf{f} \in V$  with respect to the orthonormal set

$$\mathcal{U}_N = \left\{ \mathbf{u}_0 = \frac{1}{\sqrt{2l}}, \mathbf{u}_n = \frac{1}{\sqrt{l}} \cos\left(\frac{n\pi x}{l}\right), \mathbf{v}_n = \frac{1}{\sqrt{l}} \sin\left(\frac{n\pi x}{l}\right) \right\}_{n=1}^N.$$

**SOLUTION:** Using Eq. (7.1) on a function  $\mathbf{f} \in V$  we get

$$\mathbf{f}_\perp = \langle \mathbf{u}_0, \mathbf{f} \rangle \mathbf{u}_0 + \sum_{n=1}^N \left[ \langle \mathbf{u}_n, \mathbf{f} \rangle \mathbf{u}_n + \langle \mathbf{v}_n, \mathbf{f} \rangle \mathbf{v}_n \right].$$

Introduce the vectors in  $\mathcal{U}_N$  explicitly in the expression above,

$$\mathbf{f}_{||} = \frac{1}{\sqrt{2\ell}} \langle \mathbf{u}_0, \mathbf{f} \rangle + \sum_{n=1}^N \left[ \frac{1}{\sqrt{\ell}} \langle \mathbf{u}_n, \mathbf{f} \rangle \cos\left(\frac{n\pi x}{\ell}\right) + \frac{1}{\sqrt{\ell}} \langle \mathbf{v}_n, \mathbf{f} \rangle \sin\left(\frac{n\pi x}{\ell}\right) \right].$$

Denoting by

$$a_0 = \frac{1}{\sqrt{2\ell}} \langle \mathbf{u}_0, \mathbf{f} \rangle \quad a_n = \frac{1}{\sqrt{\ell}} \langle \mathbf{u}_n, \mathbf{f} \rangle \quad b_n = \frac{1}{\sqrt{\ell}} \langle \mathbf{v}_n, \mathbf{f} \rangle,$$

then we get that

$$\mathbf{f}_{||}(x) = a_0 + \sum_{n=1}^N \left[ a_n \cos\left(\frac{n\pi x}{\ell}\right) + b_n \sin\left(\frac{n\pi x}{\ell}\right) \right],$$

where the coefficients  $a_0, a_n$  and  $b_n$ , for  $n = 1, \dots, N$  are the usual Fourier coefficients,

$$a_0 = \frac{1}{2\ell} \int_{-\ell}^{\ell} \mathbf{f}(x) dx, \quad a_n = \frac{1}{\ell} \int_{-\ell}^{\ell} \mathbf{f}(x) \cos\left(\frac{n\pi x}{\ell}\right) dx, \quad b_n = \frac{1}{\ell} \int_{-\ell}^{\ell} \mathbf{f}(x) \sin\left(\frac{n\pi x}{\ell}\right) dx.$$

◁

The example above is a good reason to name Fourier expansion the orthogonal projection of a vector onto a subspace. We already know that the Fourier expansion  $\mathbf{x}_{||}$  of a vector  $\mathbf{x}$  with respect to an orthonormal set  $\mathcal{U}_p$  has a particular property, that is,  $(\mathbf{x} - \mathbf{x}_{||}) \perp \mathbf{x}_{||} \in U_p^\perp$ . This property means that  $\mathbf{x}_{||}$  is the best approximation of the vector  $\mathbf{x}$  from within the subspace  $\text{Span}(\mathcal{U}_p)$ . See Fig. 45 for the case  $V = \mathbb{R}^3$ , with  $\langle \cdot, \cdot \rangle = \cdot$ , and  $\text{Span}(\mathcal{U}_2)$  two-dimensional. We highlight this property in the following result.

**Theorem 7.1.2 (Best approximation).** *The Fourier expansion  $\mathbf{x}_{||}$  of a vector  $\mathbf{x}$  with respect to an orthonormal set  $\mathcal{U}_p$  in an inner product space, is the unique vector in the subspace  $\text{Span}(\mathcal{U}_p)$  that is closest to  $\mathbf{x}$ , in the sense that*

$$\|\mathbf{x} - \mathbf{x}_{||}\| < \|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{y} \in \text{Span}(\mathcal{U}_p) - \{\mathbf{x}_{||}\}.$$

**Proof of Theorem 7.1.2:** Recall that  $\mathbf{x}_\perp = \mathbf{x} - \mathbf{x}_{||}$  is orthogonal to  $\text{Span}(U)$ , that is  $(\mathbf{x} - \mathbf{x}_{||}) \perp (\mathbf{x}_{||} - \mathbf{y})$  for all  $\mathbf{y} \in \text{Span}(U)$ . Hence,

$$\|\mathbf{x} - \mathbf{y}\|^2 = \|(\mathbf{x} - \mathbf{x}_{||}) + (\mathbf{x}_{||} - \mathbf{y})\|^2 = \|\mathbf{x} - \mathbf{x}_{||}\|^2 + \|\mathbf{x}_{||} - \mathbf{y}\|^2, \tag{7.2}$$

where the last equality comes from Pythagoras Theorem. Eq. (7.2) says that  $\|\mathbf{x} - \mathbf{y}\|$  is the smallest iff  $\mathbf{y} = \mathbf{x}_{||}$  and the smallest value is  $\|\mathbf{x} - \mathbf{x}_{||}\|$ . This establishes the Theorem.  $\square$

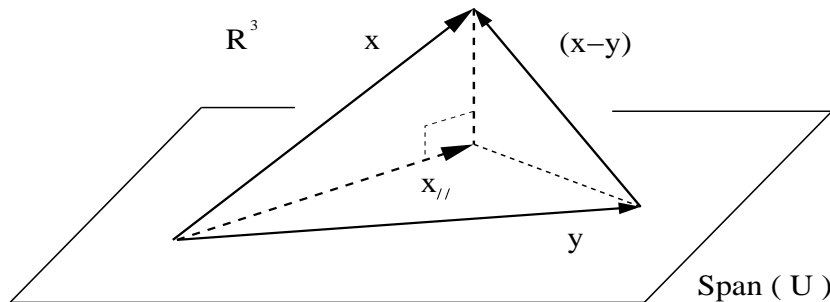


FIGURE 45. The Fourier expansion  $\mathbf{x}_{||}$  of the vector  $\mathbf{x} \in \mathbb{R}^3$  is the best approximation of  $\mathbf{x}$  from within  $\text{Span}(U)$ .



**EXAMPLE 7.1.2:** Given the vector space of continuous functions  $V = C([-1, 1], \mathbb{R})$  with the inner product given by  $\langle \mathbf{f}, \mathbf{g} \rangle = \int_{-1}^1 \mathbf{f}(x)\mathbf{g}(x) dx$ , find the Fourier expansion of the function  $\mathbf{f}(x) = x$  with respect to the orthonormal set

$$\mathcal{U}_N = \left\{ \mathbf{u}_0 = \frac{1}{\sqrt{2\ell}}, \mathbf{u}_n = \frac{1}{\sqrt{\ell}} \cos\left(\frac{n\pi x}{\ell}\right), \mathbf{v}_n = \frac{1}{\sqrt{\ell}} \sin\left(\frac{n\pi x}{\ell}\right) \right\}_{n=1}^N.$$

**SOLUTION:** We use the formulas in Example 7.1.1 above to the function  $\mathbf{f}(x) = x$  on the interval  $[-1, 1]$ , since  $\ell = 1$ . The coefficient  $a_0$  above is given by

$$a_0 = \frac{1}{2} \int_{-1}^1 x dx = \frac{1}{4} (x^2|_{-1}^1) \Rightarrow a_0 = 0.$$

The coefficients  $a_n, b_n$ , for  $n = 1, \dots, N$  computed with one integration by parts,

$$\begin{aligned} \int x \cos(n\pi x) dx &= \frac{x}{n\pi} \sin(n\pi x) + \frac{1}{n^2\pi^2} \cos(n\pi x), \\ \int x \sin(n\pi x) dx &= -\frac{x}{n\pi} \cos(n\pi x) + \frac{1}{n^2\pi^2} \sin(n\pi x). \end{aligned}$$

The coefficients  $a_n$  vanish, since

$$a_n = \int_{-1}^1 x \cos(n\pi x) dx = \left[ \frac{x}{n\pi} \sin(n\pi x) \right]_{-1}^1 + \frac{1}{n^2\pi^2} \cos(n\pi x) \Big|_{-1}^1 \Rightarrow a_n = 0.$$

The coefficients  $b_n$  are given by

$$b_n = \int_{-1}^1 x \sin(n\pi x) dx = -\left[ \frac{x}{n\pi} \cos(n\pi x) \right]_{-1}^1 + \frac{1}{n^2\pi^2} \sin(n\pi x) \Big|_{-1}^1 \Rightarrow b_n = \frac{2(-1)^{(n+1)}}{n\pi}.$$

Therefore, the Fourier expansion of  $\mathbf{f}(x) = x$  with respect to  $\mathcal{U}_N$  is given by

$$\mathbf{f}_{||}(x) = \frac{2}{\pi} \sum_{n=1}^N \frac{(-1)^{(n+1)}}{n} \sin(n\pi x).$$

**REMARK:** First, a simpler proof that the coefficients  $a_0$  and  $a_n$  vanish is to realized that we are integrating an odd function on the interval  $[-1, 1]$ . The odd function is the product of an odd function times an even function. Second, Theorem 7.1.2 tells us that the function  $\mathbf{f}_{||}$  above is only combination of the sine and cosine functions in  $\mathcal{U}_N$  that approximates best the function  $\mathbf{f}(x) = x$  on the interval  $[-1, 1]$ .  $\triangleleft$

**7.1.2. Null and range spaces of a matrix.** The null and range spaces associated with a matrix  $\mathbf{A} \in \mathbb{F}^{m,n}$  and its adjoint matrix  $\mathbf{A}^*$  are deeply related.

**Theorem 7.1.3.** For every matrix  $\mathbf{A} \in \mathbb{F}^{m,n}$  holds  $N(\mathbf{A}) = R(\mathbf{A}^*)^\perp$  and  $N(\mathbf{A}^*) = R(\mathbf{A})^\perp$ .

Since for every subspace  $W$  on a finite dimensional inner product space holds  $(W^\perp)^\perp = W$ , we also have the relations

$$N(\mathbf{A})^\perp = R(\mathbf{A}^*), \quad N(\mathbf{A}^*)^\perp = R(\mathbf{A}).$$

In the case of real-valued matrices, the Theorem above says that

$$N(\mathbf{A}) = R(\mathbf{A}^T)^\perp \quad \text{and} \quad N(\mathbf{A}^T) = R(\mathbf{A})^\perp.$$

Before we state the proof of this Theorem let us review the following notation: Given an  $m \times n$  matrix  $A \in \mathbb{F}^{m,n}$ , we write it either in terms of column vectors  $A_{:j} \in \mathbb{F}^m$  for  $j = 1, \dots, n$ , or in terms of row vectors  $A_{i:} \in \mathbb{F}^n$  for  $i = 1, \dots, m$ , as follows,

$$A = [A_{:1} \quad \dots \quad A_{:n}], \quad A = \begin{bmatrix} A_{1:} \\ \vdots \\ A_{m:} \end{bmatrix}.$$

Since the same type of definition holds for the  $n \times m$  matrix  $A^*$ , that is,

$$A^* = [(A^*)_{:1} \quad \dots \quad (A^*)_{:m}], \quad A^* = \begin{bmatrix} (A^*)_{1:} \\ \vdots \\ (A^*)_{n:} \end{bmatrix},$$

then we have the relations

$$(A_{:j})^* = (A^*)_{j:}, \quad (A_{i:})^* = (A^*)_{:i}.$$

For example, consider the  $2 \times 3$  matrix

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \Rightarrow A_{:1} = \begin{bmatrix} 1 \\ 4 \end{bmatrix}, A_{:2} = \begin{bmatrix} 2 \\ 5 \end{bmatrix}, A_{:3} = \begin{bmatrix} 3 \\ 6 \end{bmatrix}, \quad \begin{matrix} A_{1:} = [1 & 2 & 3], \\ A_{2:} = [4 & 5 & 6]. \end{matrix}$$

The transpose is a  $3 \times 2$  matrix that can be written as follows

$$A^T = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix} \Rightarrow (A^T)_{:1} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, (A^T)_{:2} = \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix}, \quad \begin{matrix} (A^T)_{1:} = [1 & 4], \\ (A^T)_{2:} = [2 & 5], \\ (A^T)_{3:} = [3 & 6]. \end{matrix}$$

So, for example we have the relation

$$(A_{:3})^T = [3 \quad 6] = (A^T)_{3:}, \quad (A_{2:})^T = \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix} = (A^T)_{:2}.$$

**Proof of Theorem 7.1.3:** We first show that the  $N(A) = R(A^*)^\perp$ . A vector  $x \in \mathbb{F}^n$  belongs to  $N(A)$  iff holds

$$Ax = 0 \Leftrightarrow \begin{bmatrix} A_{1:} \\ \vdots \\ A_{m:} \end{bmatrix} x = 0 \Leftrightarrow \begin{bmatrix} [(A^*)_{:1}]^* \\ \vdots \\ [(A^*)_{:m}]^* \end{bmatrix} x = 0 \Leftrightarrow \begin{cases} (A^*)_{:1} \cdot x = 0, \\ \vdots \\ (A^*)_{:m} \cdot x = 0. \end{cases}$$

So,  $x \in N(A)$  iff  $x$  is orthogonal to every column vector in  $A^*$ , that is,  $x \in R(A^*)^\perp$ .

The equation  $N(A^*) = R(A)^\perp$  comes from  $N(B) = R(B^*)^\perp$  taking  $B = A^*$ . Nevertheless, we repeat the proof above, just to understand the previous argument. A vector  $y \in \mathbb{F}^m$  belongs to  $N(A^*)$  iff

$$A^*y = 0 \Leftrightarrow \begin{bmatrix} (A^*)_{1:} \\ \vdots \\ (A^*)_{n:} \end{bmatrix} y = 0 \Leftrightarrow \begin{bmatrix} (A_{:1})^* \\ \vdots \\ (A_{:n})^* \end{bmatrix} y = 0 \Leftrightarrow \begin{cases} A_{:1} \cdot y = 0, \\ \vdots \\ A_{:n} \cdot y = 0. \end{cases}$$

So,  $y \in N(A^*)$  iff  $y$  is orthogonal to every column vector in  $A$ , that is,  $y \in R(A)^\perp$ . This establishes the Theorem.  $\square$

**EXAMPLE 7.1.3:** Verify Theorem 7.1.3 for the matrix  $A = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{bmatrix}$ .

**SOLUTION:** We first find the  $N(A)$ , that is, all  $x \in \mathbb{R}^3$  solutions of  $Ay = 0$ . Gauss operations on matrix  $A$  imply

$$\begin{bmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & 2 \end{bmatrix} \Rightarrow \begin{cases} x_1 = x_3, \\ x_2 = -2x_3, \\ x_3 \text{ free,} \end{cases} \Rightarrow N(A) = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}\right\}\right).$$

It is simple to find  $R(A^T)$ , since

$$R(A^T) = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix}\right\}\right).$$

Theorem 7.1.3 is verified, since

$$\begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} = 1 - 4 + 3 = 0, \quad \begin{bmatrix} 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} = 3 - 4 + 1 = 0 \Rightarrow N(A) = R(A^T)^\perp.$$

Let us verify the same Theorem for  $A^T$ . We first find  $N(A^T)$ , that is, all  $y \in \mathbb{R}^2$  solutions of  $A^T y = 0$ . Gauss operations on matrix  $A^T$  imply

$$\begin{bmatrix} 1 & 3 \\ 2 & 2 \\ 3 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \Rightarrow y = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow N(A^T) = \{0\}.$$

The space  $R(A)$  is given by

$$R(A) = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 3 \end{bmatrix}, \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \begin{bmatrix} 3 \\ 1 \end{bmatrix}\right\}\right) = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 3 \end{bmatrix}, \begin{bmatrix} 2 \\ 2 \end{bmatrix}\right\}\right) = \mathbb{R}^2.$$

Since  $(\mathbb{R}^2)^\perp = \{0\}$ , Theorem 7.1.3 is verified.  $\triangleleft$

Theorem 7.1.3 provides a simple proof for a result we used in Chapter 2.

**Theorem 7.1.4.** For every matrix  $A \in \mathbb{F}^{m,n}$  holds that  $\text{rank}(A) = \text{rank}(A^*)$ .

**Proof of Theorem 7.1.4:** Recall the Nullity-Rank result in Corollary 5.1.8, which says that for all matrix  $A \in \mathbb{F}^{m,n}$  holds  $\dim N(A) + \dim R(A) = n$ . Equivalently,

$$\dim R(A) = n - \dim N(A) = n - \dim R(A^*)^\perp,$$

since  $N(A) = R(A^*)^\perp$ . From the orthogonal decomposition  $\mathbb{F}^n = R(A^*) \oplus R(A^*)^\perp$  we know that  $\dim R(A^*) = n - \dim R(A^*)^\perp$ . We then conclude that

$$\dim R(A) = \dim R(A^*).$$

This establishes the Theorem.  $\square$

## 7.1.3. Exercises.

7.1.1.- Consider the inner product space  $(\mathbb{R}^3, \cdot)$  and the orthonormal set  $\mathcal{U}$ ,

$$\left\{ \mathbf{u}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix}, \mathbf{u}_2 = \frac{1}{\sqrt{3}} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right\}.$$

Find the best approximation of  $\mathbf{x}$  below in the subspace  $\text{Span}(\mathcal{U})$ , where

$$\mathbf{x} = \begin{bmatrix} 1 \\ 0 \\ -2 \end{bmatrix}.$$

7.1.2.- Consider the inner product space  $(\mathbb{R}^{2,2}, \langle \cdot, \cdot \rangle_F)$  and the orthonormal set  $\mathcal{U} = \{\mathbf{E}_1, \mathbf{E}_2\}$ , where

$$\mathbf{E}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \mathbf{E}_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}.$$

Find the best approximation of matrix  $\mathbf{A}$  below in the subspace  $\text{Span}(\mathcal{U})$ , where

$$\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}.$$

7.1.3.- Consider the inner product space  $\mathbb{P}_2([0, 1])$ , with  $\langle \mathbf{p}, \mathbf{q} \rangle = \int_0^1 \mathbf{p}(x)\mathbf{q}(x) dx$ , and the subspace  $U = \text{Span}(\mathcal{U})$ , where  $\mathcal{U} = \{\mathbf{q}_0 = 1, \mathbf{q}_1 = (x - \frac{1}{2})\}$ .

- Show that  $\mathcal{U}$  is an orthogonal set.
- Find  $\mathbf{r}_n$ , the best approximation with respect to  $U$  of the polynomial  $\mathbf{r}(x) = 2x + 3x^2$ .
- Verify whether  $(\mathbf{r} - \mathbf{r}_n) \in U^\perp$  or not.

7.1.4.- Consider the space  $C^\infty([-\ell, \ell], \mathbb{R})$  with inner product

$$\langle \mathbf{f}, \mathbf{g} \rangle = \int_{-\ell}^{\ell} \mathbf{f}(x)\mathbf{g}(x) dx,$$

and the orthonormal set  $\mathcal{U}$  given by

$$\begin{aligned} \mathbf{u}_0 &= \frac{1}{\sqrt{2\ell}} \\ \mathbf{u}_1 &= \frac{1}{\sqrt{\ell}} \cos\left(\frac{\pi x}{\ell}\right) \\ \mathbf{v}_1 &= \frac{1}{\sqrt{\ell}} \sin\left(\frac{\pi x}{\ell}\right). \end{aligned}$$

Find the best approximation of

$$\mathbf{f}(x) = \begin{cases} x & 0 \leq x \leq \ell, \\ -x & -\ell \leq x < 0. \end{cases}$$

in the space  $\text{Span}(\mathcal{U})$

7.1.5.- For the matrix  $\mathbf{A} \in \mathbb{R}^{3,3}$  below, verify that  $N(\mathbf{A}) = R(\mathbf{A}^T)^\perp$  and that  $N(\mathbf{A}^T) = R(\mathbf{A})^\perp$ , where

$$\mathbf{A} = \begin{bmatrix} 2 & 1 & 1 \\ -1 & -1 & 0 \\ -2 & -1 & -1 \end{bmatrix}.$$

## 7.2. LEAST SQUARES

**7.2.1. The normal equation.** We describe the least squares method to find approximate solutions to inconsistent linear systems. The method is often used to find the best parameters that fit experimental data. The parameters are the unknowns of the linear system, and the experimental data determines the matrix of coefficients and the source vector of the system. Such a linear system usually contains more equations than unknowns, and it is inconsistent, since there are no parameters that fit all the data exactly. We start introducing the notion of least squares solution of a possibly inconsistent linear system.

**Definition 7.2.1.** Given a matrix  $A \in \mathbb{F}^{m,n}$  and a vector  $\mathbf{b} \in (\mathbb{F}^m, \cdot)$ , the vector  $\hat{\mathbf{x}} \in \mathbb{F}^n$  is called a **least squares solution** of the linear system  $A\mathbf{x} = \mathbf{b}$  iff holds

$$\|A\hat{\mathbf{x}} - \mathbf{b}\| \leq \|y - \mathbf{b}\| \quad \forall y \in R(A).$$

The problem we study is to find the least squares solution to an  $m \times n$  linear system  $A\mathbf{x} = \mathbf{b}$ . In the case that  $\mathbf{b} \in R(A)$  the linear system  $A\mathbf{x} = \mathbf{b}$  is consistent and the least squares solution  $\hat{\mathbf{x}}$  is the actual solution of the system, hence  $\|A\hat{\mathbf{x}} - \mathbf{b}\| = 0$ . In the case that  $\mathbf{b}$  does not belong to  $R(A)$ , the linear system  $A\mathbf{x} = \mathbf{b}$  is inconsistent. In such a case the least squares solution  $\hat{\mathbf{x}}$  is the vector in  $\mathbb{R}^n$  with the property that  $A\hat{\mathbf{x}}$  is a vector in  $R(A)$  closest to  $\mathbf{b}$  in the inner product space  $(\mathbb{R}^m, \cdot)$ . A sketch of this situation for a matrix  $A \in \mathbb{R}^{3,2}$  is given in Fig. 46.

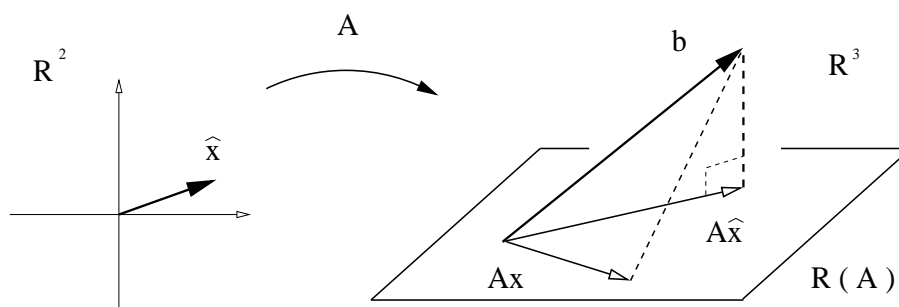


FIGURE 46. The meaning of the least squares solution  $\hat{\mathbf{x}} \in \mathbb{R}^2$  for the  $3 \times 2$  inconsistent linear system  $A\mathbf{x} = \mathbf{b}$  is that the vector  $A\hat{\mathbf{x}}$  is the closest to  $\mathbf{b}$  in the inner product space  $(\mathbb{R}^3, \cdot)$ .

The solution to the problem of finding a least squares solution to a linear system is summarized in the following result.

**Theorem 7.2.2.** Given a matrix  $A \in \mathbb{F}^{m,n}$  and a vector  $\mathbf{b}$  in the inner product space  $(\mathbb{F}^m, \cdot)$ , the vector  $\hat{\mathbf{x}} \in \mathbb{F}^n$  is a least squares solution of the  $m \times n$  linear system  $A\mathbf{x} = \mathbf{b}$  iff  $\hat{\mathbf{x}}$  is solution to the  $n \times n$  linear system, called **normal equation**,

$$A^*A\hat{\mathbf{x}} = A^*\mathbf{b}. \quad (7.3)$$

Furthermore, the least squares solution  $\hat{\mathbf{x}}$  is unique iff the column vectors of matrix  $A$  form a linearly independent set.

**REMARK:** In the case that  $\mathbb{F} = \mathbb{R}$ , the normal equation reduces to  $A^T A \hat{\mathbf{x}} = A^T \mathbf{b}$ .

**Proof of Theorem 7.2.2:** We are interested in finding a vector  $\hat{\mathbf{x}} \in \mathbb{F}^n$  such that  $A\hat{\mathbf{x}}$  is the best approximation in  $R(A)$  of vector  $\mathbf{b} \in \mathbb{F}^m$ . That is, we want to find  $\hat{\mathbf{x}} \in \mathbb{F}^n$  such that

$$\|A\hat{\mathbf{x}} - \mathbf{b}\| \leq \|y - \mathbf{b}\| \quad \forall y \in R(A).$$

Theorem 7.1.2 says that the best approximation of  $\mathbf{b}$  is when  $A\hat{\mathbf{x}} = \mathbf{b}_H$ , where  $\mathbf{b}_H$  is the orthogonal projection of  $\mathbf{b}$  onto the subspace  $R(A)$ . This means that

$$(A\hat{\mathbf{x}} - \mathbf{b}) \in R(A)^\perp = N(A^*) \Leftrightarrow A^*(A\hat{\mathbf{x}} - \mathbf{b}) = \mathbf{0}.$$

We then conclude that  $\hat{\mathbf{x}}$  must be solution of the normal equation

$$A^*A\hat{\mathbf{x}} = A^*\mathbf{b}.$$

The furthermore can be shown as follows. The column vectors of matrix  $A$  form a linearly independent set iff  $N(A) = \{\mathbf{0}\}$ . Lemma 7.2.3 stated below establishes that, for all matrix  $A$  holds that  $N(A) = N(A^*A)$ . This result in our case implies that  $N(A^*A) = \{\mathbf{0}\}$ . Since matrix  $A^*A$  is a square,  $n \times n$ , matrix, we conclude that it is invertible. This is equivalent to say that the solution  $\hat{\mathbf{x}}$  to the normal equation is unique; moreover, it is given by  $\hat{\mathbf{x}} = (A^*A)^{-1}A^*\mathbf{b}$ . This establishes the Theorem.  $\square$

In the proof of Theorem 7.2.2 above we used the following result:

**Lemma 7.2.3.** *If  $A \in \mathbb{F}^{m,n}$ , then  $N(A) = N(A^*A)$ .*

**Proof of Lemma 7.2.3:** We first show that  $N(A) \subset N(A^*A)$ . Indeed,

$$\mathbf{x} \in N(A) \Rightarrow A\mathbf{x} = \mathbf{0} \Rightarrow A^*A\mathbf{x} = \mathbf{0} \Rightarrow \mathbf{x} \in N(A^*A).$$

Now, suppose that there exists  $\mathbf{x} \in N(A^*A)$  such that  $\mathbf{x} \notin N(A)$ . Therefore,  $\mathbf{x} \neq \mathbf{0}$  and  $A\mathbf{x} \neq \mathbf{0}$ , which imply that

$$0 \neq \|A\mathbf{x}\|^2 = \mathbf{x}^*A^*A\mathbf{x} \Rightarrow A^*A\mathbf{x} \neq \mathbf{0}.$$

However, this last equation contradicts the assumption that  $\mathbf{x} \in N(A^*A)$ . Therefore, we conclude that  $N(A) = N(A^*A)$ . This establishes the Lemma.  $\square$

**EXAMPLE 7.2.1:** Show that the  $3 \times 2$  linear system  $A\mathbf{x} = \mathbf{b}$  is inconsistent; then find a least squares solutions  $\hat{\mathbf{x}} = \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \end{bmatrix}$  to that system, where

$$A = \begin{bmatrix} 1 & 3 \\ 2 & 2 \\ 3 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix}.$$

**SOLUTION:** We first show that the linear system above is inconsistent, since Gauss operation on the augmented matrix  $[A|\mathbf{b}]$  imply

$$\left[ \begin{array}{cc|c} 1 & 3 & -1 \\ 2 & 2 & 1 \\ 3 & 1 & -1 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 3 & -1 \\ 0 & -4 & 3 \\ 0 & -8 & 2 \end{array} \right] \rightarrow \left[ \begin{array}{cc|c} 1 & 3 & -1 \\ 0 & -4 & 3 \\ 0 & 0 & 1 \end{array} \right].$$

In order to find the least squares solution to the system above we first construct the normal equation. We need to compute

$$A^T A = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 2 & 2 \\ 3 & 1 \end{bmatrix} = \begin{bmatrix} 14 & 10 \\ 10 & 14 \end{bmatrix}, \quad A^T \mathbf{b} = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} = \begin{bmatrix} -2 \\ -2 \end{bmatrix}.$$

Therefore, the normal equation is given by

$$\begin{bmatrix} 14 & 10 \\ 10 & 14 \end{bmatrix} \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \end{bmatrix} = \begin{bmatrix} -2 \\ -2 \end{bmatrix}.$$

Since the column vectors of  $A$  form a linearly independent set, matrix  $A^T A$  is invertible,

$$(A^T A)^{-1} = \frac{1}{96} \begin{bmatrix} 14 & -10 \\ -10 & 14 \end{bmatrix} = \frac{1}{48} \begin{bmatrix} 7 & -5 \\ -5 & 7 \end{bmatrix}.$$

The least squares solution is unique and given by

$$\hat{\mathbf{x}} = \frac{1}{24} \begin{bmatrix} 7 & -5 \\ -5 & 7 \end{bmatrix} \begin{bmatrix} -1 \\ -1 \end{bmatrix} \Rightarrow \hat{\mathbf{x}} = \frac{1}{12} \begin{bmatrix} -1 \\ -1 \end{bmatrix}.$$

**REMARK:** We now verify that  $(A\hat{\mathbf{x}} - \mathbf{b}) \in R(A)^\perp$ . Indeed,

$$A\hat{\mathbf{x}} - \mathbf{b} = \begin{bmatrix} 1 & 3 \\ 2 & 2 \\ 3 & 1 \end{bmatrix} \frac{1}{12} \begin{bmatrix} -1 \\ -1 \end{bmatrix} - \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} = -\frac{1}{3} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} \Rightarrow A\hat{\mathbf{x}} - \mathbf{b} = \frac{2}{3} \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}.$$

Since

$$\begin{bmatrix} 1 & -2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = 1 - 4 + 3 = 0, \quad \begin{bmatrix} 1 & -2 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix} = 3 - 4 + 1 = 0,$$

we have verified that  $(A\hat{\mathbf{x}} - \mathbf{b}) \in R(A)^\perp$ .  $\triangleleft$

We finish this Subsection with an alternative proof of Theorem 7.2.2 in the particular case that involves real-valued matrices, that is,  $\mathbb{F} = \mathbb{R}$ . The proof is interesting in its own, since it is based in solving a constrained minimization problem.

**Alternative proof of Theorem 7.2.2 for  $\mathbb{F} = \mathbb{R}$ :** The vector  $\hat{\mathbf{x}} \in \mathbb{R}^n$  is a least squares solution of the system  $A\mathbf{x} = \mathbf{b}$  iff the function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  given by  $f(\mathbf{x}) = \|A\hat{\mathbf{x}} - \mathbf{b}\|^2$  has a minimum at  $\mathbf{x} = \hat{\mathbf{x}}$ . We then find all minima of function  $f$ . We first express  $f$  as follows,

$$\begin{aligned} f(\mathbf{x}) &= (A\mathbf{x} - \mathbf{b}) \cdot (A\mathbf{x} - \mathbf{b}) \\ &= (A\mathbf{x}) \cdot (A\mathbf{x}) - 2\mathbf{b} \cdot (A\mathbf{x}) + \mathbf{b} \cdot \mathbf{b} \\ &= \mathbf{x}^T A^T A \mathbf{x} - 2\mathbf{b}^T A \mathbf{x} + \mathbf{b}^T \mathbf{b}. \end{aligned}$$

We now need to find all solutions to the equation  $\nabla_{\mathbf{x}} f(\mathbf{x}) = 0$ . Recalling the definition of a gradient vector

$$\nabla_{\mathbf{x}} f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix},$$

it is simple to see that, for any vector  $\mathbf{a} \in \mathbb{R}^n$ , holds,

$$\nabla_{\mathbf{x}}(\mathbf{a}^T \mathbf{x}) = \mathbf{a}, \quad \nabla_{\mathbf{x}}(\mathbf{x}^T \mathbf{a}) = \mathbf{a}.$$

Therefore, the gradient of  $f$  is given by

$$\nabla_{\mathbf{x}} f = 2A^T A \mathbf{x} - 2A^T \mathbf{b}.$$

We are interested in the stationary points, the  $\hat{\mathbf{x}}$  solutions of

$$\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}) = 0 \Leftrightarrow A^T A \hat{\mathbf{x}} = A^T \mathbf{b}.$$

We conclude that all stationary points  $\hat{\mathbf{x}}$  are solutions of the normal equation, Eq. (7.3). These stationary points must be a minimum of  $f$ , since  $f$  is quadratic on the vector components  $x_i$  having the degree two terms all positive coefficients. This establishes the first part of Theorem 7.2.2 in the case that  $\mathbb{F} = \mathbb{R}$ .  $\square$

**7.2.2. Least squares fit.** It is often desirable to construct a mathematical model to describe the results of an experiment. This may involve fitting an algebraic curve to the given experimental data. The least squares method can be used to find the best parameters that fit the data.

**EXAMPLE 7.2.2: (Linear fit)** The simplest situation is the case where the best curve fitting the data is a straight line. More precisely, suppose that the result of an experiment is the following collection of ordered numbers

$$\{(t_1, b_1), \dots, (t_m, b_m)\}, \quad m \geq 2,$$

and suppose that a plot on a plane of the result of this experiment is given in Fig. 47. (For example, from measuring the vertical displacement  $b_i$  in a spring when a weight  $t_i$  is attached to it.) Find the best line  $y(t) = \hat{x}_2 t + \hat{x}_1$  that approximate these points in least squares sense. The latter means to find the numbers  $\hat{x}_2, \hat{x}_1 \in \mathbb{R}$  such that  $\sum_{i=1}^m |\Delta b_i|^2$  is the smallest possible, where

$$\Delta b_i = b_i - y(t_i) \quad \Leftrightarrow \quad \Delta b_i = b_i - (\hat{x}_2 t_i + \hat{x}_1), \quad i = 1, \dots, m.$$

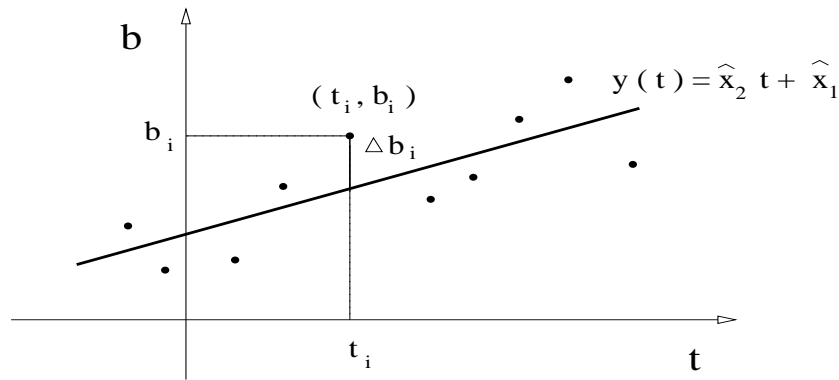


FIGURE 47. Sketch of the best line  $y(t) = \hat{x}_2 t + \hat{x}_1$  fitting the set of points  $(t_i, b_i)$ , for  $i = 1, \dots, 10$ .

**SOLUTION:** Let us rewrite this problem as the least squares solution of an  $m \times 2$  linear system, which in general is inconsistent. We are interested to find  $\hat{x}_2, \hat{x}_1$  solution of the linear system

$$\begin{array}{rcl} y(t_1) = b_1 & \hat{x}_1 + t_1 \hat{x}_2 = b_1 & \\ \vdots & \vdots & \\ y(t_m) = b_m & \hat{x}_1 + t_m \hat{x}_2 = b_m & \end{array} \quad \Leftrightarrow \quad \begin{bmatrix} 1 & t_1 \\ \vdots & \vdots \\ 1 & t_m \end{bmatrix} \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}.$$

Introducing the notation

$$\mathbf{A} = \begin{bmatrix} 1 & t_1 \\ \vdots & \vdots \\ 1 & t_m \end{bmatrix}, \quad \hat{\mathbf{x}} = \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix},$$



we are then interested in finding the solution  $\hat{\mathbf{x}}$  of the  $m \times 2$  linear system  $\mathbf{A}\hat{\mathbf{x}} = \mathbf{b}$ . Introducing also the vector

$$\Delta\mathbf{b} = \begin{bmatrix} \Delta b_1 \\ \vdots \\ \Delta b_m \end{bmatrix},$$

it is clear that  $\mathbf{A}\hat{\mathbf{x}} - \mathbf{b} = \Delta\mathbf{b}$ , and so we obtain the important relation

$$\|\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}\|^2 = \|\Delta\mathbf{b}\|^2 = \sum_{i=1}^m (\Delta b_i)^2.$$

Therefore, the vector  $\hat{\mathbf{x}}$  that minimizes the square of the deviation from the line,  $\sum_{i=1}^m (\Delta b_i)^2$ , is precisely the same vector  $\hat{\mathbf{x}} \in \mathbb{R}^2$  that minimizes the number  $\|\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}\|^2$ . We studied the latter problem at the beginning of this Section. We called it a least squares problem, and the solution  $\hat{\mathbf{x}}$  is the solution of the normal equation

$$\mathbf{A}^T \mathbf{A} \hat{\mathbf{x}} = \mathbf{A}^T \mathbf{b}.$$

It is simple to see that

$$\mathbf{A}^T \mathbf{A} = \begin{bmatrix} 1 & \cdots & 1 \\ t_1 & \cdots & t_m \end{bmatrix} \begin{bmatrix} 1 & t_1 \\ \vdots & \vdots \\ 1 & t_m \end{bmatrix} = \begin{bmatrix} m & \sum t_i \\ \sum t_i & \sum t_i^2 \end{bmatrix},$$

$$\mathbf{A}^T \mathbf{b} = \begin{bmatrix} 1 & \cdots & 1 \\ t_1 & \cdots & t_m \end{bmatrix} \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix} = \begin{bmatrix} \sum b_i \\ \sum t_i b_i \end{bmatrix}.$$

Therefore, we are interested in finding the solution to the  $2 \times 2$  linear system

$$\begin{bmatrix} m & \sum t_i \\ \sum t_i & \sum t_i^2 \end{bmatrix} \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \end{bmatrix} = \begin{bmatrix} \sum b_i \\ \sum t_i b_i \end{bmatrix}.$$

Suppose that at least one of the  $t_i$  is different from 1, then matrix  $\mathbf{A}^T \mathbf{A}$  is invertible and the inverse is

$$(\mathbf{A}^T \mathbf{A})^{-1} = \frac{1}{m \sum t_i^2 - (\sum t_i)^2} \begin{bmatrix} \sum t_i^2 & -\sum t_i \\ -\sum t_i & m \end{bmatrix}.$$

We conclude that the solution to the normal equation is

$$\begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \end{bmatrix} = \frac{1}{m \sum t_i^2 - (\sum t_i)^2} \begin{bmatrix} \sum t_i^2 & -\sum t_i \\ -\sum t_i & m \end{bmatrix} \begin{bmatrix} \sum b_i \\ \sum t_i b_i \end{bmatrix}.$$

So, the slope  $\hat{x}_2$  and vertical intercept  $\hat{x}_1$  of the best fitting line are given by

$$\hat{x}_2 = \frac{m \sum t_i b_i - (\sum t_i)(\sum b_i)}{m \sum t_i^2 - (\sum t_i)^2}, \quad \hat{x}_1 = \frac{(\sum t_i^2)(\sum b_i) - (\sum t_i)(\sum t_i b_i)}{m \sum t_i^2 - (\sum t_i)^2}.$$

◁

**EXAMPLE 7.2.3: (Polynomial fit)** Find the best polynomial of degree  $(n - 1) \geq 0$ , say  $p(t) = \hat{x}_n t^{(n-1)} + \cdots + \hat{x}_1$ , that approximates in least squares sense the set of points

$$\{(t_1, b_1), \dots, (t_m, b_m)\}, \quad m \geq n.$$

(See Fig. 48 for an example in the case that the fitting curve is a parabola,  $n = 3$ .) Following Example 7.2.2, the least squares approximation means to find the numbers  $\hat{x}_n, \dots, \hat{x}_1 \in \mathbb{R}$  such that  $\sum_{i=1}^m |\Delta b_i|^2$  is the smallest possible, where

$$\Delta b_i = b_i - p(t_i) \quad \Leftrightarrow \quad \Delta b_i = b_i - (\hat{x}_n t_i^{(n-1)} + \cdots + \hat{x}_1), \quad i = 1, \dots, m.$$

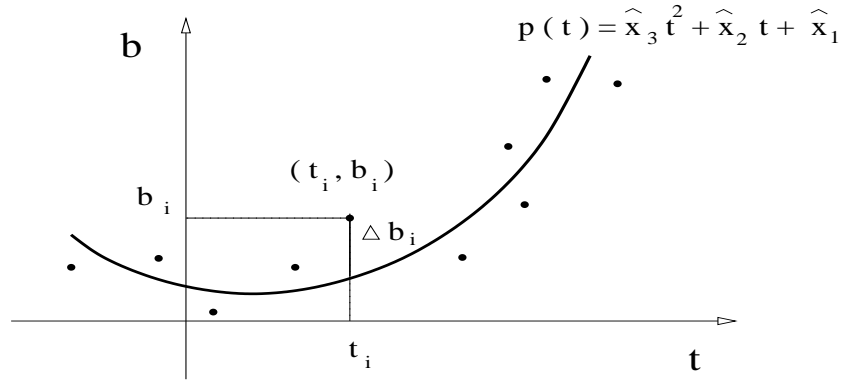


FIGURE 48. Sketch of the best parabola  $p(t) = \hat{x}_3 t^2 + \hat{x}_2 t + \hat{x}_1$  fitting the set of points  $(t_i, b_i)$ , for  $i = 1, \dots, 10$ .

**SOLUTION:** We rewrite this problem as the least squares solution of an  $m \times n$  linear system, which in general is inconsistent. We are interested to find  $\hat{x}_n, \dots, \hat{x}_1$  solution of the linear system

$$\begin{array}{rcl} p(t_1) = b_1 & \hat{x}_1 + \dots + t_1^{(n-1)} \hat{x}_n = b_1 & \Leftrightarrow \begin{bmatrix} 1 & \dots & t_1^{(n-1)} \end{bmatrix} \begin{bmatrix} \hat{x}_1 \\ \vdots \\ \hat{x}_n \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix} \\ \vdots & \vdots & \\ p(t_m) = b_m & \hat{x}_1 + \dots + t_m^{(n-1)} \hat{x}_n = b_m & \Leftrightarrow \begin{bmatrix} 1 & \dots & t_m^{(n-1)} \end{bmatrix} \begin{bmatrix} \hat{x}_1 \\ \vdots \\ \hat{x}_n \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix} \end{array}$$

Introducing the notation

$$A = \begin{bmatrix} 1 & \dots & t_1^{(n-1)} \\ \vdots & & \vdots \\ 1 & \dots & t_m^{(n-1)} \end{bmatrix}, \quad \hat{\mathbf{x}} = \begin{bmatrix} \hat{x}_1 \\ \vdots \\ \hat{x}_n \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix},$$

we are then interested in finding the solution  $\hat{\mathbf{x}}$  of the  $m \times n$  linear system  $A\hat{\mathbf{x}} = \mathbf{b}$ . Introducing also the vector

$$\Delta \mathbf{b} = \begin{bmatrix} \Delta b_1 \\ \vdots \\ \Delta b_m \end{bmatrix},$$

it is clear that  $A\hat{\mathbf{x}} - \mathbf{b} = \Delta \mathbf{b}$ , and so we obtain the important relation

$$\|A\hat{\mathbf{x}} - \mathbf{b}\|^2 = \|\Delta \mathbf{b}\|^2 = \sum_{i=1}^m (\Delta b_i)^2.$$

Therefore, the vector  $\hat{\mathbf{x}}$  that minimizes the square of the deviation from the line,  $\sum_{i=1}^m (\Delta b_i)^2$ , is precisely the same vector  $\hat{\mathbf{x}} \in \mathbb{R}^2$  that minimizes the number  $\|A\hat{\mathbf{x}} - \mathbf{b}\|^2$ . We studied the latter problem at the beginning of this Section. We called it a least squares problem, and the solution  $\hat{\mathbf{x}}$  is the solution of the normal equation

$$A^T A \hat{\mathbf{x}} = A^T \mathbf{b}. \quad (7.4)$$

It is simple to see that Eq. (7.4) is an  $n \times n$  linear system, since

$$\mathbf{A}^T \mathbf{A} = \begin{bmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ t_1^{(n-1)} & \cdots & t_m^{(n-1)} \end{bmatrix} \begin{bmatrix} 1 & \cdots & t_1^{(n-1)} \\ \vdots & & \vdots \\ 1 & \cdots & t_m^{(n-1)} \end{bmatrix},$$

$$\mathbf{A}^T \mathbf{b} = \begin{bmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ t_1^{(n-1)} & \cdots & t_m^{(n-1)} \end{bmatrix} \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}.$$

We do not compute these expressions explicitly here. In the case that the columns of  $\mathbf{A}$  form a linearly independent set, the solution  $\hat{\mathbf{x}}$  to the normal equation is

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}.$$

The components of  $\hat{\mathbf{x}}$  provide the parameters for the best polynomial fitting the data in least squares sense.  $\triangleleft$

**7.2.3. Linear correlation.** In statistics a correlation coefficient measures the departure of two random variables from independence. For centered data, that is, for data with zero average, the correlation coefficient can be viewed as the cosine of the angle in an abstract  $\mathbb{R}^n$  space between two vectors constructed with the random variables data. We now define and find the correlation coefficient for two variables as given in Example 7.2.2.

Once again, suppose that the result of an experiment is the following collection of ordered numbers

$$\{(t_1, b_1), \dots, (t_m, b_m)\}, \quad m \geq 2. \quad (7.5)$$

Introduce the vectors  $\mathbf{e}$ ,  $\mathbf{t}$ ,  $\mathbf{b} \in \mathbb{R}^m$  as follows,

$$\mathbf{e} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}, \quad \mathbf{t} = \begin{bmatrix} t_1 \\ \vdots \\ t_m \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}.$$

Before introducing the correlation coefficient, let us use these vectors above to write down the least squares coefficients  $\mathbf{x}$  found in Example 7.2.2. The matrix of coefficients can be written as  $\mathbf{A} = [\mathbf{e}, \mathbf{t}]$ , therefore,

$$\mathbf{A}^T \mathbf{A} = \begin{bmatrix} \mathbf{e}^T \\ \mathbf{t}^T \end{bmatrix} [\mathbf{e}, \mathbf{t}] = \begin{bmatrix} \mathbf{e} \cdot \mathbf{e} & \mathbf{t} \cdot \mathbf{e} \\ \mathbf{e} \cdot \mathbf{t} & \mathbf{t} \cdot \mathbf{t} \end{bmatrix}, \quad \mathbf{A}^T \mathbf{b} = \begin{bmatrix} \mathbf{e}^T \\ \mathbf{t}^T \end{bmatrix} \mathbf{b} = \begin{bmatrix} \mathbf{e} \cdot \mathbf{b} \\ \mathbf{t} \cdot \mathbf{b} \end{bmatrix}.$$

The least squares solution can be written as follows,

$$\begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \end{bmatrix} = \frac{1}{(\mathbf{e} \cdot \mathbf{e})(\mathbf{t} \cdot \mathbf{t}) - (\mathbf{t} \cdot \mathbf{e})^2} \begin{bmatrix} \mathbf{t} \cdot \mathbf{t} & -\mathbf{e} \cdot \mathbf{t} \\ -\mathbf{e} \cdot \mathbf{t} & \mathbf{e} \cdot \mathbf{e} \end{bmatrix} \begin{bmatrix} \mathbf{e} \cdot \mathbf{b} \\ \mathbf{t} \cdot \mathbf{b} \end{bmatrix},$$

that is,

$$\hat{x}_2 = \frac{(\mathbf{e} \cdot \mathbf{e})(\mathbf{t} \cdot \mathbf{b}) - (\mathbf{e} \cdot \mathbf{t})(\mathbf{e} \cdot \mathbf{b})}{(\mathbf{e} \cdot \mathbf{e})(\mathbf{t} \cdot \mathbf{t}) - (\mathbf{t} \cdot \mathbf{e})^2}, \quad \hat{x}_1 = \frac{(\mathbf{e} \cdot \mathbf{b})(\mathbf{t} \cdot \mathbf{t}) - (\mathbf{e} \cdot \mathbf{t})(\mathbf{t} \cdot \mathbf{b})}{(\mathbf{e} \cdot \mathbf{e})(\mathbf{t} \cdot \mathbf{t}) - (\mathbf{t} \cdot \mathbf{e})^2}.$$

Introduce the average values

$$\bar{t} = \frac{\mathbf{e} \cdot \mathbf{t}}{\mathbf{e} \cdot \mathbf{e}}, \quad \bar{b} = \frac{\mathbf{e} \cdot \mathbf{b}}{\mathbf{e} \cdot \mathbf{e}}.$$

These are indeed the average values of  $t_i$  and  $b_i$ , since

$$\mathbf{e} \cdot \mathbf{e} = m, \quad \mathbf{e} \cdot \mathbf{t} = \sum_{i=1}^m t_i, \quad \mathbf{e} \cdot \mathbf{b} = \sum_{i=1}^m b_i.$$

Introduce the zero-average vectors  $\hat{\mathbf{t}} = (\mathbf{t} - \bar{t}\mathbf{e})$  and  $\hat{\mathbf{b}} = (\mathbf{b} - \bar{b}\mathbf{e})$ . The *correlation coefficient* of the data given in (7.5) is given by

$$\text{cor}(\mathbf{t}, \mathbf{b}) = \frac{\hat{\mathbf{t}} \cdot \hat{\mathbf{b}}}{\|\hat{\mathbf{t}}\| \|\hat{\mathbf{b}}\|}.$$

Therefore, the correlation coefficient between the data vectors  $\mathbf{t}$  and  $\mathbf{b}$  is the angle between the zero-average vectors  $\hat{\mathbf{t}}$  and  $\hat{\mathbf{b}}$  in  $\mathbb{R}^m$ .

In order to understand what measures this angle, let us consider the case where all the ordered pairs in (7.5) lies on a line, that is, there exists a solution  $\hat{\mathbf{x}}$  of the linear system  $\mathbf{A}\hat{\mathbf{x}} = \mathbf{b}$  (a solution, not only a least squares solution). In that case we have

$$\hat{x}_1\mathbf{e} + \hat{x}_2\mathbf{t} = \mathbf{b} \quad \Rightarrow \quad \hat{x}_1 + \hat{x}_2\bar{t} = \bar{b},$$

and this implies that

$$\hat{x}_2(\mathbf{t} - \bar{t}\mathbf{e}) = (\mathbf{b} - \bar{b}\mathbf{e}) \quad \Leftrightarrow \quad \hat{x}_2\hat{\mathbf{t}} = \hat{\mathbf{b}} \quad \Leftrightarrow \quad \text{cor}(\mathbf{t}, \mathbf{b}) = 1.$$

That is, in the case that  $\mathbf{t}$  is linearly related to  $\mathbf{b}$  we obtain that the zero-average vectors  $\hat{\mathbf{t}}$  and  $\hat{\mathbf{b}}$  are parallel, so the correlation coefficient is equal one.

**7.2.4. QR-factorization.** The Gram-Schmidt method can be used to factor any  $m \times n$  matrix  $\mathbf{A}$  into a product of an  $m \times n$  matrix  $\mathbf{Q}$  with orthonormal column vectors and an upper triangular  $n \times n$  matrix  $\mathbf{R}$ . We will see that the QR-factorization is useful to solve the normal equation associated to a least squares problem.

**Theorem 7.2.4.** *If the column vectors of matrix  $\mathbf{A} \in \mathbb{F}^{m,n}$  form a linearly independent set, then there exist matrices  $\mathbf{Q} \in \mathbb{F}^{m,n}$  and  $\mathbf{R} \in \mathbb{F}^{n,n}$  satisfying that  $\mathbf{Q}^*\mathbf{Q} = \mathbf{I}_n$ , matrix  $\mathbf{R}$  is upper triangular with positive diagonal elements, and the following equation holds*

$$\mathbf{A} = \mathbf{QR}.$$

**Proof of Theorem 7.2.4:** Use the Gram-Schmidt method to obtain an orthonormal set  $\{\mathbf{q}_1, \dots, \mathbf{q}_n\}$  from the column vectors of the  $m \times n$  matrix  $\mathbf{A} = [\mathbf{A}_{:1}, \dots, \mathbf{A}_{:n}]$ , that is,

$$\begin{aligned} \mathbf{p}_1 &= \mathbf{A}_{:1} & \mathbf{q}_1 &= \frac{\mathbf{p}_1}{\|\mathbf{p}_1\|}, \\ \mathbf{p}_2 &= \mathbf{A}_{:2} - (\mathbf{A}_{:2} \cdot \mathbf{q}_1) \mathbf{q}_1 & \mathbf{q}_2 &= \frac{\mathbf{p}_2}{\|\mathbf{p}_2\|}, \\ &\vdots & &\vdots \\ \mathbf{p}_n &= \mathbf{A}_{:n} - (\mathbf{A}_{:n} \cdot \mathbf{q}_1) \mathbf{q}_1 - \dots - (\mathbf{A}_{:n} \cdot \mathbf{q}_{n-1}) \mathbf{q}_{n-1} & \mathbf{q}_n &= \frac{\mathbf{p}_n}{\|\mathbf{p}_n\|}. \end{aligned}$$

Define matrix  $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_n]$ , which then satisfies the equation  $\mathbf{Q}^*\mathbf{Q} = \mathbf{I}_n$ . Notice that the equations above can be expressed as follows,

$$\begin{aligned} \mathbf{A}_{:1} &= \|\mathbf{p}_1\| \mathbf{q}_1, \\ \mathbf{A}_{:2} &= \|\mathbf{p}_2\| \mathbf{q}_2 + (\mathbf{q}_1 \cdot \mathbf{A}_{:2}) \mathbf{q}_1 \\ &\vdots \\ \mathbf{A}_{:n} &= \|\mathbf{p}_n\| \mathbf{q}_n + (\mathbf{q}_1 \cdot \mathbf{A}_{:n}) \mathbf{q}_1 + (\mathbf{q}_2 \cdot \mathbf{A}_{:n}) \mathbf{q}_2 + \dots + (\mathbf{q}_{n-1} \cdot \mathbf{A}_{:n}) \mathbf{q}_{n-1}. \end{aligned}$$

After some time staring at the equations above, one can rewrite it as a matrix product

$$[A_{:1}, \dots, A_{:n}] = [q_1, \dots, q_n] \begin{bmatrix} \|p_1\| & (q_1 \cdot A_{:2}) & \cdots & (q_1 \cdot A_{:n}) \\ 0 & \|p_2\| & \cdots & (q_2 \cdot A_{:n}) \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & (q_{n-1} \cdot A_{:n}) \\ 0 & 0 & \cdots & \|p_n\| \end{bmatrix} \quad (7.6)$$

Define matrix  $R$  by equation above as the matrix satisfying  $A = QR$ . Then, Eq. (7.6) says that matrix  $R$  is  $n \times n$ , upper triangular, with positive diagonal elements. This establishes the Theorem.  $\square$

**EXAMPLE 7.2.4:** Find the QR-factorization of matrix  $A = \begin{bmatrix} 1 & 2 & 1 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$ .

**SOLUTION:** First use the Gram-Schmidt method to transform the column vectors of matrix  $A$  into an orthonormal set. This was done in Example 6.5.1. The result defines the matrix  $Q$  as follows

$$Q = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

Having matrix  $A$  and  $Q$ , and knowing that Theorem 7.2.4 is true, then we can compute matrix  $R$  by the equation  $R = Q^T A$ . Since the column vectors of  $Q$  form an orthonormal set, we have that  $Q^T = Q^{-1}$ , and in this particular case  $Q^{-1} = Q$ , so matrix  $R$  is given by

$$R = QA = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 1 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix} \Rightarrow R = \begin{bmatrix} \sqrt{2} & \frac{3}{\sqrt{2}} & \sqrt{2} \\ 0 & \frac{1}{\sqrt{2}} & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

The QR-factorization of matrix  $A$  is then given by

$$A = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \sqrt{2} & \frac{3}{\sqrt{2}} & \sqrt{2} \\ 0 & \frac{1}{\sqrt{2}} & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

$\triangleleft$

The QR-factorization is useful to solve the normal equation in a least squares problem.

**Theorem 7.2.5.** Assume that the matrix  $A \in \mathbb{F}^{m,n}$  admits the QR-factorization  $A = QR$ . The vector  $\hat{x} \in \mathbb{F}^n$  is solution of the normal equation  $A^*A\hat{x} = A^*b$  iff it is solution of

$$R\hat{x} = Q^*b.$$

**Proof of Theorem 7.2.5:** Just introduce the QR-factorization into the normal equation  $A^*A\hat{x} = A^*b$  as follows,

$$(R^*Q^*)(QR)\hat{x} = R^*Q^*b \Leftrightarrow R^*R\hat{x} = R^*Q^*b \Leftrightarrow R^*(R\hat{x} - Q^*b) = 0.$$

Since  $R$  is a square, upper triangular matrix with non-zero coefficients, we conclude that  $R$  is invertible. Therefore, from the last equation above we conclude that  $\hat{x}$  is solution of the normal equation iff holds

$$R\hat{x} = Q^*b.$$

This establishes the Theorem.  $\square$

## 7.2.5. Exercises.

**7.2.1.-** Consider the matrix  $A$  and the vector  $\mathbf{b}$  given by

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}.$$

- (a) Find the least-squares solution  $\hat{\mathbf{x}}$  to the linear system  $A\mathbf{x} = \mathbf{b}$ .  
 (b) Verify that the solution  $\hat{\mathbf{x}}$  satisfies

$$(A\hat{\mathbf{x}} - \mathbf{b}) \in R(A)^\perp.$$

**7.2.2.-** Consider the matrix  $A$  and the vector  $\mathbf{b}$  given by

$$A = \begin{bmatrix} 2 & 2 \\ 0 & -1 \\ -2 & 0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix}.$$

- (a) Find the least-squares solution  $\hat{\mathbf{x}}$  to the linear system  $A\mathbf{x} = \mathbf{b}$ .  
 (b) Find the orthogonal projection of the source vector  $\mathbf{b}$  onto the subspace  $R(A)$ .

**7.2.3.-** Find all the least-squares solutions  $\hat{\mathbf{x}}$  to the linear system  $A\mathbf{x} = \mathbf{b}$ , where

$$A = \begin{bmatrix} 1 & 2 \\ 2 & 4 \\ 3 & 6 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}.$$

**7.2.4.-** Find the best line in least-squares sense that fits the measurements, where  $t_1$  is the independent variable and  $b_i$  is the dependent variable,

$$\begin{aligned} t_1 &= -2, & b_1 &= 4, \\ t_2 &= -1, & b_2 &= 3, \\ t_3 &= 0, & b_3 &= 1, \\ t_4 &= 2, & b_4 &= 0. \end{aligned}$$

**7.2.5.-** Find the correlation coefficient corresponding to the measurements given in Exercise **7.2.4** above.

**7.2.6.-** Use Gram-Schmidt method on the columns of matrix  $A$  below to find its QR factorization, where

$$A = \begin{bmatrix} 1 & 1 \\ 2 & 3 \\ 2 & 1 \end{bmatrix}.$$

**7.2.7.-** Find the QR factorization of matrix

$$A = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix}.$$

## 7.3. FINITE DIFFERENCE METHOD

A differential equation is an equation where the unknown is a function and both the function itself and its derivatives appear in the equation. The differential equation is called linear iff the unknown function and its derivatives appear linearly in the equation. Solutions of linear differential equations can be approximated by solutions of appropriate  $n \times n$  algebraic linear systems in the limit that  $n$  approaches infinity. The finite difference method is a way to obtain an  $n \times n$  linear system from the original differential equation. *Derivatives are approximated by difference quotients, thus reducing a differential equation to an algebraic linear system.* Since a derivative can be approximated by infinitely many difference quotients, there are infinitely many  $n \times n$  linear systems that approximate a differential equation. One tries to choose the linear system whose solution is the best approximation of the solution of the original differential equation. Computers are used to find the vector in  $\mathbb{R}^n$  solution of the  $n \times n$  linear system. Many approximations of the solution to the differential equation are obtained from this array of  $n$  numbers. One way to obtain a function from a vector in  $\mathbb{R}^n$  is to find a degree  $n$  polynomial that contains all these  $n$  points. This is called a polynomial interpolation of the algebraic solution. In this Section we only show how to obtain  $n \times n$  algebraic linear systems that approximate a simple differential equation.

**7.3.1. Differential equations.** A differential equation is an equation where the unknown is a function and both the function and its derivatives appear in the equation. A simple example is the following: Given a continuously differentiable function  $f : [0, 1] \rightarrow \mathbb{R}$ , find a function  $u : [0, 1] \rightarrow \mathbb{R}$  solution of the differential equation

$$\frac{du}{dx}(x) = f(x),$$

To find a solution to a differential equation requires to perform appropriate integrations, thus integration constants are introduced in the solution. This suggests that the solution of a differential equation is not unique, and extra conditions must be added to the problem to select only one solution. In the differential equation above the solutions  $u$  are given by

$$u(x) = \int_0^x f(t) dt + c,$$

with  $c \in \mathbb{R}$ . An extra condition is needed to obtain a unique solution, for example the condition  $u(0) = 1$ . Then, the unique solution  $u$  is computed as follows

$$1 = u(0) = \int_0^0 f(t) dt + c \Rightarrow c = 1 \Rightarrow u(x) = \int_0^x f(t) dt + 1.$$

The example above is simple enough that no approximation is needed to obtain the solution.

An *ordinary differential equation* is a differential equation where the unknown function  $u$  depends on one variable, as in the example above. A *partial differential equation* is a differential equation where the unknown function depends on more than one variable and the equation contains derivatives of more than one variable. In this Section we use the finite difference method to find a solution to two different problems. The first one involves an ordinary differential equation while the second one involves a partial differential equation, called the heat equation.

The first problem is to find an approximate solution to a boundary value problem for an ordinary differential equation: Given a continuously differentiable function  $f : [0, 1] \rightarrow \mathbb{R}$ ,

find a function  $u : [0, 1] \rightarrow \mathbb{R}$  solution of the boundary value problem

$$\frac{d^2 u}{dx^2}(x) + \frac{du}{dx}(x) = f(x), \quad (7.7)$$

$$u(0) = u(1) = 0. \quad (7.8)$$

Finding the function  $u$  involves doing two integrations that introduce two integration constants. These constants are determined by the two conditions on  $x = 0$  and  $x = 1$  above, called boundary conditions, since they are conditions on the boundaries of the interval  $[0, 1]$ . Because of this extra condition on the ordinary differential equation this problem is called a boundary value problem.

The second problem we study in this Section is to find an approximate solution to an initial-boundary value problem for the heat equation: Given the set  $D = [0, \pi] \times [0, T]$ , the positive constant  $\kappa$ , and infinitely differentiable functions  $f : D \rightarrow \mathbb{R}$  and  $g : [0, \pi] \rightarrow \mathbb{R}$ , find the function  $u : D \rightarrow \mathbb{R}$  solution of the problem

$$\frac{\partial u}{\partial t}(x, t) - \kappa \frac{\partial^2 u}{\partial x^2}(x, t) = f(x, t) \quad (x, t) \in D, \quad (7.9)$$

$$u(0, t) = u(\pi, t) = 0 \quad t \in [0, T], \quad (7.10)$$

$$u(x, 0) = g(x) \quad x \in [0, \pi]. \quad (7.11)$$

The partial differential equation in Eq. (7.9) is called the one-dimensional heat equation, since in the case that function  $u$  is the temperature of a material that depends on time  $t$  and one spatial direction  $x$ , the equation describes how the material temperature changes in time due to heat propagation. The positive constant  $\kappa$  is called the thermal diffusivity of the material. The condition given in Eq. (7.10) is called a boundary condition, since they are conditions on  $x = 0$  and  $x = \pi$  that hold for all  $t \in [0, T]$ , see Fig. 49. The condition given in Eq. (7.11) is called an initial condition, since it is a condition on the initial time  $t = 0$  for all  $x \in [0, \pi]$ , see Fig. 49. Because of these two extra conditions on the partial differential equation this problem is called an initial-boundary value problem for the heat equation.

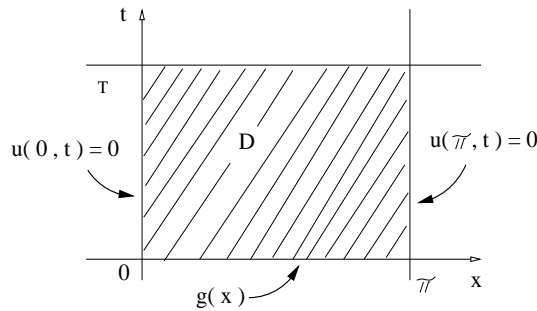


FIGURE 49. The domain  $D = [0, \pi] \times [0, T]$  where the initial-boundary value problem for the heat equation is set up. We indicate the boundary data conditions  $u(0, t) = 0$  and  $u(\pi, t) = 0$ , and the initial data function  $g$ .

**7.3.2. Difference quotients.** Finite difference methods transform a linear differential equation into an  $n \times n$  algebraic linear system by replacing derivatives by difference quotients.



Derivatives can be approximated by difference quotients in many different ways. For example, a derivative of a function  $u$  can be expressed in the following equivalent ways,

$$\begin{aligned}\frac{du}{dx}(x) &= \lim_{\Delta x \rightarrow 0} \frac{u(x + \Delta x) - u(x)}{\Delta x}, \\ &= \lim_{\Delta x \rightarrow 0} \frac{u(x) - u(x - \Delta x)}{\Delta x}, \\ &= \lim_{\Delta x \rightarrow 0} \frac{u(x + \Delta x) - u(x - \Delta x)}{2\Delta x}.\end{aligned}$$

However, for a fixed nonzero value of  $\Delta x$ , the expressions below are, in general, different.

**Definition 7.3.1.** The *forward difference quotient*  $d_+$  and the *backward difference quotient*  $d_-$  of a continuous function  $u : \mathbb{R} \rightarrow \mathbb{R}$  at  $x \in \mathbb{R}$  are given by

$$d_+u(x) = \frac{u(x + \Delta x) - u(x)}{\Delta x}, \quad d_-u(x) = \frac{u(x) - u(x - \Delta x)}{\Delta x}.$$

The *centered difference quotient*  $d_c$  of a continuous function  $u$  at  $x \in \mathbb{R}$  is given by

$$d_cu(x) = \frac{u(x + \Delta x) - u(x - \Delta x)}{2\Delta x}. \quad (7.12)$$

In the case that the function  $u$  has second continuous derivative, the Taylor Expansion Theorem implies that the forward and backward differences differ from the actual derivative by terms order  $\Delta x$ . The proof is not difficult, since

$$\begin{aligned}u(x + \Delta x) &= u(x) + \frac{du}{dx}(x) \Delta x + O((\Delta x)^2) \quad \Rightarrow \quad d_+u(x) = \frac{du}{dx}(x) + O(\Delta x), \\ u(x - \Delta x) &= u(x) - \frac{du}{dx}(x) \Delta x + O((\Delta x)^2) \quad \Rightarrow \quad d_-u(x) = \frac{du}{dx}(x) + O(\Delta x),\end{aligned}$$

where  $O((\Delta x)^n)$  denotes a function satisfying  $[O((\Delta x)^n)]/(\Delta x)^n$  approaches a constant as  $\Delta x \rightarrow 0$ . In the case that the function  $u$  has third continuous derivative, the Taylor Expansion Theorem implies that the centered difference quotient differs from the actual derivative by terms of order  $(\Delta x)^2$ . Again, the proof is not difficult, and it is based on the Taylor expansion of the function  $u$ . Compute this expansion in two different ways,

$$u(x + \Delta x) = u(x) + \frac{du}{dx}(x) \Delta x + \frac{1}{2} \frac{d^2u}{dx^2}(x) (\Delta x)^2 + O((\Delta x)^3), \quad (7.13)$$

$$u(x - \Delta x) = u(x) - \frac{du}{dx}(x) \Delta x + \frac{1}{2} \frac{d^2u}{dx^2}(x) (\Delta x)^2 + O((\Delta x)^3), \quad (7.14)$$

Subtracting the two expressions above we obtain that

$$u(x + \Delta x) - u(x - \Delta x) = 2 \frac{du}{dx}(x) \Delta x + O((\Delta x)^3) \quad \Rightarrow \quad d_cu(x) = \frac{du}{dx}(x) + O((\Delta x)^2),$$

which establishes that centered difference quotients differ from the derivative by order  $(\Delta x)^2$ . If a function is infinitely continuously differentiable, centered difference quotients are more accurate than forward or backward differences.

Second and higher derivatives of a function can also be approximated by difference quotients. Again, there are infinitely many ways to approximate second derivatives by difference quotients. The freedom to choose difference quotients is higher for second derivatives than for first derivatives. We now present two difference quotients to give an idea of this freedom. On the one hand, one possible approximation is to use the centered difference quotient twice, since a second derivative is the derivative of the derivative function. Using the more precise notation

$$d_{c\Delta x}u(x) = \frac{u(x + \Delta x) - u(x - \Delta x)}{2\Delta x},$$

it is not difficult to check that  $(d_{c\Delta x})^2 = d_{c\Delta x}(d_{c\Delta x})$  is given by

$$(d_{c\Delta x})^2 u(x) = \frac{1}{(2\Delta x)^2} [u(x + 2\Delta x) + u(x - 2\Delta x) - 2u(x)]. \quad (7.15)$$

On the other hand, another approximation for the second derivative of a function can be obtained directly from the Taylor expansion formulas in (7.13)-(7.14). Indeed, add Eqs. (7.13)-(7.14), that is,

$$u(x + \Delta x) + u(x - \Delta x) = 2u(x) + \frac{d^2 u}{dx^2}(x) (\Delta x)^2 + O((\Delta x)^3).$$

This equation can be rewritten as

$$\frac{d^2 u}{dx^2}(x) = \frac{u(x + \Delta x) + u(x - \Delta x) - 2u(x)}{(\Delta x)^2} + O(\Delta x).$$

This equation suggests to introduce a second-order centered difference quotient  $(d^2)_{c\Delta x}$  as

$$(d^2)_{c\Delta x} u(x) = \frac{1}{(\Delta x)^2} [u(x + \Delta x) + u(x - \Delta x) - 2u(x)]. \quad (7.16)$$

Using this notation, the equation above is given by

$$\frac{d^2 u}{dx^2}(x) = (d^2)_{c\Delta x} u(x) + O(\Delta x).$$

Therefore, both  $(d_{c\Delta x})^2$  and  $(d^2)_{c\Delta x}$  are approximations of the second derivative of a function. However, they are not the same approximation, since comparing Eqs. (7.15) and (7.16) it is not difficult to see that

$$(d_{c(\Delta x)/2})^2 = (d^2)_{c\Delta x}.$$

We conclude that there are many different ways to approximate second derivatives by difference quotients, many more ways than those to approximate first order derivatives. In this Section we use the difference quotient in Eq. (7.16), and we use the simplified notation given by  $d_c^2 = (d^2)_{c\Delta x}$ .

**7.3.3. Method of finite differences.** We now describe the finite difference method using two examples. In the first example we find an approximate solution for the boundary value problem in Eqs. (7.7)-(7.8). In the second example we find an approximate solution for the initial-boundary value problem in Eqs. (7.9)-(7.11).

**EXAMPLE 7.3.1:** Consider the boundary value problem for the ordinary differential equation given in Eqs. (7.7)-(7.8).

- Divide the interval  $[0, 1]$  into  $n > 1$  equal intervals and use the finite difference method to find a vector  $\mathbf{u} = [u_i] \in \mathbb{R}^{n+1}$  that approximates the function  $u : [0, 1] \rightarrow \mathbb{R}$  solution of that boundary value problem. Use centered difference quotients to approximate the first and second derivatives of the unknown function  $u$ .
- Find the explicit form of the linear system in the case  $n = 6$ .
- Find the degree  $n$  polynomial  $p_n$  that interpolates the approximate solution vector  $\mathbf{u} = [u_i] \in \mathbb{R}^{n+1}$ , which includes the boundary points.

**SOLUTION:**

**Part (a):** Fix a positive integer  $n \in \mathbb{N}$ , define the grid step size  $h = 1/n$ , and introduce the uniform grid

$$\{x_i\}, \quad x_i = ih, \quad i = 0, 1, \dots, n, \quad \text{on } [0, 1].$$

(See Fig. 50.) Introduce the numbers  $f_i = f(x_i)$ . Finally, denote  $u_i = u(x_i)$ . The numbers  $x_i$  and  $f_i$  are known from the problem, and so are the  $u_0 = u(0) = 0$  and  $u_n = u(1) = 0$ ,

while the  $u_i$  for  $i = 1, \dots, n - 1$  are the unknowns. We now use the original differential equation to construct an  $(n - 1) \times (n - 1)$  linear system for these unknowns  $u_i$ .

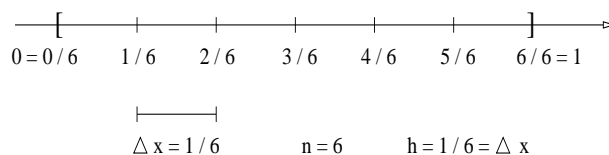


FIGURE 50. A uniform grid for  $n = 6$  on the domain  $[0, 1]$ .

We use centered difference quotient given in Eqs. (7.12) and (7.16) to approximate the first and second derivatives of the function  $u$ , respectively. We choose  $\Delta x = h$ , and we denote the difference quotients evaluated at grid points  $x_i$  as follows,

$$d_c u(x_i) = d_c u_i, \quad d_c^2 u(x_i) = d_c^2 u_i.$$

Therefore, we obtain the following formulas for the difference quotients,

$$d_c u_i = \frac{u_{i+1} - u_{i-1}}{2h}, \quad d_c^2 u_i = \frac{u_{i+1} + u_{i-1} - 2u_i}{h^2}. \tag{7.17}$$

Now we state the approximate problem we will solve: Given the constants  $\{f_i\}_{i=1}^{n-1}$ , find the vector  $\mathbf{u} = [u_i] \in \mathbb{R}^{n+1}$  solution of the  $(n - 1) \times (n - 1)$  linear system and boundary conditions, respectively,

$$d_c^2 u_i + d_c u_i = f_i, \quad i = 1, \dots, n - 1, \tag{7.18}$$

$$u_0 = u_n = 0. \tag{7.19}$$

Eq. (7.18) is indeed a linear system for  $\mathbf{u} = [u_i]$ , since it is equivalent to the system

$$(2 + h)u_{i+1} - 4u_i + (2 - h)u_{i-1} = 2h^2 f_i, \quad i = 1, \dots, n - 1.$$

When the boundary conditions in Eq. (7.19) are introduced in the equation above, we obtain an  $(n - 1) \times (n - 1)$  linear system for the unknowns  $u_i$ , where  $i = 1, \dots, (n - 1)$ .

**Part (b):** In the case  $n = 6$ , we have  $h = 1/6$ , so we denote  $a = 2 - 1/6$ ,  $b = 2 + 1/6$  and  $c = 2/36$ . Also recall the boundary conditions  $u_0 = u_6 = 0$ . Then, the system above and its augmented matrix are given by, respectively,

$$\begin{aligned} -4u_1 + bu_2 &= cf_1, \\ au_1 - 4u_2 + bu_3 &= cf_2, \\ au_2 - 4u_3 + bu_4 &= cf_3, \\ au_3 - 4u_4 + bu_5 &= cf_4, \\ au_4 - 4u_5 &= cf_5, \end{aligned} \quad \Leftrightarrow \quad \left[ \begin{array}{ccccc|c} -4 & b & 0 & 0 & 0 & cf_1 \\ a & -4 & b & 0 & 0 & cf_2 \\ 0 & a & -4 & b & 0 & cf_3 \\ 0 & 0 & a & -4 & b & cf_4 \\ 0 & 0 & 0 & a & -4 & cf_5 \end{array} \right].$$

We then conclude that the solution  $u$  of the boundary value problem in Eq. (7.7) can be approximated by the solution  $\mathbf{u} = [u_i] \in \mathbb{R}^7$  of the  $5 \times 5$  linear system above plus the two boundary conditions. The same type of approximate solution can be found for all  $n \in \mathbb{N}$ .

**Part (c):** The output of the finite difference method is a vector  $\mathbf{u} = [u_i] \in \mathbb{R}^{(n+1)}$ . An approximate solution to the ordinary differential equation in Eqs. (7.7) can be constructed from the vector  $\mathbf{u}$  in many different ways. One way is polynomial interpolation, that is,

to construct a polynomial of degree  $n$  whose graph contains all the points  $(x_i, u_i)$ . Such polynomial is given by

$$p_n(x) = \sum_{i=0}^n u_i q_i(x), \quad q_i(x) = \prod_{j \neq i} \frac{(x - x_j)}{(x_i - x_j)}.$$

It can be verified that the degree  $n$  polynomials  $q_i$  when evaluated at grid points satisfies

$$q_i(x_j) = \begin{cases} 0 & \text{if } i \neq j, \\ 1 & \text{if } i = j. \end{cases}$$

Therefore, the polynomial  $p_n$  has degree  $n$  and satisfies that  $p_n(x_i) = u_i$ . This polynomial function  $p_n$  approximates the solution  $u$  of the boundary value problem in Eqs. (7.7)-(7.8). ◁

**EXAMPLE 7.3.2:** Consider the boundary value problem for the partial differential equation given in Eqs. (7.9)-(7.11). Use the finite difference method to find an approximate solution of the function  $u : D \rightarrow \mathbb{R}$  solution of that initial-boundary value problem.

- (a) Use centered difference quotients to approximate the spatial derivatives and forward difference quotients to approximate time derivatives of the unknown function  $u$ .
- (b) Repeat the calculations in part (a) now using backward difference quotients for the time derivatives of the unknown function  $u$ .

**SOLUTION:**

**Part (a):** Introduce a grid in the domain  $D = [0, \pi] \times [0, T]$  as follows: Fix the positive integers  $n_x, n_t \in \mathbb{N}$ , define the step sizes  $h_x = \pi/n_x$  and  $h_t = T/n_t$ , and then introduce the uniform grids

$$\begin{aligned} \{x_i\}, & \quad x_i = ih_x, & \quad i = 0, 1, \dots, n_x, & \quad \text{on } [0, \pi] \\ \{t_j\}, & \quad t_j = jh_t, & \quad j = 0, 1, \dots, n_t, & \quad \text{on } [0, T]. \end{aligned}$$

A point of the form  $(x_i, t_j) \in D$  is called a grid point. (See Fig. 51.)

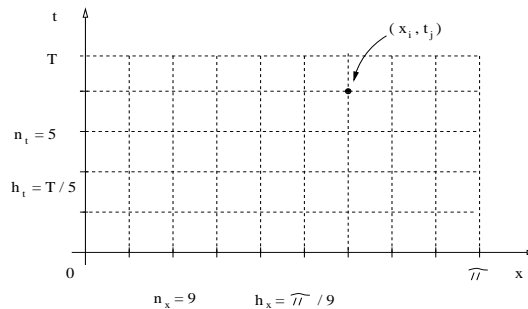


FIGURE 51. A rectangular grid with  $n_x = 9$  and  $n_t = 5$  on the domain  $D = [0, \pi] \times [0, T]$ .

We will compute the approximate solution values at grid points, and we use the notation  $u_{i,j} = u(x_i, t_j)$  and  $f_{i,j} = f(x_i, t_j)$  to denote unknown and source function values at grid points. We also denote  $g_i = g(x_i)$  for the initial data function values at grid points. In this notation the boundary conditions in Eq. (7.10) have the form  $u_{0,j} = 0$  and  $u_{n_x,j} = 0$ . We

finally introduce the forward difference quotient in time  $d_{+t}$  and the second order centered difference quotient in space  $d_{cx}^2$  as follows,

$$d_{+t}u_{i,j} = \frac{u_{i,(j+1)} - u_{i,j}}{h_t}, \quad d_{cx}^2u_{i,j} = \frac{u_{(i+1),j} + u_{(i-1),j} - 2u_{i,j}}{h_x^2}.$$

Now we state the approximate problem we will solve: Given the constants  $f_{i,j}$  and  $g_i$ , for  $i = 0, 1, \dots, n_x$  and  $j = 0, 1, \dots, n_y$ , find the constants  $u_{i,j}$  solution of the linear equations and boundary conditions, respectively,

$$d_{+t}u_{i,j} - \kappa d_{cx}^2u_{i,j} = f_{i,j}, \quad (7.20)$$

$$u_{0,j} = u_{n_x,j} = 0, \quad (7.21)$$

$$u_{i,0} = g_i. \quad (7.22)$$

This system is simpler to solve than it looks at first sight. Let us rewrite it in as follows,

$$\frac{1}{h_t}(u_{i,(j+1)} - u_{i,j}) - \frac{\kappa}{h_x^2}(u_{(i+1),j} + u_{(i-1),j} - 2u_{i,j}) = f_{i,j},$$

which is equivalent to the equation

$$u_{i,(j+1)} = r(u_{(i+1),j} + u_{(i-1),j}) + (1 - 2r)u_{i,j} + h_t f_{i,j}, \quad (7.23)$$

where  $r = \kappa h_t / h_x^2$ . This last equation says that the solution at the time  $t_{j+1}$  can be computed if the solution at the previous time  $t_j$  is known. Since the solution at the initial time  $t_0 = 0$  is known an equal to  $g_i$ , and since solution at the boundary of the domain  $u_{0,j}$  and  $u_{n_x,j}$  is known from the boundary conditions, then the solution  $u_{i,j}$  can be computed from Eq. (7.23) time step after time step. For this reason the linear system in Eqs. (7.20)-(7.22) above is an example of an *explicit method*.

**Part (b):** We now repeat the calculations in part (a) using a backward difference quotient in time. We introduce the notation  $d_{-t}$  for the backward difference quotient in time, and we keep the notation  $d_{cx}^2$  for the second order centered difference quotient in space,

$$d_{-t}u_{i,j} = \frac{u_{i,j} - u_{i,(j-1)}}{h_t}, \quad d_{cx}^2u_{i,j} = \frac{u_{(i+1),j} + u_{(i-1),j} - 2u_{i,j}}{h_x^2}.$$

Now we state the approximate problem we will solve: Given the constants  $f_{i,j}$  and  $g_i$ , for  $i = 0, 1, \dots, n_x$  and  $j = 0, 1, \dots, n_y$ , find the constants  $u_{i,j}$  solution of the linear equations and boundary conditions, respectively,

$$d_{-t}u_{i,j} - \kappa d_{cx}^2u_{i,j} = f_{i,j}, \quad (7.24)$$

$$u_{0,j} = u_{n_x,j} = 0, \quad (7.25)$$

$$u_{i,0} = g_i. \quad (7.26)$$

This system is simpler to solve than it looks at first sight. Let us rewrite it in as follows,

$$\frac{1}{h_t}(u_{i,j} - u_{i,(j-1)}) - \frac{\kappa}{h_x^2}(u_{(i+1),j} + u_{(i-1),j} - 2u_{i,j}) = f_{i,j},$$

which is equivalent to the equation

$$(+2r)u_{i,j} - r(u_{(i+1),j} + u_{(i-1),j}) = u_{i,(j-1)} + h_t f_{i,j}, \quad (7.27)$$

where  $r = \kappa h_t / h_x^2$ , as above. This last equation says that the solution at the time  $t_{j+1}$  can be computed if the solution at the previous time  $t_j$  is known. However, in this case we need to solve an  $(n_x - 1) \times (n_x - 1)$  linear system at each time step. Such system is similar to the one that appeared in Example 7.3.1. Since the solution at the initial time  $t_0 = 0$  is known an equal to  $g_i$ , and since solution at the boundary of the domain  $u_{0,j}$  and  $u_{n_x,j}$  is known from the boundary conditions, then the solution  $u_{i,j}$  can be computed from Eq. (7.27) time step

after time step. We emphasize that the solution is computed by solving an  $(n_x - 1) \times (n_x - 1)$  system at every time step. For this reason the linear system in Eqs. (7.24)-(7.26) above is an example of an *implicit method*.  $\triangleleft$

What we have seen in this Section is just the first part of the story. We have seen how we can use linear algebra to obtain approximate solutions of few problems involving differential equations. The second part is to study how the solutions of the approximate problems approaches the solution of the original problem as the grid step size approaches zero. Consider the approximate solution  $\mathbf{u} \in \mathbb{R}^{n+1}$  found in Example 7.3.1. Does the interpolation polynomial  $p_n$  constructed with the components of  $\mathbf{u}$  approximate the function  $u : [0, 1] \rightarrow \mathbb{R}$  solution to the boundary value problem in Eq. (7.7) in the limit  $n \rightarrow \infty$ ? A similar question can be asked for the solutions  $\{u_{i,j}\}$  obtained in parts (b) and (c) in Example 7.3.2. We will study the answers to these questions in the following Chapters.

One last remark is the following. Comparing parts (a) and (b) in Example 7.3.2 we see that explicit methods are simpler to solve than implicit methods. A matrix must be inverted to solve an implicit method, while this is not needed to solve an explicit method. So, why are implicit methods studied at all? The reason is that in the limit  $n \rightarrow \infty$  the approximate solutions of explicit methods do not approximate the solution of the original differential equation as good as a solution of an implicit method. Moreover, the solution of an explicit method may not converge at all, while the solutions of implicit methods always converges.

**Further reading.** See Section 1.4 in Meyer's book [3] for a detailed discussion on discretizations of two-point boundary values problems.

## 7.3.4. Exercises.

- 7.3.1.- Consider the boundary value problem for the function  $u$  given by

$$\begin{aligned}\frac{d^2 u}{dx^2}(x) &= 25x, \\ u(0) &= 0, \quad u(1) = 0, \\ x &\in [0, 1].\end{aligned}$$

Divide the interval  $[0, 1]$  into five equal subintervals and use the finite difference method to find an approximate solution vector  $\mathbf{u} = [u_i]$  to the boundary value problem above, where  $i = 0, \dots, 5$ . Use **centered** difference quotients given in Eq. (7.17) to approximate the derivatives of function  $u$ .

- 7.3.2.- Given an infinite differentiable function  $u : \mathbb{R} \rightarrow \mathbb{R}$ . apply twice the forward difference quotient  $d_+$  to show that the second order forward difference quotient has the form

$$\frac{d_+^2 u(x) = \frac{u(x + 2\Delta x) - 2u(x + \Delta x) + u(x)}{(\Delta x)^2}.$$

- 7.3.3.- Given an infinite differentiable function  $u : \mathbb{R} \rightarrow \mathbb{R}$ . apply twice the backward difference quotient  $d_-$  to show that the second order backward difference quotient has the form

$$\frac{d_-^2 u(x) = \frac{u(x) - 2u(x - \Delta x) + u(x - 2\Delta x)}{(\Delta x)^2}.$$

- 7.3.4.- Consider the boundary value problem given in the Exercise 7.3.1. Divide again the interval  $[0, 1]$  into five equal subintervals and find the  $4 \times 4$  linear system that approximates the original problem using **forward** difference quotients to approximate derivatives of function  $u$ . You do not need to solve the linear system.

- 7.3.5.- Consider the boundary value problem given Problem 7.3.1. Divide again the interval  $[0, 1]$  into five equal subintervals and find the  $4 \times 4$  linear system that approximates the original problem using **backward** difference quotients to approximate derivatives of function  $u$ . You do not need to solve the linear system.

- 7.3.6.- Consider the boundary value problem for the function  $u$  given by

$$\begin{aligned}\frac{d^2 u}{dx^2}(x) + 2\frac{du}{dx}(x) &= 25x, \\ u(0) &= 0, \quad u(1) = 0, \\ x &\in [0, 1].\end{aligned}$$

Divide the interval  $[0, 1]$  into five equal subintervals and use the finite difference method to find an algebraic linear system that approximates the original boundary value problem above. Use centered difference quotients given in Eq. (7.17) to approximate derivatives of function  $u$ . You do not need to solve the linear system.

## 7.4. FINITE ELEMENT METHOD

The finite element method permits the computation of approximate solutions to differential equations. Boundary value problems involving a differential equations are first transformed into integral equations. The boundary conditions are included into the integral equation by performing integration by parts. The original problem is transformed into inverting a bilinear form on a vector space. The approximate problem is obtained when the integral equation is solved not on the whole vector space but on a finite dimensional subspace. By a careful choice of the subspace, the calculations needed to obtain the approximate solution can be simplified. In this Section we study the same differential equations we have seen in Sect. 7.3 when we described the finite difference methods. Finite difference methods are different from finite element methods. The former method approximates derivatives by difference quotients in order to obtain the approximate problem involving an algebraic linear system. The latter method produces an algebraic linear system by restricting an integral version of the original differential equation onto a finite dimensional subspace of the original vector space where the integral equation is defined.

**7.4.1. Differential equations.** We now recall the boundary value problem we are interested to study. This is the first problem we studied in Sect. 7.3, that is, to find an approximate solution to a boundary value problem for an ordinary differential equation: Given an infinitely differentiable function  $f : [0, 1] \rightarrow \mathbb{R}$ , find a function  $u : [0, 1] \rightarrow \mathbb{R}$  solution of the boundary value problem

$$\frac{d^2 u}{dx^2}(x) + \frac{du}{dx}(x) = f(x), \quad (7.28)$$

$$u(0) = u(1) = 0. \quad (7.29)$$

Recall that finding the function  $u$  involves doing two integrations that introduce two integration constants. These constants are determined by the two conditions on  $x = 0$  and  $x = 1$  above, called boundary conditions, since they are conditions on the boundaries of the interval  $[0, 1]$ . Because of this extra condition on the ordinary differential equation this problem is called a boundary value problem.

This problem can be expressed using linear transformations on vector spaces. Consider the inner product spaces  $(V, \langle \cdot, \cdot \rangle)$  and  $(W, \langle \cdot, \cdot \rangle)$ , where  $V = C_0^\infty([0, 1], \mathbb{R})$  is the space of infinitely many differentiable functions that vanish at  $x = 0$  and  $x = 1$ ,  $W = C^\infty([0, 1], \mathbb{R})$  is the space of infinitely many differentiable functions, while the inner product is defined as

$$\langle \mathbf{f}, \mathbf{g} \rangle = \int_0^1 \mathbf{f}(x)\mathbf{g}(x) dx.$$

Introduce the linear transformation  $L : V \rightarrow W$  defined by

$$L(\mathbf{v}) = \frac{d^2 \mathbf{v}}{dx^2} + \frac{d\mathbf{v}}{dx}$$

Then, the boundary value problem in Eqs. (7.28)-(7.29) can be expressed in the following way: Given a vector  $\mathbf{f} \in W$ , find a vector  $\mathbf{u} \in V$  solution of the equation

$$L(\mathbf{u}) = \mathbf{f}. \quad (7.30)$$

Notice that the boundary conditions in Eq. (7.29) have been included into the definition of the vector space  $V$ . The problem defined by Eqs. (7.28)-(7.29), which is equivalently defined by Eq. (7.30), is called the **strong formulation** of the problem.

So, we have expressed a boundary value problem for a differential equation in terms of linear transformations between appropriate infinite dimensional vector spaces. The next step is to use the inner product defined on  $V$  and  $W$  to transform the differential equation



in (7.30) into an integro-differential equation, which will be called the weak formulation of the problem. This idea is summarized below.

**7.4.2. The Galerkin method.** The Galerkin method refers to a collection of ideas to transform a problem involving a linear transformation between infinite dimensional inner product spaces into a problem involving a matrix as a function between finite dimensional subspaces. We describe in this Section the original idea, introduced by Boris Galerkin around 1915. Galerkin worked only with partial differential equations, but we now know that his idea works in the more general context of a linear transformation between infinite dimensional inner product spaces. For this reason we describe Galerkin's idea in this more general context. The Galerkin method is to transform the strong formulation of the problem in (7.30) into what is called the weak formulation of the problem. This transformation is done using the inner product defined on the infinite dimensional vector spaces. Before we describe this transformation we need few definitions.

**7.4.3. Finite element method.**

**7.4.4. Exercises.**

7.4.1.- .

7.4.2.- .

## CHAPTER 8. NORMED SPACES

8.1. THE  $p$ -NORM

An inner product in a vector space always determines an inner product norm, which satisfies the properties (a)-(c) in Theorem 6.2.5. However, the inner product norm is not the only function  $V \rightarrow \mathbb{R}$  satisfying these properties.

**Definition 8.1.1.** A norm on a vector space  $V$  over the field  $\mathbb{F}$  is any function  $\|\cdot\| : V \rightarrow \mathbb{R}$  satisfying the following properties,

- (a) **(Positive definiteness)** For all  $\mathbf{x} \in V$  holds  $\|\mathbf{x}\| \geq 0$ , and  $\|\mathbf{x}\| = 0$  iff  $\mathbf{x} = 0$ ;
- (b) **(Scaling)** For all  $\mathbf{x} \in V$  and all  $a \in \mathbb{F}$  holds  $\|a\mathbf{x}\| = |a| \|\mathbf{x}\|$ ;
- (c) **(Triangle inequality)** For all  $\mathbf{x}, \mathbf{y} \in V$  holds  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ .

A normed space is a pair  $(V, \|\cdot\|)$  of a vector space with a norm.

The inner product norm,  $\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$  for all  $\mathbf{x} \in V$ , defined in an inner product space  $(V, \langle \cdot, \cdot \rangle)$ , is thus an obvious example of a norm, since it satisfies the properties (a)-(c) in Theorem 6.2.5, which are precisely the conditions given in Definition 8.1.1. It is important to notice that alternative norms exist on inner product spaces. Moreover, a norm can be introduced in a vector space without having an inner product structure. One particularly important example of the former case is given by the  $p$ -norms defined on  $V = \mathbb{F}^n$ .

**Definition 8.1.2.** The  $p$ -norm on the vector space  $\mathbb{F}^n$ , with  $1 \leq p \leq \infty$ , is the function  $\|\cdot\|_p : \mathbb{F}^n \rightarrow \mathbb{R}$  defined as follows,

$$\begin{aligned} \|\mathbf{x}\|_p &= (|x_1|^p + \cdots + |x_n|^p)^{1/p}, & p \in [1, \infty), \\ \|\mathbf{x}\|_\infty &= \max\{|x_1|, \dots, |x_n|\}, & (p = \infty), \end{aligned}$$

with  $\mathbf{x} = [x_i]$ , for  $i = 1, \dots, n$ , the vector components in the standard ordered basis of  $\mathbb{F}^n$ .

Since the dot product norm is given by  $\|\mathbf{x}\| = (|x_1|^2 + \cdots + |x_n|^2)^{1/2}$ , it is simple to see that  $\|\cdot\| = \|\cdot\|_2$ , that is, the case  $p = 2$  coincides with the dot product norm on  $\mathbb{F}^n$ . The most commonly used norms, besides  $p = 2$ , are the cases  $p = 1$  and  $p = \infty$ ,

$$\|\mathbf{x}\|_1 = |x_1| + \cdots + |x_n|, \quad \|\mathbf{x}\|_\infty = \max\{|x_1|, \dots, |x_n|\}.$$

**Theorem 8.1.3.** For each value of  $p \in [1, \infty]$  the  $p$ -norm function  $\|\cdot\|_p : \mathbb{F}^n \rightarrow \mathbb{R}$  introduced in Definition 8.1.2 is a norm on  $\mathbb{F}^n$ .

The Theorem above states that the function  $\|\cdot\|_p$  satisfies the properties (a)-(c) in Definition 8.1.1. Therefore, for each value  $p \in [1, \infty]$  the space  $(\mathbb{F}^n, \|\cdot\|_p)$  is a normed space. We will see at the end of this Section that in the case  $p \in [1, \infty]$  and  $p \neq 2$  these norms are not inner product norms. In other words, for these values of  $p$  there is **no** inner product  $\langle \cdot, \cdot \rangle_p$  defined on  $\mathbb{F}^n$  such that  $\|\cdot\|_p = \sqrt{\langle \cdot, \cdot \rangle_p}$ .

In order to prove that the  $p$ -norm is indeed a norm, we first need to establish the following generalization of the Cauchy-Schwarz inequality, called Hölder inequality.

**Theorem 8.1.4. (Hölder inequality)** For all vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{F}^n$  and  $p \in [1, \infty)$  holds that

$$|\mathbf{x}^* \mathbf{y}| \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q, \quad \text{with} \quad \frac{1}{p} + \frac{1}{q} = 1. \quad (8.1)$$

**Proof of Theorem 8.1.4:** We first show that for all real numbers  $a \geq 0, b > 0$  holds

$$a^\lambda b^{1-\lambda} \leq (1-\lambda)b + \lambda a \quad \lambda \in [0, 1]. \quad (8.2)$$

This inequality can be shown using the auxiliary function

$$f(t) = (1 - \lambda) + \lambda t - t^\lambda.$$

Its derivative is  $f'(t) = \lambda[1 - t^{(\lambda-1)}]$ , so the function  $f$  satisfies

$$f(1) = 0, \quad f'(t) > 0 \quad \text{for } t > 1, \quad f'(t) < 0 \quad \text{for } 0 \leq t < 1.$$

We conclude that  $f(t) \geq 0$  for  $t \in [0, \infty)$ . Given two real numbers  $a \geq 0$ ,  $b > 0$ , we have proven that  $f(a/b) \geq 0$ , that is,

$$\frac{a^\lambda}{b^\lambda} \leq (1 - \lambda) + \lambda \frac{a}{b} \quad \Leftrightarrow \quad a^\lambda b^{1-\lambda} \leq (1 - \lambda)b + \lambda a.$$

Having established Eq. (8.2), we use it to prove Hölder inequality. Let  $\mathbf{x} = [x_i]$ ,  $\mathbf{y} = [y_j] \in \mathbb{F}^n$  be arbitrary vectors, where  $i, j = 1, \dots, n$ , and introduce the rescaled vectors

$$\hat{\mathbf{x}} = \frac{\mathbf{x}}{\|\mathbf{x}\|_p}, \quad \hat{\mathbf{y}} = \frac{\mathbf{y}}{\|\mathbf{y}\|_q}, \quad \frac{1}{p} + \frac{1}{q} = 1.$$

These rescaled vectors satisfy that  $\|\hat{\mathbf{x}}\|_p = 1$  and  $\|\hat{\mathbf{y}}\|_q = 1$ . Denoting  $\hat{x}_i = [\hat{x}_i]$  and  $\hat{y}_i = [\hat{y}_i]$ , use the inequality in Eq. (8.2) in the case that  $a = \hat{x}_i$ ,  $b = \hat{y}_i$  and  $\lambda = 1/p$ , as follows,

$$\begin{aligned} |\bar{\hat{x}}_i \hat{y}_i| &= \left( \frac{|x_i|^p}{\sum_{j=1}^n |x_j|^p} \right)^{\frac{1}{p}} \left( \frac{|y_i|^q}{\sum_{j=1}^n |y_j|^q} \right)^{1-\frac{1}{p}} \\ &\leq \left( 1 - \frac{1}{p} \right) \left( \frac{|y_i|^q}{\sum_{j=1}^n |y_j|^q} \right) + \frac{1}{p} \left( \frac{|x_i|^p}{\sum_{j=1}^n |x_j|^p} \right) \\ &\leq \frac{1}{q} |\hat{y}_i|^q + \frac{1}{p} |\hat{x}_i|^p. \end{aligned}$$

Adding up over all components,

$$|\hat{\mathbf{x}}^* \hat{\mathbf{y}}| \leq \sum_{i=1}^n |\bar{\hat{x}}_i \hat{y}_i| \leq \frac{1}{q} \|\hat{\mathbf{y}}\|_q^q + \frac{1}{p} \|\hat{\mathbf{x}}\|_p^p = \frac{1}{q} \|\hat{\mathbf{y}}\|_q^q + \frac{1}{p} \|\hat{\mathbf{x}}\|_p^p = \frac{1}{q} + \frac{1}{p} = 1.$$

Therefore,  $|\hat{\mathbf{x}}^* \hat{\mathbf{y}}| \leq 1$ , which is equivalent to

$$\left| \frac{\mathbf{x}^* \mathbf{y}}{\|\mathbf{x}\|_p \|\mathbf{y}\|_q} \right| \leq 1 \quad \Leftrightarrow \quad |\mathbf{x}^* \mathbf{y}| \leq \|\mathbf{x}\|_p \|\mathbf{y}\|_q.$$

This establishes the Theorem.  $\square$

The Hölder inequality plays an important role to show that the  $p$ -norms satisfy the triangle inequality. We saw the same situation when the Cauchy-Schwarz inequality played a crucial role to prove that the inner product norm satisfied the triangle inequality.

**Proof of Theorem 8.1.3:** We show the proof for  $p \in [1, \infty)$ . The case  $p = \infty$  is left as an exercise. So, we assume that  $p \in [1, \infty)$ , we introduce  $q$  by the equation  $\frac{1}{p} + \frac{1}{q} = 1$ . In order to show that the  $p$ -norm is a norm we need to show that the  $p$ -norm satisfies the properties (a)-(c) in Definition 8.1.1. The first two properties are simple to prove. The  $p$ -norm is positive, since for all  $\mathbf{x} \in \mathbb{F}^n$  holds

$$\|\mathbf{x}\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}} \geq 0, \quad \text{and} \quad \left( \sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}} = 0 \quad \Leftrightarrow \quad |x_i| = 0 \quad \Leftrightarrow \quad \mathbf{x} = \mathbf{0}.$$

The  $p$ -norm satisfies the scaling property, since for all  $\mathbf{x} \in \mathbb{F}^n$  and all  $a \in \mathbb{F}$  holds

$$\|a\mathbf{x}\|_p = \left( \sum_{i=1}^n |ax_i|^p \right)^{\frac{1}{p}} = \left( \sum_{i=1}^n |a|^p |x_i|^p \right)^{\frac{1}{p}} = \left( |a|^p \sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}} = |a| \|\mathbf{x}\|_p.$$

The difficult part is to establish that the  $p$ -norm satisfies the triangle inequality. We start proving the following statement: For all real numbers  $a, b$  holds

$$|a + b|^p \leq |a| |a + b|^{\frac{p}{q}} + |b| |a + b|^{\frac{p}{q}}. \quad (8.3)$$

This is indeed the case, since

$$|a + b|^p = |a + b| |a + b|^{p-1}, \quad \text{and} \quad \frac{1}{q} = 1 - \frac{1}{p} = \frac{p-1}{p} \quad \Leftrightarrow \quad \frac{p}{q} = p-1,$$

so we conclude that

$$|a + b| = |a + b| |a + b|^{\frac{p}{q}} \leq (|a| + |b|) |a + b|^{\frac{p}{q}},$$

which is the inequality in Eq. (8.3). This inequality will be used in the following calculation. Given arbitrary vectors  $\mathbf{x} = [x_i], \mathbf{y} = [y_i] \in \mathbb{F}^n$ , for  $i = 1, \dots, n$ , the following inequalities hold

$$\|\mathbf{x} + \mathbf{y}\|_p^p = \sum_{i=1}^n |x_i + y_i|^p \leq \sum_{i=1}^n \left( |x_i| |x_i + y_i|^{\frac{p}{q}} + |y_i| |x_i + y_i|^{\frac{p}{q}} \right), \quad (8.4)$$

where Eq. (8.3) was used to obtain the last inequality. Now the Hölder inequality implies both

$$\begin{aligned} \sum_{i=1}^n |x_i| |x_i + y_i|^{\frac{p}{q}} &\leq \left( \sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}} \left( \sum_{i=1}^n |x_i + y_i|^p \right)^{\frac{1}{q}}, \\ \sum_{i=1}^n |y_i| |x_i + y_i|^{\frac{p}{q}} &\leq \left( \sum_{i=1}^n |y_i|^p \right)^{\frac{1}{p}} \left( \sum_{i=1}^n |x_i + y_i|^p \right)^{\frac{1}{q}}. \end{aligned}$$

Inserting these expressions in Eq. (8.4) we obtain

$$\|\mathbf{x} + \mathbf{y}\|_p^p \leq \|\mathbf{x}\|_p \left( \|\mathbf{x} + \mathbf{y}\|_p \right)^{\frac{p}{q}} + \|\mathbf{y}\|_p \left( \|\mathbf{x} + \mathbf{y}\|_p \right)^{\frac{p}{q}}.$$

Recalling that  $\frac{p}{q} = p-1$ , we obtain the inequality

$$\|\mathbf{x} + \mathbf{y}\|_p^p \leq (\|\mathbf{x}\|_p + \|\mathbf{y}\|_p) \|\mathbf{x} + \mathbf{y}\|_p^{(p-1)} \quad \Leftrightarrow \quad \|\mathbf{x} + \mathbf{y}\|_p \leq \|\mathbf{x}\|_p + \|\mathbf{y}\|_p.$$

This establishes the Theorem for  $p \in [1, \infty)$ .  $\square$

**EXAMPLE 8.1.1:** Find the length of  $\mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ -3 \end{bmatrix}$  in the norms  $\|\cdot\|_1, \|\cdot\|_2$  and  $\|\cdot\|_\infty$ .

**SOLUTION:** A simple calculations shows,

$$\begin{aligned} \|\mathbf{x}\|_1 &= |1| + |2| + |-3| = 6, \\ \|\mathbf{x}\|_2 &= [1^2 + 2^2 + (-3)^2]^{1/2} = \sqrt{14} = 3.74, \\ \|\mathbf{x}\|_\infty &= \max\{|1|, |2|, |-3|\} = 3. \end{aligned}$$

Notice that for the vector  $\mathbf{x}$  above holds

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_1. \quad (8.5)$$

One can prove that the inequality in Eq. (8.5) holds for all  $\mathbf{x} \in \mathbb{F}^n$ .  $\triangleleft$

**EXAMPLE 8.1.2:** Sketch on  $\mathbb{R}^2$  the set of vectors  $B_p = \{\mathbf{x} \in \mathbb{R}^2 : \|\mathbf{x}\|_p = 1\}$  for the cases  $p = 1, p = 2$ , and  $p = \infty$ .

**SOLUTION:** Recall we use the standard basis to express  $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ . We start with the set  $B_2$ , which is the circle of radius one in Fig. 52, that is,

$$(x_1)^2 + (x_2)^2 = 1.$$

The set  $B_1$  is the square of side one given by

$$|x_1| + |x_2| = 1.$$

The sides are given by the lines  $\pm x_1 \pm x_2 = 1$ . See Fig. 52. The set  $B_\infty$  is the square of side two given by

$$\max\{|x_1|, |x_2|\} = 1.$$

The sides are given by the lines  $x_1 = \pm 1$  and  $x_2 = \pm 1$ . See Fig. 52.

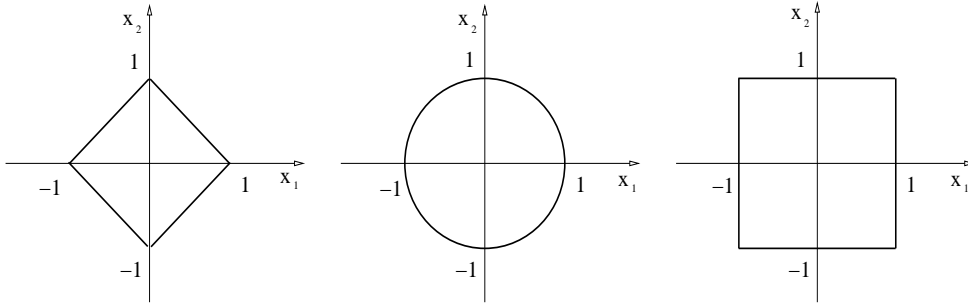


FIGURE 52. Unit sets in  $\mathbb{R}^2$  for the  $p$ -norms, with  $p = 1, 2, \infty$ , respectively.

◀

**EXAMPLE 8.1.3:** Show that for every  $\mathbf{x} \in \mathbb{R}^2$  holds

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_1.$$

**SOLUTION:** Introducing the unit disks  $D_p = \{\mathbf{x} \in \mathbb{R}^2 : \|\mathbf{x}\|_p \leq 1\}$ , with  $p = 1, 2, \infty$ , then Fig. 53 shows that  $D_1 \subset D_2 \subset D_\infty$ . Let us choose  $\mathbf{y} \in \mathbb{R}^2$  such that  $\|\mathbf{y}\|_2 = 1$ , that is, a vector on the circle, for example the vector given in second picture in Fig. 53. Since this vector is outside the disk  $D_1$ , that implies  $\|\mathbf{y}\|_1 \geq 1$ , and since this vector is inside the disk  $D_\infty$ , that implies  $\|\mathbf{y}\|_\infty \leq 1$ . The three conditions together say

$$\|\mathbf{y}\|_\infty \leq \|\mathbf{y}\|_2 \leq \|\mathbf{y}\|_1.$$

The equal signs correspond to the cases where  $\mathbf{y}$  is a horizontal or a vertical vector. Since any vector  $\mathbf{x} \in \mathbb{R}^2$  is a scaling of an appropriate vector  $\mathbf{y}$  on the border of  $D_2$ , that is,  $\mathbf{x} = c\mathbf{y}$ , with  $0 \leq c \in \mathbb{R}$ , then, multiplying the inequality above by  $c$  we obtain

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_1, \quad \forall \mathbf{x} \in \mathbb{R}^2.$$

◀

The  $p$ -norms can be defined on infinite dimensional vector spaces like function spaces.

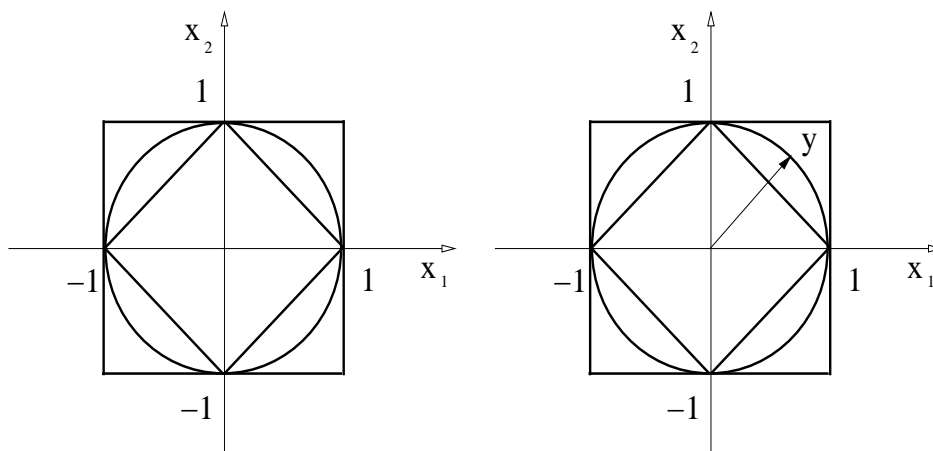


FIGURE 53. A comparison of the unit sets in  $\mathbb{R}^2$  for the  $p$ -norms, with  $p = 1, 2, \infty$ .

**Definition 8.1.5.** The  $p$ -norm on the vector space  $V = C^k([a, b], \mathbb{R})$ , with  $1 \leq p \leq \infty$ , is the function  $\|\cdot\|_p : V \rightarrow \mathbb{R}$  defined as follows,

$$\begin{aligned} \|\mathbf{f}\|_p &= \left( \int_a^b |\mathbf{f}(x)|^p dx \right)^{1/p}, & p \in [1, \infty), \\ \|\mathbf{f}\|_\infty &= \max_{x \in [a, b]} |\mathbf{f}(x)|, & (p = \infty). \end{aligned}$$

One can show that the  $p$ -norms introduced in Definition 8.1.5 are indeed norms on the vector space  $C^k([a, b], \mathbb{R})$ . The proof of this statement follows the same ideas given in the proofs of Theorems 8.1.4 and 8.1.3 above, and we do not present it in these notes.

**EXAMPLE 8.1.4:** Consider the normed space  $(C^0([-1, 1], \mathbb{R}), \|\cdot\|_p)$  for any  $p \in [1, \infty]$  and find the  $p$ -norm of the element  $\mathbf{f}(x) = x$ .

**SOLUTION:** In the case of  $p \in [1, \infty)$  we obtain

$$\|\mathbf{f}\|_p^p = \int_{-1}^1 |x|^p dx = 2 \int_0^1 x^p dx = 2 \frac{x^{p+1}}{p+1} \Big|_0^1 = \frac{2}{p+1} \Rightarrow \|\mathbf{f}\|_p = \left( \frac{2}{p+1} \right)^{\frac{1}{p}}.$$

In the case of  $p = \infty$  we obtain

$$\|\mathbf{f}\|_\infty = \max_{x \in [-1, 1]} |x| \Rightarrow \|\mathbf{f}\|_\infty = 1.$$

Relations between the  $p$ -norms of a given vector  $\mathbf{f}$  analogous to those relations found in Example 8.1.3 are not longer true. The volume of the integration region, the interval  $[a, b]$ , appears in any relation between  $p$ -norms of a fixed vector. We do not address such relations in these notes.  $\triangleleft$

**8.1.1. Not every norm is an inner product norm.** Since every inner product in a vector space determines a norm, the inner product norm, it is natural to ask whether the converse property holds: Does every norm in a vector space determine an inner product? The answer is no. Only those norms satisfying an extra condition, called the parallelogram identity, define an inner product.

**Definition 8.1.6.** A norm  $\|\cdot\|$  in a vector space  $V$  satisfies the **polarization identity** iff for all vectors  $\mathbf{x}, \mathbf{y} \in V$  holds

$$\|\mathbf{x} + \mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{y}\|^2 = 2(\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2).$$

The polarization identity (also referred as parallelogram identity) is a well-known property of the dot product norm, which is geometrically described in Fig. 54 in the case of the vector space  $\mathbb{R}^2$ . It turns out that this property is crucial to determine whether a norm is an inner product norm for some inner product.

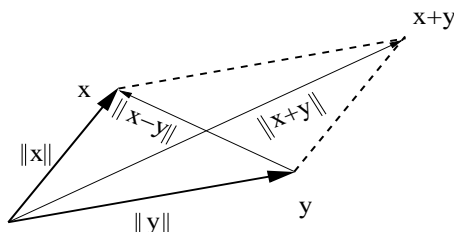


FIGURE 54. The polarization identity says that the sum of the squares of the diagonals in a parallelogram is twice the sum of the squares of the sides.

**Theorem 8.1.7.** Given a normed space  $(V, \|\cdot\|)$ , the norm  $\|\cdot\|$  is an inner product norm iff the norm  $\|\cdot\|$  satisfies the polarization identity.

It is not difficult to see that an inner product norm satisfies the polarization identity; see the first part in the proof below. It is rather involved to show the converse statement. If the norm  $\|\cdot\|$  satisfies the polarization identity and  $V$  is a real vector space, then one shows that the function  $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$  given by

$$\langle \mathbf{x}, \mathbf{y} \rangle = \frac{1}{4} (\|\mathbf{x} + \mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2)$$

is an inner product on  $V$ ; in the case that  $V$  is a complex vector space, then one shows that the function  $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{C}$  given by

$$\langle \mathbf{x}, \mathbf{y} \rangle = \frac{1}{4} (\|\mathbf{x} + \mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2) + \frac{i}{4} (\|\mathbf{x} + i\mathbf{y}\|^2 - \|\mathbf{x} - i\mathbf{y}\|^2)$$

is an inner product on  $V$ .

**Proof of Theorem 8.1.7:**

( $\Rightarrow$ ) Consider the inner product space  $(V, \langle \cdot, \cdot \rangle)$  with inner product norm  $\|\cdot\|$ . For all vectors  $\mathbf{x}, \mathbf{y} \in V$  holds

$$\begin{aligned} \|\mathbf{x} + \mathbf{y}\|^2 &= \langle (\mathbf{x} + \mathbf{y}), (\mathbf{x} + \mathbf{y}) \rangle = \|\mathbf{x}\|^2 + \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{y}, \mathbf{x} \rangle + \|\mathbf{y}\|^2, \\ \|\mathbf{x} - \mathbf{y}\|^2 &= \langle (\mathbf{x} - \mathbf{y}), (\mathbf{x} - \mathbf{y}) \rangle = \|\mathbf{x}\|^2 - \langle \mathbf{x}, \mathbf{y} \rangle - \langle \mathbf{y}, \mathbf{x} \rangle + \|\mathbf{y}\|^2. \end{aligned}$$

Adding both equations above up we obtain that

$$\|\mathbf{x} + \mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{y}\|^2 = 2(\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2).$$

( $\Leftarrow$ ) We only give the proof for real vector spaces. The proof for complex vector spaces is left as an exercise. Consider the normed space  $(V, \|\cdot\|)$  and assume that  $V$  is a real vector space. In this case, introduce the function  $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$  as follows,

$$\langle \mathbf{x}, \mathbf{y} \rangle = \frac{1}{4} (\|\mathbf{x} + \mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2).$$



Notice that this function satisfies  $\langle \mathbf{x}, \mathbf{0} \rangle = \frac{1}{4} (\|\mathbf{x}\|^2 - \|\mathbf{x}\|^2) = 0$  for all  $\mathbf{x} \in V$ . We now show that this function  $\langle \cdot, \cdot \rangle$  is an inner product on  $V$ . It is positive definite, since

$$\langle \mathbf{x}, \mathbf{x} \rangle = \frac{1}{4} \|2\mathbf{x}\|^2 = \|\mathbf{x}\|^2$$

and the norm is positive definite, so we obtain that

$$\langle \mathbf{x}, \mathbf{x} \rangle \geq 0, \quad \text{and} \quad \langle \mathbf{x}, \mathbf{x} \rangle = 0 \quad \Leftrightarrow \quad \mathbf{x} = \mathbf{0}.$$

The function  $\langle \cdot, \cdot \rangle$  is symmetric, since

$$\langle \mathbf{x}, \mathbf{y} \rangle = \frac{1}{4} (\|\mathbf{x} + \mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2) = \frac{1}{4} (\|\mathbf{y} + \mathbf{x}\|^2 - \|\mathbf{y} - \mathbf{x}\|^2) = \langle \mathbf{y}, \mathbf{x} \rangle.$$

The difficult part is to show that the function  $\langle \cdot, \cdot \rangle$  is linear in the second argument. Here is where we need the polarization identity. We start with the following expressions, which are obtained from the polarization identity,

$$\begin{aligned} \|\mathbf{x} + \mathbf{y} + \mathbf{z}\|^2 + \|\mathbf{x} + \mathbf{y} - \mathbf{z}\|^2 &= 2\|\mathbf{x} + \mathbf{y}\|^2 + 2\|\mathbf{z}\|^2, \\ \|\mathbf{x} - \mathbf{y} + \mathbf{z}\|^2 + \|\mathbf{x} - \mathbf{y} - \mathbf{z}\|^2 &= 2\|\mathbf{x} - \mathbf{y}\|^2 + 2\|\mathbf{z}\|^2. \end{aligned}$$

If we subtract the second equation from the first one, and ordering the terms conveniently, we obtain,

$$\left[ \|\mathbf{x} + (\mathbf{y} + \mathbf{z})\|^2 - \|\mathbf{x} - (\mathbf{y} + \mathbf{z})\|^2 \right] + \left[ \|\mathbf{x} + (\mathbf{y} - \mathbf{z})\|^2 - \|\mathbf{x} - (\mathbf{y} - \mathbf{z})\|^2 \right] = 2 \left[ \|\mathbf{x} + \mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2 \right]$$

which can be written in terms of the function  $\langle \cdot, \cdot \rangle$  as follows,

$$\langle \mathbf{x}, (\mathbf{y} + \mathbf{z}) \rangle + \langle \mathbf{x}, (\mathbf{y} - \mathbf{z}) \rangle = 2\langle \mathbf{x}, \mathbf{y} \rangle. \quad (8.6)$$

Several relations come from this equation. For the first relation, take  $\mathbf{y} = \mathbf{z}$ , and recall that  $\langle \mathbf{x}, \mathbf{0} \rangle = 0$ , then we obtain

$$\langle \mathbf{x}, 2\mathbf{y} \rangle = 2\langle \mathbf{x}, \mathbf{y} \rangle. \quad (8.7)$$

The second relation derived from Eq. (8.6) is obtained renaming the vectors  $\mathbf{y} + \mathbf{z} = \mathbf{u}$  and  $\mathbf{y} - \mathbf{z} = \mathbf{v}$ , that is,

$$\langle \mathbf{x}, \mathbf{u} \rangle + \langle \mathbf{x}, \mathbf{v} \rangle = 2\left\langle \mathbf{x}, \frac{(\mathbf{u} + \mathbf{v})}{2} \right\rangle = \langle \mathbf{x}, (\mathbf{u} + \mathbf{v}) \rangle,$$

where the equation on the far right comes from Eq. (8.7). We have shown that for all  $\mathbf{x}, \mathbf{u}, \mathbf{v} \in V$  holds

$$\langle \mathbf{x}, (\mathbf{u} + \mathbf{v}) \rangle = \langle \mathbf{x}, \mathbf{u} \rangle + \langle \mathbf{x}, \mathbf{v} \rangle$$

which is a particular case of the linearity in the second argument property of an inner product. We only need to show that for all  $\mathbf{x}, \mathbf{y} \in V$  and all  $a \in \mathbb{R}$  holds

$$\langle \mathbf{x}, a\mathbf{y} \rangle = a\langle \mathbf{x}, \mathbf{y} \rangle.$$

We have proven the case  $a = 2$  in Eq. (8.7). The case  $a = n \in \mathbb{N}$  is proven by induction: If  $\langle \mathbf{x}, n\mathbf{y} \rangle = n\langle \mathbf{x}, \mathbf{y} \rangle$ , then the same relation holds for  $(n + 1)$ , since

$$\begin{aligned} \langle \mathbf{x}, (n + 1)\mathbf{y} \rangle &= \langle \mathbf{x}, (n\mathbf{y} + \mathbf{y}) \rangle \\ &= \langle \mathbf{x}, n\mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{y} \rangle \\ &= n\langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{y} \rangle \\ &= (n + 1)\langle \mathbf{x}, \mathbf{y} \rangle. \end{aligned}$$

The case  $a = 1/n$  with  $n \in \mathbb{N}$  comes from the relation

$$\langle \mathbf{x}, \mathbf{y} \rangle = \left\langle \mathbf{x}, \frac{n}{n}\mathbf{y} \right\rangle = n\left\langle \mathbf{x}, \frac{1}{n}\mathbf{y} \right\rangle \quad \Rightarrow \quad \left\langle \mathbf{x}, \frac{1}{n}\mathbf{y} \right\rangle = \frac{1}{n}\langle \mathbf{x}, \mathbf{y} \rangle.$$

These two cases show that for any rational number  $a = p/q \in \mathbb{Q}$  holds

$$\left\langle \mathbf{x}, \frac{p}{q} \mathbf{y} \right\rangle = p \left\langle \mathbf{x}, \frac{1}{q} \mathbf{y} \right\rangle = \frac{p}{q} \langle \mathbf{x}, \mathbf{y} \rangle. \quad (8.8)$$

Finally, the same property holds for all  $a \in \mathbb{R}$  by the following continuity argument. Fix arbitrary vectors  $\mathbf{x}, \mathbf{y} \in V$  and define the function  $f: \mathbb{R} \rightarrow \mathbb{R}$  by

$$a \mapsto f(a) = \langle \mathbf{x}, a \mathbf{y} \rangle.$$

We left as an exercise to show that  $f$  is a continuous function. Now, using Eq. (8.8) we know that this function satisfies

$$f(a) = a f(1) \quad \forall a \in \mathbb{Q}.$$

Let  $\{a_k\}_{k=1}^{\infty} \subset \mathbb{Q}$  be a sequence of rational numbers that converges to  $a \in \mathbb{R}$ . Since  $f$  is a continuous function we know that  $\lim_{k \rightarrow \infty} f(a_k) = f(a)$ . Since the sequence is constructed with rational numbers, for every element in this sequence holds

$$f(a_k) = a_k f(1)$$

which implies that

$$\lim_{k \rightarrow \infty} f(a_k) = \left( \lim_{k \rightarrow \infty} a_k \right) f(1) = a f(1).$$

So we have shown that for all  $a \in \mathbb{R}$  holds  $f(a) = a f(1)$ , that is,

$$\langle \mathbf{x}, a \mathbf{y} \rangle = a \langle \mathbf{x}, \mathbf{y} \rangle.$$

Since  $\mathbf{x}, \mathbf{y} \in V$  are fixed but arbitrary, we have shown that the function  $\langle \cdot, \cdot \rangle$  is linear in the second argument. We conclude that this function is an inner product on  $V$ . This establishes the Theorem in the case that  $V$  is a real vector space.  $\square$

**EXAMPLE 8.1.5:** Show that for  $p \neq 2$  the  $p$ -norm on the vector space  $\mathbb{F}^n$  introduced in Definition 8.1.2 does not satisfy the polarization identity.

**SOLUTION:** Consider the vectors  $\mathbf{e}_1$  and  $\mathbf{e}_2$ , the first two columns of the identity matrix  $\mathbf{I}_n$ . It is simple to compute,

$$\|\mathbf{e}_1 + \mathbf{e}_2\|_p^2 = 2^{2/p}, \quad \|\mathbf{e}_1 - \mathbf{e}_2\|_p^2 = 2^{2/p}, \quad \|\mathbf{e}_1\|_p^2 = \|\mathbf{e}_2\|_p^2 = 1$$

therefore,

$$\|\mathbf{e}_1 + \mathbf{e}_2\|_p^2 + \|\mathbf{e}_1 - \mathbf{e}_2\|_p^2 = 2 \frac{2^{2/p}}{p} \quad \text{and} \quad 2(\|\mathbf{e}_1\|_p^2 + \|\mathbf{e}_2\|_p^2) = 4.$$

We conclude that the  $p$ -norm satisfies the polarization identity only in the case  $p = 2$ .  $\triangleleft$

**8.1.2. Equivalent norms.** We have seen that a norm in a vector space determines a notion of distance between vectors given by the norm distance introduced in Definition 6.2.6. The notion of distance is the structure needed to define the convergence of an infinite sequence of vectors. A sequence of vectors  $\{\mathbf{x}_i\}_{i=1}^{\infty}$  in a normed space  $(V, \|\cdot\|)$  converges to a vector  $\mathbf{x} \in V$  iff

$$\lim_{k \rightarrow \infty} \|\mathbf{x} - \mathbf{x}_k\| = 0,$$

that is, for every  $\epsilon > 0$  there exists  $n \in \mathbb{N}$  such that for all  $k > n$  holds that  $\|\mathbf{x} - \mathbf{x}_k\| < \epsilon$ . With the notion of convergence of a sequence it is possible to introduce concepts like the continuous and differentiable functions defined on the vector space. Therefore, the calculus can be extended from  $\mathbb{R}^n$  to any normed vector space.

We have also seen that there is no unique norm in a vector space. This implies that there is no unique notion of norm distance. It is important to know whether two different norms provide the same notion of convergence of a sequence. By this we mean that every sequence that converges (diverges) with respect to one norm distance also converges (diverges) with

respect to the other norm distance. In order to answer this question it is useful the following notion.

**Definition 8.1.8.** *The norms  $\|\cdot\|_a$  and  $\|\cdot\|_b$  defined on a vector space  $V$  are called to be **equivalent** iff there exist real constants  $K \geq k > 0$  such that for all  $\mathbf{x} \in V$  holds*

$$k \|\mathbf{x}\|_b \leq \|\mathbf{x}\|_a \leq K \|\mathbf{x}\|_b.$$

It is simple to see that if two norms defined on a vector space are equivalent, then they have the same notion of convergence. What it is non-trivial to prove is that the converse also holds. Since two norms have the same notion of convergence iff they are equivalent, it is important to know whether a vector space can have non-equivalent norms. The following result addresses the case of finite dimensional vector spaces.

**Theorem 8.1.9.** *If  $V$  is a finite dimensional, then all norms defined on  $V$  are equivalent.*

This result says that there is a unique notion of convergence in any finite dimensional vector space. So, functions that are continuous or differentiable with respect to one norm are also continuous or differentiable with respect to any other norm. This is not the case of infinite dimensional vector spaces. It is possible to find non-equivalent norms on infinite dimensional vector spaces. Therefore, functions defined on such a vector space can be continuous with respect to one norm and discontinuous with respect to the other norm.

**Proof of Theorem 8.1.9:** Let  $\|\cdot\|_a$  and  $\|\cdot\|_b$  be two norms defined on a finite dimensional vector space  $V$ . Let  $\dim V = n$ , and fix a basis  $\mathcal{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  of  $V$ . Then, any vector  $\mathbf{x} \in V$  can be decomposed in terms of the basis vectors as

$$\mathbf{x} = x_1 \mathbf{v}_1 + \dots + x_n \mathbf{v}_n.$$

Since  $\|\cdot\|_a$  is a norm, we can obtain the following bound on  $\|\mathbf{x}\|_a$  for every  $\mathbf{x} \in V$  in terms of  $(\sum_{i=1}^n |x_i|)$  as follows.

$$\begin{aligned} \|\mathbf{x}\|_a &= \|x_1 \mathbf{v}_1 + \dots + x_n \mathbf{v}_n\|_a \\ &\leq |x_1| \|\mathbf{v}_1\|_a + \dots + |x_n| \|\mathbf{v}_n\|_a \\ &\leq (|x_1| + \dots + |x_n|) a_{\max} \quad \Rightarrow \quad \|\mathbf{x}\|_a \leq \left( \sum_{i=1}^n |x_i| \right) a_{\max}. \end{aligned}$$

where  $a_{\max} = \max\{\|\mathbf{v}_1\|_a, \dots, \|\mathbf{v}_n\|_a\}$ . A lower bound can also be found as follows. Introduce the set

$$S = \left\{ [\hat{x}_i] \in \mathbb{F}^n : \sum_{i=1}^n |\hat{x}_i| = 1 \right\} \subset \mathbb{F}^n,$$

and then introduce the function  $f_a : S \rightarrow \mathbb{R}$  as follows

$$f_a([\hat{x}_i]) = \left\| \sum_{i=1}^n \hat{x}_i \mathbf{v}_i \right\|_a.$$

Since function  $f_a$  is a continuous and defined on a closed and bounded set of  $\mathbb{F}^n$ , then  $f_a$  has attains a maximum and a minimum values on  $S$ . We are here interested only in the minimum value. If  $[\hat{y}_i] \in S$  is a point where  $f_a$  takes its minimum value, let let us denote

$$f_a([\hat{y}_i]) = a_{\min}.$$

Since the norm  $\|\cdot\|_a$  is positive, we know that  $a_{\min} > 0$ . The existence of this minimum value of  $f_a$  implies the following bound on  $\|\mathbf{x}\|_a$  for all  $\mathbf{x} \in V$ , namely,

$$\begin{aligned} \|\mathbf{x}\|_a &= \left\| \sum_{i=1}^n x_i \mathbf{v}_i \right\|_a \\ &= \frac{\left( \sum_{j=1}^n |x_j| \right)}{\left( \sum_{j=1}^n |x_j| \right)} \left\| \sum_{i=1}^n x_i \mathbf{v}_i \right\|_a \\ &= \left( \sum_{j=1}^n |x_j| \right) \left\| \sum_{i=1}^n \left[ \frac{x_i}{\left( \sum_{j=1}^n |x_j| \right)} \right] \mathbf{v}_i \right\|_a \\ &= \left( \sum_{j=1}^n |x_j| \right) \left\| \sum_{i=1}^n \hat{x}_i \mathbf{v}_i \right\|_a \\ &= \left( \sum_{j=1}^n |x_j| \right) f_a([\hat{x}_i]) \quad \Rightarrow \quad \|\mathbf{x}\|_a \geq \left( \sum_{j=1}^n |x_j| \right) a_{\min}. \end{aligned}$$

Summarizing, we have found real numbers  $a_{\max} \geq a_{\min} > 0$  such that the following inequality holds for all  $\mathbf{x} \in V$ ,

$$a_{\min} \left( \sum_{j=1}^n |x_j| \right) \leq \|\mathbf{x}\|_a \leq \left( \sum_{j=1}^n |x_j| \right) a_{\max}.$$

Since no special property of norm  $\|\cdot\|_a$  has been used, the same type of inequality holds for norm  $\|\cdot\|_b$ , that is, there exist real numbers  $b_{\max} \geq b_{\min} > 0$  such that the following inequality holds for all  $\mathbf{x} \in B$ ,

$$b_{\min} \left( \sum_{j=1}^n |x_j| \right) \leq \|\mathbf{x}\|_b \leq \left( \sum_{j=1}^n |x_j| \right) b_{\max}.$$

These inequalities imply that norms  $\|\cdot\|_a$  and  $\|\cdot\|_b$  are equivalent, since

$$\|\mathbf{x}\|_b \frac{a_{\min}}{b_{\max}} \leq \frac{b_{\max}}{b_{\max}} a_{\min} \left( \sum_{j=1}^n |x_j| \right) \leq \|\mathbf{x}\|_a \leq \left( \sum_{j=1}^n |x_j| \right) a_{\max} \frac{b_{\min}}{b_{\min}} \leq \frac{a_{\max}}{b_{\min}} \|\mathbf{x}\|_b.$$

(Start reading the inequality from the center, at  $\|\mathbf{x}\|_a$ , and first see the inequalities to the right; then go to the center again and read the inequalities to the left.) Denoting  $k = a_{\min}/b_{\max}$  and  $K = a_{\max}/b_{\min}$ , we have obtained that

$$k \|\mathbf{x}\|_b \leq \|\mathbf{x}\|_a \leq K \|\mathbf{x}\|_b \quad \forall \mathbf{x} \in V.$$

This establishes the Theorem. □

## 8.1.3. Exercises.

8.1.1.- Consider the vector space  $\mathbb{C}^4$  and for  $p = 1, 2, \infty$  find the  $p$ -norms of the vectors

$$\mathbf{x} = \begin{bmatrix} 2 \\ 1 \\ -4 \\ -2 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 1+i \\ 1-i \\ 1 \\ 4i \end{bmatrix}.$$

8.1.2.- Determine which of the following functions  $\|\cdot\| : \mathbb{R}^2 \rightarrow \mathbb{R}$  defines a norm on  $\mathbb{R}^2$ . We denote  $\mathbf{x} = [x_i]$  the components of  $\mathbf{x}$  in a standard basis of  $\mathbb{R}^2$ . Justify your answers.

- (a)  $\|\mathbf{x}\| = |x_1|$ ;
- (b)  $\|\mathbf{x}\| = |x_1 + x_2|$ ;
- (c)  $\|\mathbf{x}\| = |x_1|^2 + |x_2|^2$ ;
- (d)  $\|\mathbf{x}\| = 2|x_1| + 3|x_2|$ .

8.1.3.- True or false? Justify your answer:

If  $\|\cdot\|_a$  and  $\|\cdot\|_b$  are two norms on a vector space  $V$ , then  $\|\cdot\|$  defined as

$$\|\mathbf{x}\| = \|\mathbf{x}\|_a + \|\mathbf{x}\|_b \quad \forall \mathbf{x} \in V$$

is also a norm on  $V$ .

8.1.4.- Consider the space  $\mathbb{P}_2([0, 1])$  with the  $p$ -norm

$$\|\mathbf{q}\|_p = \left( \int_0^1 |\mathbf{q}(x)|^p dx \right)^{\frac{1}{p}},$$

for  $p \in [1, \infty)$ . Find the  $p$ -norm of the vector  $\mathbf{q}(x) = -3x^2$ . Also find the supremum norm of  $\mathbf{q}$ , defined as

$$\|\mathbf{q}\|_\infty = \max_{x \in [0, 1]} |\mathbf{q}(x)|.$$

## 8.2. OPERATOR NORMS

We have seen that the space of all linear transformations  $L(V, W)$  between the vector spaces  $V$  and  $W$  is itself a vector space. As in any vector space, it is possible to define norm functions  $\| \cdot \| : L(V, W) \rightarrow \mathbb{R}$  on the vector space  $L(V, W)$ . For example, We will see that in the case where  $V = \mathbb{F}^n$  and  $W = \mathbb{F}^m$ , one of such norms is the Frobenius norm, which is the inner product norm corresponding to the Frobenius inner product introduced in Section 6.2. More precisely, the **Frobenius norm** is given by

$$\|A\|_F = \sqrt{\langle A, A \rangle_F} \quad \forall A \in \mathbb{F}^{m,n}.$$

It is simple to see that

$$\|A\|_F = \sqrt{\operatorname{tr}(A^*A)} = \left( \sum_{i=1}^m \sum_{j=1}^n |A_{ij}|^2 \right)^{\frac{1}{2}}.$$

However, in the case where  $V$  and  $W$  are normed spaces with norms  $\| \cdot \|_v$  and  $\| \cdot \|_w$ , there exists a particular norm on  $L(V, W)$  which is induced from the norms on  $V$  and  $W$ . This induced norm can be defined on elements  $L(V, W)$  by recalling that the element  $T \in L(V, W)$  as a function  $T : V \rightarrow W$ .

**Definition 8.2.1.** Let  $(V, \| \cdot \|_v)$  and  $(W, \| \cdot \|_w)$  be finite dimensional normed spaces. The **induced norm** on the space  $L(V, W)$  of all linear transformations  $T : V \rightarrow W$  is a function  $\| \cdot \| : L(V, W) \rightarrow \mathbb{R}$  given by

$$\|T\| = \max_{\|x\|_v=1} \|T(x)\|_w.$$

In the particular case that  $V = W$  and  $\| \cdot \|_v = \| \cdot \|_w = \| \cdot \|$ , the induced norm on  $L(V)$  is called the **operator norm**, and is given by

$$\|T\| = \max_{\|x\|=1} \|T(x)\|.$$

The definition above says that given arbitrary norms on the vector spaces  $V$  and  $W$ , they induce a norm on the vector space  $L(V, W)$  of linear transformations  $T : V \rightarrow W$ . A particular case of this definition is when  $V = \mathbb{F}^n$  and  $W = \mathbb{F}^m$ , we fix standard bases on  $V$  and  $W$ , and we introduce  $p$ -norms on these spaces. So, the normed spaces  $(\mathbb{F}^n, \| \cdot \|_p)$  and  $(\mathbb{F}^m, \| \cdot \|_q)$  induce a norm on the vector space  $\mathbb{F}^{m,n}$  as follows.

**Definition 8.2.2.** Consider the normed spaces  $(\mathbb{F}^n, \| \cdot \|_p)$  and  $(\mathbb{F}^m, \| \cdot \|_q)$ , with  $p, q \in [1, \infty]$ . The **induced  $(p, q)$ -norm** on the space  $L(\mathbb{F}^n, \mathbb{F}^m)$  is the function  $\| \cdot \|_{p,q} : L(\mathbb{F}^n, \mathbb{F}^m) \rightarrow \mathbb{R}$  given by

$$\|T\|_{p,q} = \max_{\|x\|_p=1} \|T(x)\|_q \quad \forall T \in L(\mathbb{F}^n, \mathbb{F}^m).$$

In the particular case of  $p = q$  we denote  $\| \cdot \|_{p,p} = \| \cdot \|_p$ . In the case  $p = q$  and  $n = m$  the induced norm on  $L(\mathbb{F}^n)$  is called the  **$p$ -operator norm** and is given by

$$\|T\|_p = \max_{\|x\|_p=1} \|T(x)\|_p \quad \forall T \in L(\mathbb{F}^n).$$

In the case  $p = q$  above we use the same notation  $\| \cdot \|_p$  for the  $p$ -norm on  $\mathbb{F}^n$ , the  $(p, p)$ -norm on  $L(\mathbb{F}^n, \mathbb{F}^m)$  and the  $p$ -operator norm in  $L(\mathbb{F}^n)$ . The context should help to decide which norm we use on every particular situation.

**EXAMPLE 8.2.1:** Consider the particular case  $V = W = \mathbb{R}^2$  with standard ordered bases and the  $p = 2$  norms in both,  $V$  and  $W$ . In this case, the space  $L(\mathbb{R}^2)$  can be identified with

the space  $\mathbb{R}^{2,2}$  of all  $2 \times 2$  matrices. The induced norm on  $\mathbb{R}^{2,2}$ , denoted as  $\|\cdot\|_2$  and called the operator norm, is the following:

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2, \quad \forall A \in \mathbb{R}^{2,2}.$$

The meaning of this norm is deeply related with the interpretation of the matrix  $A$  as a linear operator  $A : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ . Suppose that the action of the operator  $A$  on the unit circle ( $B_2$  in the notation of Example 8.1.2) is given in Fig. 55. Then the value of the operator norm  $\|A\|_2$  is the size measured in the 2-norm of the maximum deformation by  $A$  of the unit circle.

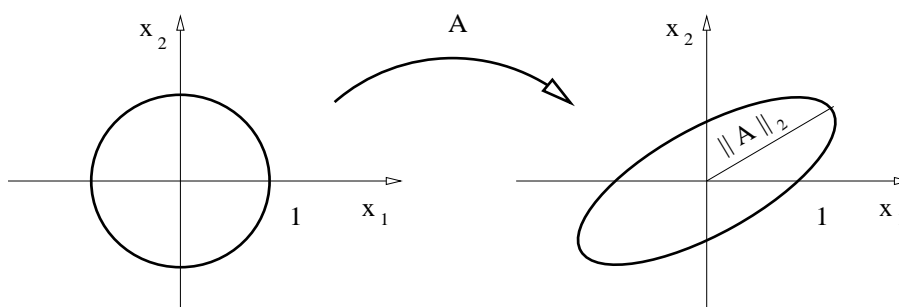


FIGURE 55. Geometrical meaning of the operator norm on  $\mathbb{R}^{2,2}$  induced by the 2-norm on  $\mathbb{R}^2$ .

◁

**EXAMPLE 8.2.2:** Consider the normed space  $(\mathbb{R}^2, \|\cdot\|_2)$ , and let  $A : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be the matrix

$$A = \begin{bmatrix} A_{11} & 0 \\ 0 & A_{22} \end{bmatrix} \quad \text{with} \quad |A_{11}| \neq |A_{22}|.$$

Find the 2-operator norm induced on  $A$ .

**SOLUTION:** Since the norm on  $\mathbb{R}^2$  is the 2-norm, the induced norm on  $A$  is given by

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2.$$

We need to find the maximum of  $\|Ax\|_2$  among all  $x$  subject to the constraint  $\|x\|_2 = 1$ . So, this is a constrained maximization problem, that is, a maxima-minima problem where the variable  $x$  is restricted by a constraint equation. In general, this type of problems can be solved using the Lagrange multipliers method. Since this example is simple enough, we solve it in a simpler way. We first solve the constraint equation on  $x$ , and we then find the maxima of  $\|Ax\|_2$  among these solutions only. The general solution of the equation

$$\|x\|_2 = 1 \quad \Leftrightarrow \quad (x_1)^2 + x_2^2 = 1$$

is given by

$$x(\theta) = \begin{bmatrix} \cos(\theta) \\ \sin(\theta) \end{bmatrix} \quad \text{with} \quad \theta \in [0, 2\pi).$$

Introduce this general solution into  $\|Ax\|_2$  we obtain

$$\|Ax\|_2 = \sqrt{(A_{11})^2 \cos^2(\theta) + (A_{22})^2 \sin^2(\theta)}.$$

Since the maximum in  $\theta$  of the function  $\|\mathbf{Ax}\|_2$  is the same as the maximum of  $\|\mathbf{Ax}\|_2^2$ , we need to find the maximum on  $\theta$  of the function  $f(\theta) = \|\mathbf{Ax}(\theta)\|_2^2$ , that is,

$$f(\theta) = (A_{11})^2 \cos^2(\theta) + (A_{22})^2 \sin^2(\theta).$$

The solution is simple, find  $\theta$  solutions of

$$f'(\theta) = \frac{df}{d\theta}(\theta) = 0 \quad \Rightarrow \quad 2[-(A_{11})^2 + (A_{22})^2] \sin(\theta) \cos(\theta) = 0.$$

Since we assumed that  $|A_{11}| \neq |A_{22}|$ , then

$$\begin{aligned} \sin(\theta) = 0 &\Rightarrow \theta_1 = 0, \quad \theta_2 = \pi, \\ \cos(\theta) = 0 &\Rightarrow \theta_3 = \frac{\pi}{2}, \quad \theta_4 = \frac{3\pi}{2}. \end{aligned}$$

We obtained four solutions for  $\theta$ . We evaluate  $f$  at these solutions,

$$f(0) = f(\pi) = (A_{11})^2, \quad f\left(\frac{\pi}{2}\right) = f\left(\frac{3\pi}{2}\right) = (A_{22})^2.$$

Recalling that  $\|\mathbf{Ax}(\theta)\|_2 = \sqrt{f(\theta)}$ , we obtain

$$\|\mathbf{A}\|_2 = \max\{|A_{11}|, |A_{22}|\}.$$

◁

It was mentioned in Example 8.2.2 above that finding the operator norm requires solving a constrained maximization problem. In the case that the operator norm is induced from a dot product norm, the constrained maximization problem can be solved in an explicit form.

**Proposition 8.2.3.** *Consider the vector spaces  $\mathbb{R}^n$  and  $\mathbb{R}^m$  with inner product given by the dot products, inner product norms  $\|\cdot\|_2$ , and the vector space  $L(\mathbb{R}^n, \mathbb{R}^m)$  with the induced 2-norm.*

$$\|\mathbf{A}\|_2 = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{Ax}\|_2 \quad \forall \mathbf{A} \in L(\mathbb{R}^n, \mathbb{R}^m).$$

Introduce the scalars  $\lambda_i \in \mathbb{R}$ , with  $i = 1, \dots, k \leq n$ , as all the roots of the polynomial

$$p(\lambda) = \det(\mathbf{A}^T \mathbf{A} - \lambda \mathbf{I}_n).$$

Then, all scalars  $\lambda_i$  are non-negative real numbers and the induced 2-norm of the transformation  $\mathbf{A}$  is given by

$$\|\mathbf{A}\|_2 = \max\{\lambda_1, \dots, \lambda_k\}.$$

**Proof of Proposition 8.2.3:** Introduce the functions  $f : \mathbb{R}^m \rightarrow \mathbb{R}$  and  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  as follows,

$$f(\mathbf{x}) = \|\mathbf{Ax}\|_2^2 = (\mathbf{Ax}) \cdot (\mathbf{Ax}) = \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x}, \quad g(\mathbf{x}) = \|\mathbf{x}\|_2^2 = \mathbf{x} \cdot \mathbf{x} = \mathbf{x}^T \mathbf{x}.$$

To find the induced norm of  $\mathbf{T}$  is then equivalent to solve the constrained maximization problem: Find the maximum of  $f(\mathbf{x})$  for  $\mathbf{x} \in \mathbb{R}^n$  subject to the constraint  $g(\mathbf{x}) = 1$ . The vectors  $\mathbf{x}$  that provide solutions to the constrained maximization problem must be solutions of the Euler-Lagrange equations

$$\nabla f = \lambda \nabla g \tag{8.9}$$

where  $\lambda \in \mathbb{R}$ , and we introduced the gradient row vectors

$$\nabla f = \left[ \frac{\partial f}{\partial x_1} \quad \dots \quad \frac{\partial f}{\partial x_n} \right], \quad \nabla g = \left[ \frac{\partial g}{\partial x_1} \quad \dots \quad \frac{\partial g}{\partial x_n} \right].$$

In order to understand why the solution  $\mathbf{x}$  must satisfy the Euler-Lagrange equations above we need to recall two properties of the gradient vector. First, the gradient of a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a vector that determines the direction on  $\mathbb{R}^n$  where  $f$  has the maximum increase. Second, which is deeply related to the first property, the gradient vector of a



function is perpendicular to the surfaces where the function has a constant value. The surfaces of constant value of a function are called level surfaces of the function. Therefore, the function  $f$  has a maximum or minimum value at  $\mathbf{x}$  on the constraint level surface  $g = 1$  if  $\nabla f$  is perpendicular to the level surface  $g = 1$  at that  $\mathbf{x}$ . (Proof: Suppose that at a particular  $\mathbf{x}$  on the constraint surface  $g = 1$  the projection of  $\nabla f$  onto the constraint surface is nonzero; then the values of  $f$  increase along that direction on the constraint surface; this means that  $f$  does not attain a maximum value at that  $\mathbf{x}$  on  $g = 1$ .) We conclude that both gradients  $\nabla f$  and  $\nabla g$  are parallel, which is precisely what Eq. (8.9) says. In our particular problem we obtain for  $f$  and  $g$  the following:

$$\begin{aligned} f(\mathbf{x}) = \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} &\Rightarrow \nabla f(\mathbf{x}) = 2\mathbf{x}^T \mathbf{A}^T \mathbf{A}, \\ g(\mathbf{x}) = \mathbf{x}^T \mathbf{x} &\Rightarrow \nabla g(\mathbf{x}) = 2\mathbf{x}^T. \end{aligned}$$

We must look for  $\mathbf{x} \neq \mathbf{0}$  solution of the equation

$$\mathbf{x}^T \mathbf{A}^T \mathbf{A} = \lambda \mathbf{x}^T \Leftrightarrow \mathbf{A}^T \mathbf{A} \mathbf{x} = \lambda \mathbf{x},$$

where the condition  $\mathbf{x} \neq \mathbf{0}$  comes from  $\|\mathbf{x}\|_2 = 1$ . Therefore,  $\lambda$  must not be any scalar but the precise scalar or scalars such that the matrix  $(\mathbf{A}^T \mathbf{A} - \lambda \mathbf{I}_n)$  is not invertible. An equivalent condition is that

$$p(\lambda) = \det(\mathbf{A}^T \mathbf{A} - \lambda \mathbf{I}_n) = 0.$$

The function  $p$  is a polynomial of degree  $n$  in  $\lambda$  so it has at most  $n$  real roots. Let us denote these roots by  $\lambda_1, \dots, \lambda_k$ , with  $1 \leq k \leq n$ . For each of these values  $\lambda_i$  the matrix  $\mathbf{A}^T \mathbf{A} - \lambda_i \mathbf{I}_n$  is not invertible, so  $N(\mathbf{A}^T \mathbf{A} - \lambda_i \mathbf{I}_n)$  is non-trivial. Let  $\mathbf{x}_i$  be any element in  $N(\mathbf{A}^T \mathbf{A} - \lambda_i \mathbf{I}_n)$ , that is,

$$\mathbf{A}^T \mathbf{A} \mathbf{x}_i = \lambda_i \mathbf{x}_i, \quad i = 1, \dots, k.$$

At this point it is not difficult to see that  $\lambda_i \geq 0$  for  $i = 1, \dots, k$ . Indeed, multiply the equation above by  $\mathbf{x}_i^T$ , that is,

$$\mathbf{x}_i^T \mathbf{A}^T \mathbf{A} \mathbf{x}_i = \lambda_i \mathbf{x}_i^T \mathbf{x}_i.$$

Since both  $\mathbf{x}_i^T \mathbf{A}^T \mathbf{A} \mathbf{x}_i$  and  $\mathbf{x}_i^T \mathbf{x}_i$  are non-negative numbers, so is  $\lambda_i$ . Returning to  $\mathbf{x}_i$ , only these vectors are the candidates to find a solution to our maximization problem, and  $f$  has the values

$$f(\mathbf{x}_i) = \mathbf{x}_i^T \mathbf{A}^T \mathbf{A} \mathbf{x}_i = \lambda_i \mathbf{x}_i^T \mathbf{x}_i = \lambda_i \Rightarrow f(\mathbf{x}_i) = \lambda_i, \quad i = 1, \dots, k.$$

The induced 2-norm of  $\mathbf{A}$  is the maximum of these scalar  $\lambda_i$ . This establishes the Proposition.  $\square$

**EXAMPLE 8.2.3:** Consider the vector spaces  $\mathbb{R}^3$  and  $\mathbb{R}^2$  with the dot product. Find the induced 2-norm of the  $2 \times 3$  matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \end{bmatrix}.$$

**SOLUTION:** Following the Proposition 8.2.3 the value of the induced norm  $\|\mathbf{A}\|_2$  is the maximum of the  $\lambda_i$  roots of the polynomial

$$p(\lambda) = \det(\mathbf{A}^T \mathbf{A} - \lambda \mathbf{I}_3) = 0.$$

We start computing

$$\mathbf{A}^T \mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \\ 1 & -1 & 2 \end{bmatrix}.$$

The next step is to find the polynomial

$$p(\lambda) = \begin{vmatrix} (1-\lambda) & 0 & 1 \\ 0 & (1-\lambda) & -1 \\ 1 & -1 & (2-\lambda) \end{vmatrix} = (1-\lambda)[(1-\lambda)(2-\lambda) - 1] + (1-\lambda),$$

therefore, we obtain  $p(\lambda) = -\lambda(\lambda - 1)(\lambda - 3)$ . We have three roots

$$\lambda_1 = 0, \quad \lambda_2 = 1, \quad \lambda_3 = 3.$$

We then conclude that the induced norm of  $\mathbf{A}$  is given by

$$\|\mathbf{A}\|_2 = 3.$$

◁

In the case that the norm in a normed space is not an inner product norm the induced norm on linear operators is not simple to evaluate. The constrained maximization problem is in general complicated to solve. Two particular cases can be solved explicitly though, when the operator norm is induced from the  $p$ -norms with  $p = 1$  and  $p = \infty$ .

**Proposition 8.2.4.** *Consider the normed spaces  $(\mathbb{R}^n, \|\cdot\|_p)$  and  $(\mathbb{R}^m, \|\cdot\|_p)$  and the vector space  $L(\mathbb{R}^n, \mathbb{R}^m)$  with the induced  $p$ -norm.*

$$\|\mathbf{A}\|_p = \max_{\|\mathbf{x}\|_p=1} \|\mathbf{A}\mathbf{x}\|_p \quad \forall \mathbf{A} \in L(\mathbb{R}^n, \mathbb{R}^m).$$

If  $p = 1$  or  $p = \infty$ , then the following formulas hold, respectively,

$$\|\mathbf{A}\|_1 = \max_{j \in \{1, \dots, n\}} \sum_{i=1}^m |A_{ij}|, \quad \|\mathbf{A}\|_\infty = \max_{i \in \{1, \dots, m\}} \sum_{j=1}^n |A_{ij}|.$$

**Proof of proposition 8.2.4:** From the definition of the induced  $p$ -norm for  $p = 1$ ,

$$\|\mathbf{A}\|_1 = \max_{\|\mathbf{x}\|_1=1} \|\mathbf{A}\mathbf{x}\|_1.$$

From the  $p$ -norm on  $\mathbb{R}^m$  we know that  $\|\mathbf{A}\mathbf{x}\|_1 = \sum_{i=1}^m \left| \sum_{j=1}^n A_{ij}x_j \right|$ , therefore,

$$\|\mathbf{A}\mathbf{x}\|_1 \leq \sum_{i=1}^m \sum_{j=1}^n |A_{ij}| |x_j| = \sum_{j=1}^n |x_j| \sum_{i=1}^m |A_{ij}| \leq \left( \sum_{j=1}^n |x_j| \right) \max_{j \in \{1, \dots, n\}} \left( \sum_{i=1}^m |A_{ij}| \right);$$

and introducing the condition  $\|\mathbf{x}\|_1 = 1$  we obtain the inequality

$$\|\mathbf{A}\mathbf{x}\|_1 \leq \max_{j \in \{1, \dots, n\}} \sum_{i=1}^m |A_{ij}|.$$

Recalling the column vector notation  $\mathbf{A} = [\mathbf{A}_{\cdot 1}, \dots, \mathbf{A}_{\cdot n}]$ , we notice that  $\|\mathbf{A}_{\cdot j}\|_1 = \sum_{i=1}^m |A_{ij}|$ , so the inequality above can be expressed as

$$\|\mathbf{A}\mathbf{x}\|_1 \leq \max_{j \in \{1, \dots, n\}} \|\mathbf{A}_{\cdot j}\|_1.$$

Since the left hand side is independent of  $\mathbf{x}$ , the inequality also holds for the maximum in  $\|\mathbf{x}\|_1 = 1$ , that is,

$$\|\mathbf{A}\|_1 \leq \max_{j \in \{1, \dots, n\}} \|\mathbf{A}_{\cdot j}\|_1.$$

It is now clear that the equality has to be achieved, since the  $\|\mathbf{A}\mathbf{x}\|_1$  in the case  $\mathbf{x} = \mathbf{e}_j$ , with  $\mathbf{e}_j$  a standard basis vector, takes the value  $\|\mathbf{A}\mathbf{e}_j\|_1 = \|\mathbf{A}_{\cdot j}\|_1$ . Therefore,

$$\|\mathbf{A}\|_1 = \max_{j \in \{1, \dots, n\}} \|\mathbf{A}_{\cdot j}\|_1.$$

The second part of Proposition 8.2.4 can be proven as follows. From the definition of the induced  $p$ -norm for  $p = \infty$  we know that

$$\|\mathbf{A}\|_{\infty} = \max_{\|\mathbf{x}\|_{\infty}=1} \|\mathbf{Ax}\|_{\infty}.$$

From the  $p$ -norm on  $\mathbb{R}^m$  we see

$$\|\mathbf{Ax}\|_{\infty} = \max_{i \in \{1, \dots, m\}} \left| \sum_{j=1}^n A_{ij} x_j \right| \leq \max_{i \in \{1, \dots, m\}} \sum_{j=1}^n |A_{ij}| |x_j| \leq \max_{i \in \{1, \dots, m\}} \sum_{j=1}^n |A_{ij}|.$$

Since the far right hand side does not depend on  $\mathbf{x}$ , the inequality must hold for the maximum in  $\mathbf{x}$  with  $\|\mathbf{x}\|_{\infty} = 1$ , so we conclude that

$$\|\mathbf{A}\|_{\infty} \leq \max_{i \in \{1, \dots, m\}} \sum_{j=1}^n |A_{ij}|.$$

As in the previous  $p = 1$  case, the value on the right hand side above is achieved for  $\|\mathbf{Ax}\|_{\infty}$  in the case of  $\mathbf{x}$  with components  $\pm 1$  depending on the sign of  $A_{ij}$ . More precisely, choose  $\mathbf{x}$  as follows,

$$x_j = \begin{cases} 1 & \text{if } A_{ij} \geq 0, \\ -1 & \text{if } A_{ij} < 0, \end{cases} \Rightarrow \sum_{j=1}^n A_{ij} x_j = \sum_{j=1}^n |A_{ij}|, \quad i = 1, \dots, m.$$

Therefore, for that  $\mathbf{x}$  we have  $\|\mathbf{x}\|_{\infty} = 1$  and  $\|\mathbf{Ax}\|_{\infty} = \max_{i \in \{1, \dots, m\}} \sum_{j=1}^n |A_{ij}|$ . So, we conclude

$$\|\mathbf{A}\|_{\infty} = \max_{i \in \{1, \dots, m\}} \sum_{j=1}^n |A_{ij}|.$$

This establishes the Proposition □

**EXAMPLE 8.2.4:** Find the induced  $p$ -norm, where  $p = 1, \infty$ , for the  $2 \times 3$  matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 4 & 1 \\ 2 & 1 & -1 \end{bmatrix}.$$

**SOLUTION:** Proposition 8.2.4 says that  $\|\mathbf{A}\|_1$  is the largest absolute value sum of components among columns of  $\mathbf{A}$ , while  $\|\mathbf{A}\|_{\infty}$  is the largest absolute value sum among rows of  $\mathbf{A}$ . In the first case we have:

$$\|\mathbf{A}\|_1 = \max_{j=1,2,3} \sum_{i=1}^2 |A_{ij}|;$$

since

$$\sum_{i=1}^2 |A_{i1}| = 3, \quad \sum_{i=1}^2 |A_{i2}| = 5, \quad \sum_{i=1}^2 |A_{i3}| = 2,$$

therefore,  $\|\mathbf{A}\|_1 = 5$ . In the second case we have:

$$\|\mathbf{A}\|_{\infty} = \max_{i=1,2} \sum_{j=1}^3 |A_{ij}|;$$

since

$$\sum_{j=1}^3 |A_{1j}| = 6, \quad \sum_{j=1}^3 |A_{2j}| = 4,$$

therefore,  $\|\mathbf{A}\|_{\infty} = 6$ . ◁

## 8.2.1. Exercises.

8.2.1.- Evaluate the induced  $p$ -norm, where  $p = 1, 2, \infty$ , for the matrices

$$A = \begin{bmatrix} 1 & -2 \\ -1 & 2 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}.$$

8.2.2.- In the normed space  $(\mathbb{R}^2, \|\cdot\|_2)$ , find the induced norm of  $A : \mathbb{R}^2 \rightarrow \mathbb{R}^2$

$$A = \frac{1}{\sqrt{3}} \begin{bmatrix} 3 & -1 \\ 0 & \sqrt{8} \end{bmatrix}.$$

8.2.3.- Consider the space  $(\mathbb{F}^n, \|\cdot\|_p)$  and the space  $\mathbb{F}^{n,n}$  with the induced norm  $\|\cdot\|_p$ . Prove that for all  $A, B \in \mathbb{F}^{n,n}$  and all  $x \in \mathbb{F}^n$  holds

- (a)  $\|Ax\|_p \leq \|A\|_p \|x\|_p$ ;
- (b)  $\|AB\|_p \leq \|A\|_p \|B\|_p$ .

## 8.3. CONDITION NUMBERS

In Sect. 1.5 we discussed several types of approximations errors that appear when solving an  $m \times n$  linear system in a floating-point number set using rounding. In this Section we discuss a particular type of square linear systems with unique solutions that are greatly affected by small changes in the coefficients of their augmented matrices. We will call such systems ill-conditioned. When an ill-conditioned  $n \times n$  linear system is solved on a floating-point number set using rounding, a small rounding error in the coefficients of the system may produce an approximate solution that differs significantly from the exact solution.

**Definition 8.3.1.** *An  $n \times n$  linear system having a unique solution is **ill-conditioned** iff a 1% perturbation in a coefficient of its augmented matrix produces a perturbed linear system still having a unique solution which differs from the unperturbed solution in a 100% or more.*

We remark that the choice of the values 1% and 100% in the above definition is not a standard choice in the literature. While these values may change on different books, the idea behind the definition of an ill-conditioned system is still the same, that a small change in a coefficient of the linear system produces a big change in its solution.

We also remark that our definition of an ill-conditioned system applies only to square linear system having a unique solution, and such that the perturbed linear system also has a unique solution. The concept of ill-conditioned system can be generalized to other linear systems, but we do not study those cases here.

The following example gives some insight to understand what causes a  $2 \times 2$  system to be ill-conditioned.

**EXAMPLE 8.3.1:** It is not difficult to understand when a  $2 \times 2$  linear system is ill-conditioned. The solution of a  $2 \times 2$  linear system can be thought as the intersection of two lines on the plane, where each line represents the solution of each equation of the system. A  $2 \times 2$  linear system is ill-conditioned when these intersecting lines are almost parallel. Then, a small change in the coefficients of the system produces a small change in the lines representing the solution of each equation. Since the lines are almost parallel, this small change on the lines may produce a large change of the intersection point. This situation is sketched on Fig. 56.

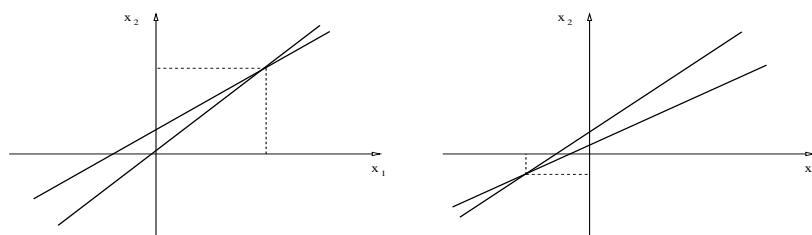


FIGURE 56. The intersection point of almost parallel lines represents a solution of an ill-conditioned  $2 \times 2$  linear system. A small perturbation in the coefficients of the system produces a small change in the lines, which in turn produces a large change in the solution of the system.

◀

**EXAMPLE 8.3.2:** Show that the following  $2 \times 2$  linear system is ill-conditioned:

$$0.835 x_1 + 0.667 x_2 = 0.168,$$

$$0.333 x_1 + 0.266 x_2 = 0.067.$$

**SOLUTION:** We first note that the solution to the linear system above is given by

$$x_1 = 1, \quad x_2 = -1.$$

In order to show that the system above is ill-conditioned we only need to find a coefficient in the system such that a small change in that coefficient produces a large change in the solution. Consider the following change on the second source coefficient:

$$0.067 \rightarrow 0.066.$$

One can check that the solution to the new system

$$0.835 x_1 + 0.667 x_2 = 0.168,$$

$$0.333 x_1 + 0.266 x_2 = 0.066,$$

is given by

$$x_1 = -666, \quad x_2 = 834.$$

We then conclude that the system is ill-conditioned.  $\triangleleft$

Summarizing, solving a linear system in a floating-point number set using rounding introduces approximation errors in the coefficients of the system. The modified Gauss-Jordan method on the floating point number set also introduces approximation errors in the solution of the system. Both type of approximation errors can be controlled, that is, be kept small, choosing a particular scheme of Gauss operations; for example partial pivoting and complete pivoting. However, if the original linear system we solve is ill-conditioned, then even very small approximation errors in the system coefficients and in the Gauss operations may result in a huge error in the solutions. Therefore, it is important to prevent solving ill-conditioned systems when approximation errors are unavoidable. How handle such situations is one important research area in numerical analysis.

**8.3.1. Exercises.**

**8.3.1.-** Consider the ill-conditioned system from Example 8.3.2,

$$0.835x_1 + 0.667x_2 = 0.168,$$

$$0.333x_1 + 0.266x_2 = 0.067.$$

- (a) Solve this system in  $\mathbb{F}_{5,10,6}$  without scaling, partial or complete pivoting.
- (b) Solve this system in  $\mathbb{F}_{6,10,6}$  without scaling, partial or complete pivoting.
- (c) Compare the results found in parts (a) and (b) with the result in  $\mathbb{R}$ .

**8.3.2.-** Perturb the ill-conditioned system given in Exercise **8.3.1** as follows,

$$0.835x_1 + 0.667x_2 = 0.1669995,$$

$$0.333x_1 + 0.266x_2 = 0.066601.$$

Find the solution of this system in  $\mathbb{R}$  and compare it with the solution in Exercise **8.3.1**.

**8.3.3.-** Find the solution in  $\mathbb{R}$  of the following system

$$8x_1 + 5x_2 + 2x_3 = 15,$$

$$21x_1 + 19x_2 + 16x_3 = 56,$$

$$39x_1 + 48x_2 + 53x_3 = 140.$$

Then, change 15 to 14 in the first equation and solve it again in  $\mathbb{R}$ . Is this system ill-conditioned?

## CHAPTER 9. SPECTRAL DECOMPOSITION

In Sect.5.4 we discussed the matrix representation of linear transformations defined on finite dimensional vector spaces. We saw that this representation is basis-dependent. The matrix of a linear transformation can be complicated in one basis and simple in another basis. In the case of linear operators sometimes there exists a special basis where the operator matrix is the simplest possible: A diagonal matrix. In this Chapter we study those linear operators on finite dimensional vector spaces having a diagonal matrix representation, called normal operators. We start introducing the eigenvalues and eigenvectors of a linear operator. Later on we present the main result, the Spectral Theorem for normal operators. We use this result to define functions of operators. We finally mention how to apply these ideas to find solutions to a linear system of ordinary differential equations.

## 9.1. EIGENVALUES AND EIGENVECTORS

**9.1.1. Main definitions.** Sometimes a linear operator has the following property: The image under the operator of a particular line in the vector space is again the same line. In that case we call the line special, eigen in German. Any non-zero vector in that line is also special and is called an eigenvector.

**Definition 9.1.1.** Let  $V$  be a finite dimensional vector space over the field  $\mathbb{F}$ . The scalar  $\lambda \in \mathbb{F}$  and the non-zero vector  $\mathbf{x} \in V$  are called **eigenvalue** and **eigenvector** of the linear operator  $\mathbf{T} \in L(V)$  iff holds

$$\mathbf{T}(\mathbf{x}) = \lambda \mathbf{x}. \quad (9.1)$$

The set  $\sigma_{\mathbf{T}} \subset \mathbb{F}$  of all eigenvalues of the operator  $\mathbf{T}$  is called the **spectrum** of  $\mathbf{T}$ . The subspace  $E_{\lambda} = N(\mathbf{T} - \lambda \mathbf{I}) \subset V$ , the null space of the operator  $(\mathbf{T} - \lambda \mathbf{I})$ , is called the **eigenspace** of  $\mathbf{T}$  corresponding to  $\lambda$ .

An eigenvector of a linear operator  $\mathbf{T}: V \rightarrow V$  is a vector that remain invariant except by scaling under the action of  $\mathbf{T}$ . The change in scaling of the eigenvector  $\mathbf{x}$  under  $\mathbf{T}$  determines the eigenvalue  $\lambda$ . Since the operator  $\mathbf{T}$  is linear, given an eigenvector  $\mathbf{x}$  and any non-zero scalar  $a \in \mathbb{F}$ , the vector  $a\mathbf{x}$  is also an eigenvector. (Proof:  $\mathbf{T}(a\mathbf{x}) = a\mathbf{T}(\mathbf{x}) = \lambda(a\mathbf{x})$ .) The elements of the eigenspace  $E_{\lambda}$  are all eigenvectors with eigenvalue  $\lambda$  and the zero vector. Indeed, a vector  $\mathbf{x} \in E_{\lambda}$  iff holds  $(\mathbf{T} - \lambda \mathbf{I})(\mathbf{x}) = \mathbf{0}$ , where  $\mathbf{I} \in L(V)$  is the identity operator, and this equation implies that  $\mathbf{x} = \mathbf{0}$  or  $\mathbf{T}(\mathbf{x}) = \lambda \mathbf{x}$ .

Eigenvalues and eigenvectors are notions defined on an operator, independently of any basis on the vector space. However, given a basis in a vector space, the eigenvalue-eigenvector equation can be expressed in terms of a matrix-vector product. This is summarized in the following result.

**Theorem 9.1.2.** If  $\mathbf{T}_{vv} \in F^{n,n}$  is the matrix of a linear operator  $\mathbf{T} \in L(V)$  in any ordered basis  $\mathcal{V} \subset V$ , then the eigenvalue  $\lambda \in \mathbb{F}$  and eigenvector components  $\mathbf{x}_v \in \mathbb{F}^n$  of the operator  $\mathbf{T}$  satisfy the eigenvalue-eigenvector equation

$$\mathbf{T}_{vv}\mathbf{x}_v = \lambda \mathbf{x}_v. \quad (9.2)$$

The eigenvalues and eigenvectors of a linear operator are the eigenvalues and eigenvectors of any matrix representation of the operator on any ordered basis of the vector space.

**Proof of Theorem 9.1.2:** The eigenvalue-eigenvector equation in (9.1) written in an ordered basis  $\mathcal{V} \subset V$  is given by

$$[\mathbf{T}(x_1 \mathbf{v}_1 + \cdots + x_n \mathbf{v}_n)]_v = \lambda [\mathbf{x}]_v \quad \Leftrightarrow \quad [[\mathbf{T}(\mathbf{v}_1)]_v, \cdots, [\mathbf{T}(\mathbf{v}_n)]_v] [\mathbf{x}]_v = \lambda [\mathbf{x}]_v,$$



that is, we obtain Eq. (9.2). It is simple to see that Eq. (9.2) is invariant under similarity transformations of the matrix  $T_{vv}$ . This property says that the components equation in (9.2) looks the same in any basis. Indeed, given any other ordered basis  $\tilde{V}$  of the vector space  $V$ , denote the change of basis matrix  $P = I_{\tilde{v}}$ . Then, multiply Eq. (9.2) by  $P^{-1}$ , that is,

$$P^{-1}T_{vv}x_v = \lambda P^{-1}x_v \iff (P^{-1}T_{vv}P)(P^{-1}x_v) = \lambda(P^{-1}x_v);$$

using the change of basis formulas  $T_{\tilde{v}\tilde{v}} = P^{-1}T_{vv}P$  and  $x_{\tilde{v}} = P^{-1}x_v$  we conclude that

$$T_{\tilde{v}\tilde{v}}x_{\tilde{v}} = \lambda x_{\tilde{v}}.$$

This establishes the Theorem. □

**EXAMPLE 9.1.1:** Consider the vector space  $\mathbb{R}^2$  with the standard basis  $\mathcal{S}$ . Find the eigenvalues and eigenvectors of the linear operator  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  with matrix in the standard basis given by

$$T = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}.$$

**SOLUTION:** This operator makes a reflection along the line  $x_1 = x_2$ , that is,

$$Tx = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ x_1 \end{bmatrix}.$$

From this definition we see that any non-zero vector proportional to  $v_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$  is left invariant by  $T$ , that is,

$$\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \iff T(v_1) = v_1.$$

So we conclude that  $v_1$  is an eigenvector of  $T$  with eigenvalue  $\lambda_1 = 1$ . Analogously, one can check that the vector  $v_2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$  satisfies the equation

$$\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} = - \begin{bmatrix} -1 \\ 1 \end{bmatrix} \iff T(v_2) = -v_2.$$

So we conclude that  $v_2$  is an eigenvector of  $T$  with eigenvalue  $\lambda_2 = -1$ . See Fig. 57.

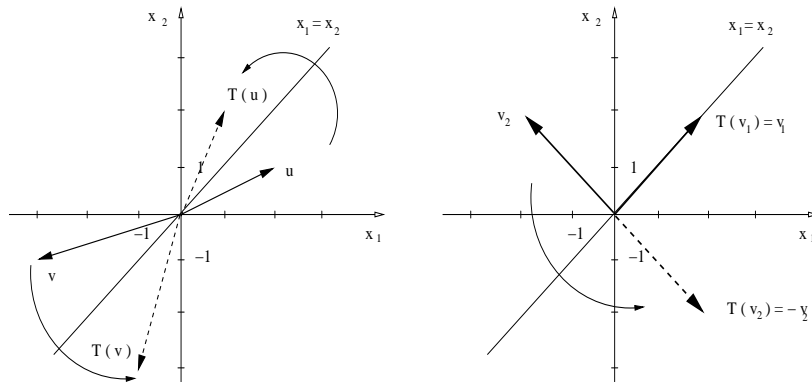


FIGURE 57. On the first picture we sketch the action of matrix  $T$  in Example 9.1.1, and on the second picture, and we sketch the eigenvectors  $v_1$  and  $v_2$  with eigenvalues  $\lambda_1 = 1$  and  $\lambda_2 = -1$ , respectively.

In this example the spectrum is  $\sigma_T = \{1, -1\}$  and the respective eigenspaces are

$$E_1 = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 1 \end{bmatrix}\right\}\right), \quad E_{-1} = \text{Span}\left(\left\{\begin{bmatrix} -1 \\ 1 \end{bmatrix}\right\}\right).$$

◁

**EXAMPLE 9.1.2:** Not every linear operator has eigenvalues and eigenvectors. Consider the vector space  $\mathbb{R}^2$  with standard basis and fix  $\theta \in (0, \pi)$ . The linear operator  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  with matrix

$$T = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}$$

acts on the plane rotating every vector by an angle  $\theta$  counterclockwise. Since  $\theta \in (0, \pi)$ , there is no line on the plane left invariant by the rotation. Therefore, **this operator has no eigenvalues and eigenvectors.**

◁

**EXAMPLE 9.1.3:** Consider the vector space  $\mathbb{R}^2$  with standard basis. Show that the linear operator  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  with matrix  $T = \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix}$  has the eigenvalues and eigenvectors

$$v_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad \lambda_1 = 4, \quad \text{and} \quad v_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \quad \lambda_2 = -2.$$

**SOLUTION:** We only verify that Eq. (9.1) holds for the vectors and scalars above, since

$$\begin{aligned} T v_1 &= \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \end{bmatrix} = 4 \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \lambda_1 v_1 \quad \Rightarrow \quad T v_1 = \lambda_1 v_1 \\ T v_2 &= \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} -2 \\ 2 \end{bmatrix} = -2 \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \lambda_2 v_2 \quad \Rightarrow \quad T v_2 = \lambda_2 v_2. \end{aligned}$$

In this example the spectrum is  $\sigma_T = \{4, -2\}$  and the respective eigenspaces are

$$E_4 = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 1 \end{bmatrix}\right\}\right), \quad E_{-2} = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ -1 \end{bmatrix}\right\}\right).$$

◁

**EXAMPLE 9.1.4:** Consider the vector space  $V = C^\infty(\mathbb{R}, \mathbb{R})$ . Show that the vector  $f(x) = e^{ax}$ , with  $a \neq 0$ , is an eigenvector with eigenvalue  $a$  of the linear operator  $D : V \rightarrow V$ , given by  $D(f)(x) = \frac{df}{dx}(x)$ .

**SOLUTION:** The proof is straightforward, since

$$D(f)(x) = \frac{df}{dx}(x) = ae^{ax} = a f(x) \quad \Rightarrow \quad D(f)(x) = a f(x).$$

◁

**EXAMPLE 9.1.5:** Consider again the vector space  $V = C^\infty(\mathbb{R}, \mathbb{R})$  and show that the vector  $f(x) = \cos(ax)$ , with  $a \neq 0$ , is an eigenvector with eigenvalue  $-a^2$  of the linear operator  $T : V \rightarrow V$ , given by  $T(f)(x) = \frac{d^2 f}{dx^2}(x)$ .

**SOLUTION:** Again, the proof is straightforward, since

$$T(f)(x) = \frac{d^2 f}{dx^2}(x) = -a^2 \cos(ax) = -a^2 f(x) \quad \Rightarrow \quad T(f)(x) = -a^2 f(x).$$

◁

We know that an eigenspace of a linear operator is not only a subset of the vector space, it is a subspace. Moreover, it is not any subspace, it is an invariant subspace under the linear operator. Given a vector space  $V$  and a linear operator  $\mathbf{T} \in L(V)$ , the subspace  $W \subset V$  is *invariant* under  $\mathbf{T}$  iff holds  $\mathbf{T}(W) \subset W$ .

**Theorem 9.1.3.** *The eigenspace  $E_\lambda$  of the linear operator  $\mathbf{T} \in L(V)$  corresponding to the eigenvalue  $\lambda$  is an invariant subspace of the vector space  $V$  under the operator  $\mathbf{T}$ .*

**Proof of Theorem 9.1.3:** We first show that  $E_\lambda$  is a subspace. Indded, pick up any vectors  $\mathbf{x}, \mathbf{y} \in E_\lambda$ , that is,  $\mathbf{T}(\mathbf{x}) = \lambda \mathbf{x}$  and  $\mathbf{T}(\mathbf{y}) = \lambda \mathbf{y}$ . Then, for all  $a, b \in \mathbb{F}$  holds,

$$\mathbf{T}(a\mathbf{x} + b\mathbf{y}) = a\mathbf{T}(\mathbf{x}) + b\mathbf{T}(\mathbf{y}) = a\lambda\mathbf{x} + b\lambda\mathbf{y} = \lambda(a\mathbf{x} + b\mathbf{y}) \Rightarrow (a\mathbf{x} + b\mathbf{y}) \in E_\lambda.$$

This shows that  $E_\lambda$  is a subspace. Now we need to show that  $\mathbf{T}(E_\lambda) \subset E_\lambda$ . This is straightforward, since for every vector  $\mathbf{x} \in E_\lambda$  we know that

$$\mathbf{T}(\mathbf{x}) = \lambda \mathbf{x} \in E_\lambda,$$

since the set  $E_\lambda$  is a subspace. Therefore,  $\mathbf{T}(E_\lambda) \subset E_\lambda$ . This establishes the Theorem.  $\square$

**9.1.2. Characteristic polynomial.** We now address the eigenvalue-eigenvector problem: Given a finite-dimensional vector space  $V$  over  $\mathbb{F}$  and a linear operator  $\mathbf{T} \in L(V)$ , find a scalar  $\lambda \in \mathbb{F}$  and a non-zero vector  $\mathbf{x} \in V$  solution of  $\mathbf{T}(\mathbf{x}) = \lambda \mathbf{x}$ . This problem is more complicated than solving a linear system of equations  $\mathbf{T}(\mathbf{x}) = \mathbf{b}$ , since in our case the source vector is  $\mathbf{b} = \lambda \mathbf{x}$ , which is not a given but part of the unknowns. One way to solve the eigenvalue-eigenvector problem is to first solve for the eigenvalues  $\lambda$  and then solve for the eigenvectors  $\mathbf{x}$ .

**Theorem 9.1.4.** *Let  $V$  be a finite-dimensional vector space over  $\mathbb{F}$  and let  $\mathbf{T} \in L(V)$ .*

(a) *The scalar  $\lambda \in \mathbb{F}$  is an eigenvalue of  $\mathbf{T}$  iff  $\lambda$  is solution of the equation*

$$\det(\mathbf{T} - \lambda \mathbf{I}) = 0.$$

(b) *Given  $\lambda \in \mathbb{F}$  eigenvalue of  $\mathbf{T}$ , the corresponding eigenvectors  $\mathbf{x} \in V$  are the non-zero solutions of the equation*

$$(\mathbf{T} - \lambda \mathbf{I})(\mathbf{x}) = \mathbf{0}.$$

**REMARK:** The determinant of a linear operator on a finite-dimensional vector space  $V$  was introduced in Def. 5.5.6 as the determinant of its associated matrix in any ordered basis of  $V$ . This definition is independent of the basis chosen in  $V$ , since the operator matrix transforms by a similarity transformation under a change of basis and the determinant is invariant under similarity transformations.

**Proof of Theorem 9.1.4:**

**Part (a):** The scalar  $\lambda$  and the vector  $\mathbf{x}$  are eigenvalue and eigenvector of  $\mathbf{T}$  iff holds

$$\mathbf{T}(\mathbf{x}) = \lambda \mathbf{x} \Leftrightarrow (\mathbf{T} - \lambda \mathbf{I})\mathbf{x} = \mathbf{0} \Leftrightarrow \det(\mathbf{T} - \lambda \mathbf{I}) = 0.$$

**Part (b):** This is simpler. Since  $\lambda$  is the scalar such that the operator  $(\mathbf{T} - \lambda \mathbf{I})$  is not invertible, this means that  $N(\mathbf{T} - \lambda \mathbf{I}) \neq \{\mathbf{0}\}$ , that is, there exists a solution  $\mathbf{x}$  to the linear equation  $(\mathbf{T} - \lambda \mathbf{I})\mathbf{x} = \mathbf{0}$ . It is simple to see that this solution is an eigenvector of  $\mathbf{T}$ . This establishes the Theorem.  $\square$

**Definition 9.1.5.** *Given a finite-dimensional vector space  $V$  and a linear operator  $\mathbf{T} \in L(V)$ , the function  $p(\lambda) = \det(\mathbf{T} - \lambda \mathbf{I})$  is called the **characteristic polynomial** of  $\mathbf{T}$ .*

The function  $p$  defined above is a polynomial in  $\lambda$ , which can be seen from the definition of determinant of a matrix. The eigenvalues of a linear operator are the roots of its characteristic polynomial.

**EXAMPLE 9.1.6:** Find the eigenvalues and eigenvectors of the linear operator  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , given by  $T = \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix}$ .

**SOLUTION:** We start computing the eigenvalues, which are the roots of the characteristic polynomial

$$p(\lambda) = \det(T - \lambda I) = \det\left(\begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}\right) = \begin{vmatrix} (1-\lambda) & 3 \\ 3 & (1-\lambda) \end{vmatrix},$$

hence  $p(\lambda) = (1 - \lambda)^2 - 9$ . The roots are

$$(\lambda - 1)^2 = 3^2 \quad \Rightarrow \quad \begin{cases} \lambda_1 = 4, \\ \lambda_2 = -2. \end{cases}$$

We now find the eigenvector for the eigenvalue  $\lambda_1 = 4$ . We solve the system  $(T - 4I)x = 0$  performing Gauss operation in the matrix

$$T - 4I = \begin{bmatrix} -3 & 3 \\ 3 & -3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -1 \\ 0 & 0 \end{bmatrix} \quad \Rightarrow \quad \begin{cases} x_1 = x_2, \\ x_2 \text{ free.} \end{cases}$$

Choosing  $x_2 = 1$  we obtain the eigenvector  $x_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ . In a similar way we find the eigenvector for the eigenvalue  $\lambda_2 = -2$ . We solve the linear system  $(T + 2I)x = 0$  performing Gauss operation in the matrix

$$T + 2I = \begin{bmatrix} 3 & 3 \\ 3 & 3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} \quad \Rightarrow \quad \begin{cases} x_1 = -x_2, \\ x_2 \text{ free.} \end{cases}$$

Choosing  $x_2 = -1$  we obtain the eigenvector  $x_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$ . These results  $\lambda_1, x_1$  and  $\lambda_2, x_2$  agree with Example 9.1.3. ◀

**9.1.3. Eigenvalue multiplicities.** We now introduce two numbers that give information regarding the size of eigenspaces. The first number determines the maximum possible size of an eigenspace, while the second number characterizes the actual size of an eigenspace.

**Definition 9.1.6.** Let  $\lambda_i \in \mathbb{F}$ , for  $i = 1, \dots, k$  be all the eigenvalues of a linear operator  $T \in L(V)$  on a vector space  $V$  over  $\mathbb{F}$ . Express the characteristic polynomial  $p$  associated with  $T$  as follows,

$$p(\lambda) = (\lambda - \lambda_1)^{r_1} \cdots (\lambda - \lambda_k)^{r_k} q(\lambda), \quad \text{with } q(\lambda_i) \neq 0,$$

and denote by  $s_i = \dim E_{\lambda_i}$ , the dimension of the eigenspaces corresponding to the eigenvalue  $\lambda_i$ . Then, the numbers  $r_i$  are called the **algebraic multiplicity** of the eigenvalue  $\lambda_i$ ; and the numbers  $s_i$  are called the **geometric multiplicity** of the eigenvalue  $\lambda_i$ .

**REMARK:** In the case that  $\mathbb{F} = \mathbb{C}$ , hence the characteristic polynomial is complex-valued, the polynomial  $q(\lambda) = 1$ . Indeed, when  $p$  is an  $n$ -degree complex-valued polynomial the Fundamental Theorem of Algebra says that  $p$  has  $n$  complex roots. Therefore, the characteristic polynomial has the form

$$p(\lambda) = (\lambda - \lambda_1)^{r_1} \cdots (\lambda - \lambda_k)^{r_k},$$

where  $r_1 + \cdots + r_k = n$ . On the other hand, in the case that  $\mathbb{F} = \mathbb{R}$  the characteristic polynomial is real-valued. In this case the polynomial  $q$  may have degree greater than zero.

**EXAMPLE 9.1.7:** For the following matrices find the algebraic and geometric multiplicities of their eigenvalues,

$$A = \begin{bmatrix} 3 & 2 \\ 0 & 3 \end{bmatrix}, \quad B = \begin{bmatrix} 3 & 1 & 1 \\ 0 & 3 & 2 \\ 0 & 0 & 1 \end{bmatrix}, \quad C = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 3 & 2 \\ 0 & 0 & 1 \end{bmatrix}.$$

**SOLUTION:** The eigenvalues of matrix  $A$  are the roots of the characteristic polynomial

$$p_a(\lambda) = \begin{vmatrix} (3-\lambda) & 2 \\ 0 & (3-\lambda) \end{vmatrix} = (\lambda-3)^2 \Rightarrow \lambda_1 = 3, \quad r_1 = 2.$$

So, the eigenvalue  $\lambda_1 = 3$  has algebraic multiplicity  $r_1 = 2$ . To find the geometric multiplicity we need to compute the eigenspace  $E_{\lambda_1}$ , which is the null space of the matrix  $(A - 3I)$ , that is,

$$A - 3I_2 = \begin{bmatrix} 0 & 2 \\ 0 & 0 \end{bmatrix} \Rightarrow x_2 = 0, \quad x_1 \text{ free} \Rightarrow x_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} x_1.$$

We have obtained

$$E_{\lambda_1} = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 0 \end{bmatrix}\right\}\right) \Rightarrow s_1 = \dim E_{\lambda_1} = 1.$$

In the case of matrix  $A$  the algebraic multiplicity is greater than the geometric multiplicity,  $2 = r_1 > s_1 = 1$ , since the greatest linearly independent set of eigenvectors for the eigenvalue  $\lambda_1$  contains only one vector.

The algebraic and geometric multiplicities for matrices  $B$  and  $C$  is computed in a similar way. These matrices differ only in one single matrix coefficient, the coefficient  $(1, 2)$ . This difference does not affect their eigenvalues, since we will show that both matrices  $B$  and  $C$  have the same eigenvalues with the same algebraic multiplicities. However, this difference is enough to change their eigenvectors, since we will show that matrices  $B$  and  $C$  have different eigenvectors with different geometric multiplicities. We start computing the characteristic polynomial of matrix  $B$ ,

$$p_b(\lambda) = \begin{vmatrix} (3-\lambda) & 1 & 1 \\ 0 & (3-\lambda) & 2 \\ 0 & 0 & (1-\lambda) \end{vmatrix} = -(\lambda-3)^2(\lambda-1),$$

which implies that  $\lambda_1 = 3$  has algebraic multiplicity  $r_1 = 2$  and  $\lambda_2 = 1$  has algebraic multiplicity  $r_2 = 1$ . To find the geometric multiplicities we need to find the corresponding eigenspaces. We start with  $\lambda_1 = 3$ ,

$$B - 3I_3 = \begin{bmatrix} 0 & 1 & 1 \\ 0 & 0 & 2 \\ 0 & 0 & -2 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \Rightarrow \begin{cases} x_1 \text{ free,} \\ x_2 = 0, \\ x_3 = 0, \end{cases}$$

which implies that

$$E_{\lambda_1} = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}\right\}\right) \Rightarrow s_1 = 1.$$

The geometric multiplicity for the eigenvalue  $\lambda_2 = 1$  is computed as follows,

$$B - I_3 = \begin{bmatrix} 2 & 1 & 1 \\ 0 & 2 & 2 \\ 0 & 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix} \Rightarrow \begin{cases} x_1 = 0, \\ x_2 = -x_3, \\ x_3 \text{ free,} \end{cases}$$

which implies that

$$E_{\lambda_2} = \text{Span}\left(\left\{\begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix}\right\}\right) \Rightarrow s_2 = 1.$$

We then conclude  $2 = r_1 > s_1 = 1$  and  $r_2 = s_2 = 1$ .

Finally, we compute the characteristic polynomial of matrix  $C$ ,

$$p_c(\lambda) = \begin{vmatrix} (3-\lambda) & 0 & 1 \\ 0 & (3-\lambda) & 2 \\ 0 & 0 & (1-\lambda) \end{vmatrix} = -(\lambda-3)^2(\lambda-1)$$

so, we obtain the same eigenvalues we had for matrix  $B$ , that is,  $\lambda_1 = 3$  has algebraic multiplicity  $r_1 = 2$  and  $\lambda_2 = 1$  has algebraic multiplicity  $r_2 = 1$ . The geometric multiplicity of  $\lambda_1$  is computed as follows,

$$C - 3I_3 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 2 \\ 0 & 0 & -2 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \Rightarrow \begin{cases} x_1 \text{ free,} \\ x_2 \text{ free,} \\ x_3 = 0, \end{cases}$$

which implies that

$$E_{\lambda_1} = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}\right\}\right) \Rightarrow s_1 = 2.$$

In this case we obtained  $2 = r_1 = s_1$ . The geometric multiplicity for the eigenvalue  $\lambda_2 = 1$  is computed as follows,

$$C - I_3 = \begin{bmatrix} 2 & 0 & 1 \\ 0 & 2 & 2 \\ 0 & 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & \frac{1}{2} \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix} \Rightarrow \begin{cases} x_1 = -\frac{1}{2}x_3, \\ x_2 = -x_3, \\ x_3 \text{ free,} \end{cases}$$

which implies that (choosing  $x_3 = 2$ ),

$$E_{\lambda_2} = \text{Span}\left(\left\{\begin{bmatrix} -1 \\ -2 \\ 2 \end{bmatrix}\right\}\right) \Rightarrow s_2 = 1.$$

We then conclude  $r_1 = s_1 = 2$  and  $r_2 = s_2 = 1$ .

**REMARK:** Comparing the results for matrix  $B$  and  $C$  we see that a change in just one matrix coefficient can change the eigenspaces even in the case where the eigenvalues do not change. In fact, it can be shown that the eigenvalues are continuous functions of the matrix coefficients while the eigenvectors are not continuous functions. This means that a small change in the matrix coefficients produces a small change in the eigenvalues but it might produce a big change in the eigenvectors, just like in this Example.  $\triangleleft$

**9.1.4. Operators with distinct eigenvalues.** The eigenvectors corresponding to distinct eigenvalues form a linearly independent set. In other words, eigenspaces corresponding to different eigenvalues have trivial intersection.

**Theorem 9.1.7.** *If the eigenvalues  $\lambda_1, \dots, \lambda_k \in \mathbb{F}$ , with  $k \geq 1$ , of the linear operator  $T \in L(V)$  are all different, then the set of corresponding eigenvectors  $\{\mathbf{x}_1, \dots, \mathbf{x}_k\} \subset V$  is linearly independent.*

**Proof of Theorem 9.1.7:** If  $k = 1$ , the Theorem is trivially true, so assume  $k \geq 2$ . Let  $c_1, \dots, c_k \in F$  be scalars such that

$$c_1 \mathbf{x}_1 + \dots + c_k \mathbf{x}_k = \mathbf{0}. \quad (9.3)$$

Now perform the following two steps: First, apply the operator  $T$  to both sides of the equation above. Since  $T(\mathbf{x}_i) = \lambda_i \mathbf{x}_i$ , we obtain,

$$c_1 \lambda_1 \mathbf{x}_1 + \dots + c_k \lambda_k \mathbf{x}_k = \mathbf{0}. \quad (9.4)$$

Second, multiply Eq. (9.3) by  $\lambda_1$  and subtract it from Eq.(9.4). The result is

$$c_2(\lambda_2 - \lambda_1) \mathbf{x}_2 + \dots + c_k(\lambda_k - \lambda_1) \mathbf{x}_k = \mathbf{0}. \quad (9.5)$$

Notice that all factors  $\lambda_i - \lambda_1 \neq 0$  for  $i = 2, \dots, k$ . Repeat these two steps: First, apply the operator  $T$  on both sides of Eq. (9.5), that is,

$$c_2(\lambda_2 - \lambda_1) \lambda_2 \mathbf{x}_2 + \dots + c_k(\lambda_k - \lambda_1) \lambda_k \mathbf{x}_k = \mathbf{0}; \quad (9.6)$$

second, multiply Eq. (9.5) by  $\lambda_2$  and subtract it from Eq. (9.6), that is,

$$c_2(\lambda_2 - \lambda_1) (\lambda_3 - \lambda_2) \mathbf{x}_3 + \dots + c_k(\lambda_k - \lambda_1) (\lambda_k - \lambda_2) \mathbf{x}_k = \mathbf{0}. \quad (9.7)$$

Repeat the idea in these two steps until one reaches the equation

$$c_k(\lambda_k - \lambda_1) \dots (\lambda_k - \lambda_{k-1}) \mathbf{x}_k = \mathbf{0}.$$

Since the eigenvalues  $\lambda_i$  are all different, we conclude that  $c_k = 0$ . Introducing this information at the very beginning we get that

$$c_1 \mathbf{x}_1 + \dots + c_{k-1} \mathbf{x}_{k-1} = \mathbf{0}.$$

Repeating the whole procedure we conclude that  $c_{k-1} = 0$ . In this way one shows that all coefficient  $c_1 = \dots = c_k = 0$ . Therefore, the set  $\{\mathbf{x}_1, \dots, \mathbf{x}_k\}$  is linearly independent. This establishes the Theorem.  $\square$

**Further reading.** A detailed presentation of eigenvalues and eigenvectors of a matrix can be found in Sections 5.1 and 5.2 in Lay's book [2], while a shorter and deeper summary can be found in Section 7.1 in Meyer's book [3]. Also see Chapter 4 in Hassani's book [1].

## 9.1.5. Exercises.

- 9.1.1.- Find the spectrum and all eigenspaces of the operators  $A : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  and  $B : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ ,

$$A = \begin{bmatrix} -10 & -7 \\ 14 & 11 \end{bmatrix},$$

$$B = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 3 & 1 \\ 1 & -1 & 2 \end{bmatrix}.$$

- 9.1.2.- Show that for  $\theta \in (0, \pi)$  the rotation operator  $R(\theta) : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  has no eigenvalues, where

$$R(\theta) = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}.$$

Consider now matrix above as a linear operator  $R(\theta) : \mathbb{C}^2 \rightarrow \mathbb{C}^2$ . Show that this linear operator has eigenvalues, and find them.

- 9.1.3.- Let  $A : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  be the linear operator given by

$$A = \begin{bmatrix} 2 & -1 & 3 \\ 0 & 1 & h \\ 0 & 0 & 2 \end{bmatrix}.$$

- Find all the eigenvalues and their corresponding algebraic multiplicities of the matrix  $A$ .
- Find the value(s) of  $h \in \mathbb{R}$  such that the matrix  $A$  above has a two-dimensional eigenspace, and find a basis for this eigenspace.
- Set  $h = 1$ , and find a basis for all the eigenspaces of matrix  $A$  above.

- 9.1.4.- Find all the eigenvalues with their corresponding algebraic multiplicities, and find all the associated eigenspaces of the matrix  $A \in \mathbb{R}^{3,3}$  given by

$$A = \begin{bmatrix} 2 & 1 & 1 \\ 0 & 2 & 3 \\ 0 & 0 & 1 \end{bmatrix}.$$

- 9.1.5.- Let  $k \in \mathbb{R}$  and consider the matrix  $A \in \mathbb{R}^{4,4}$  given by

$$A = \begin{bmatrix} 2 & -2 & 4 & -1 \\ 0 & 3 & k & 0 \\ 0 & 0 & 2 & 4 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

- Find the eigenvalues of  $A$  and their algebraic multiplicity.
- Find the number  $k$  such that matrix  $A$  has an eigenspace  $E_\lambda$  that is two dimensional, and find a basis for this  $E_\lambda$ .

- 9.1.6.- Comparing the characteristic polynomials for  $A \in \mathbb{F}^{n,n}$  and  $A^T$ , show that these two matrices have the same eigenvalues.

- 9.1.7.- Let  $A \in \mathbb{R}^{3,3}$  be an invertible matrix with eigenvalues 2,  $-1$  and 3. Find the eigenvalues of:

- $A^{-1}$ .
- $A^k$ , for any  $k \in \mathbb{N}$ .
- $A^2 - A$ .



## 9.2. DIAGONALIZABLE OPERATORS

**9.2.1. Eigenvectors and diagonalization.** In this Section we study linear operators on a finite dimensional vector space that have a complete set of eigenvectors. This means that there exists a basis of the vector space formed with eigenvectors of the linear operator. We show that these operators are diagonalizable, that is, the matrix of the operator in the basis of its own eigenvectors is diagonal. We end this Section showing that it is not difficult to define functions of operators in the case that the operator is diagonalizable.

**Definition 9.2.1.** A linear operator  $\mathbf{T} \in L(V)$  defined on an  $n$ -dimensional vector space  $V$  has a **complete set of eigenvectors** iff there exists a linearly independent set formed with  $n$  eigenvectors of  $\mathbf{T}$ .

In other words, a linear operator  $\mathbf{T}$  has a complete set of eigenvectors iff there exists a basis of  $V$  formed by eigenvectors of  $\mathbf{T}$ . Not every linear operator has a complete set of eigenvectors. For example, a linear operator without a complete set of eigenvectors is a rotation on a plane  $\mathbf{R}(\theta) : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  by an angle  $\theta \in (0, \pi)$ . This particular operator has not eigenvectors at all.

**EXAMPLE 9.2.1:** Matrix  $\mathbf{B}$  in Example 9.1.7 does not have a complete set of eigenvectors. The largest linearly independent set of eigenvectors of matrix  $\mathbf{B}$  contains only two vectors, one possibility is shown below.

$$\mathbf{B} = \begin{bmatrix} 3 & 1 & 1 \\ 0 & 3 & 2 \\ 0 & 0 & 1 \end{bmatrix}, \quad X = \left\{ \mathbf{x}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \mathbf{x}_2 = \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix} \right\}.$$

Matrix  $\mathbf{C}$  in Example 9.1.7 has a complete set of eigenvectors, indeed,

$$\mathbf{C} = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 3 & 2 \\ 0 & 0 & 1 \end{bmatrix}, \quad X = \left\{ \mathbf{x}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \mathbf{x}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \mathbf{x}_3 = \begin{bmatrix} 1 \\ -2 \\ 2 \end{bmatrix} \right\}.$$

◀

We now introduce the notion of a diagonalizable operator.

**Definition 9.2.2.** A linear operator  $\mathbf{T} \in L(V)$  defined on a finite dimensional vector space  $V$  is called **diagonalizable** iff there exists a basis  $\tilde{\mathcal{V}}$  of  $V$  such that the matrix  $\mathbf{T}_{\tilde{\mathbf{v}}\tilde{\mathbf{v}}}$  is diagonal.

Recall that a square matrix  $\mathbf{D} = [D_{ij}]$  is diagonal iff  $D_{ij} = 0$  for  $i \neq j$ . We denote an  $n \times n$  diagonal matrix by  $\mathbf{D} = \text{diag}[D_1, \dots, D_n]$ , so we use only one index to label the diagonal elements,  $D_{ii} = D_i$ . Examples of  $3 \times 3$  diagonal matrices are given by

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 4 \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} D_{11} & 0 & 0 \\ 0 & D_{22} & 0 \\ 0 & 0 & D_{33} \end{bmatrix}.$$

Our first result is to show that these two notions given in Definitions 9.2.1 and 9.2.2 are equivalent.

**Theorem 9.2.3.** A linear operator  $\mathbf{T} \in L(V)$  defined on an  $n$ -dimensional vector space is diagonalizable iff the operator  $\mathbf{T}$  has a complete set of eigenvectors. Furthermore, if  $\lambda_i, \mathbf{x}_i$ , for  $i = 1, \dots, n$ , are eigenvalue-eigenvector pairs of  $\mathbf{T}$ , then the matrix of the operator  $\mathbf{T}$  in the ordered eigenvector basis  $\tilde{\mathcal{V}} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$  is diagonal with the eigenvalues on the diagonal, that is,  $\mathbf{T}_{\tilde{\mathbf{v}}\tilde{\mathbf{v}}} = \text{diag}[\lambda_1, \dots, \lambda_n]$ .

**Proof of Theorem 9.2.3:**

( $\Rightarrow$ ) Since  $\mathbf{T}$  is diagonalizable, we know that there exists a basis  $\tilde{\mathcal{V}} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$  such that its matrix is diagonal, that is,  $\mathsf{T}_{\tilde{v}\tilde{v}} = \text{diag}[\lambda_1, \dots, \lambda_n]$ . This implies that for  $i = 1, \dots, n$  holds

$$\mathsf{T}_{\tilde{v}\tilde{v}} \mathbf{x}_{i\tilde{v}} = \lambda_i \mathbf{x}_{i\tilde{v}} \quad \Leftrightarrow \quad \mathbf{T}(\mathbf{x}_i) = \lambda_i \mathbf{x}_i.$$

We then conclude that  $\mathbf{x}_i$  is an eigenvector of the operator  $\mathbf{T}$  with eigenvalue  $\lambda_i$ . Since  $\tilde{\mathcal{V}}$  is a basis of the vector space  $V$ , this means that the operator  $\mathbf{T}$  has a complete set of eigenvectors.

( $\Leftarrow$ ) Since the set of eigenvectors  $\tilde{\mathcal{V}} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$  of the operator  $\mathbf{T}$  with corresponding eigenvalues  $\lambda_1, \dots, \lambda_n$  form a basis of  $V$ , then the eigenvalue-eigenvector equation for  $i = 1, \dots, n$

$$\mathbf{T}(\mathbf{x}_i) = \lambda_i \mathbf{x}_i.$$

imply that the matrix  $\mathsf{T}_{\tilde{v}\tilde{v}}$  is diagonal, since

$$\mathsf{T}_{\tilde{v}\tilde{v}} = [[\mathbf{T}(\mathbf{x}_1)]_{\tilde{v}}, \dots, [\mathbf{T}(\mathbf{x}_n)]_{\tilde{v}}] = [\lambda_1 [\mathbf{x}_1]_{\tilde{v}}, \dots, \lambda_n [\mathbf{x}_n]_{\tilde{v}}] = [\lambda_1 \mathbf{e}_1, \dots, \lambda_n \mathbf{e}_n],$$

so we arrive at the equation  $\mathsf{T}_{\tilde{v}\tilde{v}} = \text{diag}[\lambda_1, \dots, \lambda_n]$ . This establishes the Theorem.  $\square$

**EXAMPLE 9.2.2:** Show that the linear operator  $\mathbf{T} \in L(\mathbb{R}^2)$  with matrix  $\mathsf{T}_{ss} = \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix}$  in the standard basis of  $\mathcal{S}$  of  $\mathbb{R}^2$  is diagonalizable.

**SOLUTION:** In Example 9.1.6 we obtained that the eigenvalues and eigenvectors of matrix  $\mathsf{T}_{ss}$  are given by

$$\lambda_1 = 4, \quad \mathbf{x}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \text{and} \quad \lambda_2 = -2, \quad \mathbf{x}_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}.$$

We now affirm that the matrix of the linear operator  $\mathbf{T}$  is diagonal in the ordered basis formed by the eigenvectors above,

$$\tilde{\mathcal{V}} = \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ -1 \end{bmatrix} \right).$$

First notice that the set  $\tilde{\mathcal{V}}$  is linearly independent, so it is a basis for  $\mathbb{R}^2$ . Second, the change of basis matrix  $\mathsf{P} = \mathsf{l}_{\tilde{v}s}$  is given by

$$\mathsf{P} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \quad \Rightarrow \quad \mathsf{P}^{-1} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

Third, the result of the change of basis is a diagonal matrix:

$$\mathsf{P}^{-1} \mathsf{T}_{ss} \mathsf{P} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 4 & -2 \\ 4 & 2 \end{bmatrix} = \begin{bmatrix} 4 & 0 \\ 0 & -2 \end{bmatrix} = \mathsf{T}_{\tilde{v}\tilde{v}}.$$

As stated in Theorem 9.2.3, the diagonal elements in  $\mathsf{T}_{\tilde{v}\tilde{v}}$  are precisely the eigenvalues of the operator  $\mathbf{T}$ , in the same order as the eigenvectors in the ordered basis  $\tilde{\mathcal{V}}$ .  $\triangleleft$

The statement in Definition 9.2.2 can be expressed as a statement between matrices. A matrix  $\mathbf{A} \in \mathbb{F}^{n,n}$  is diagonalizable if there exists an invertible matrix  $\mathsf{P} \in \mathbb{F}^{n,n}$  and a diagonal matrix  $\mathbf{D} \in \mathbb{F}^{n,n}$  such that

$$\mathbf{A} = \mathsf{P} \mathbf{D} \mathsf{P}^{-1}.$$

These two notions are equivalent, since  $\mathbf{A}$  is the matrix of an linear operator  $\mathbf{T} \in L(V)$  in the standard basis  $\mathcal{S}$  of  $V$ , that is,  $\mathbf{A} = \mathsf{T}_{ss}$ . The similarity transformation above can be expressed as

$$\mathbf{D} = \mathsf{P}^{-1} \mathsf{T}_{ss} \mathsf{P}.$$

Denoting  $P = I_{\tilde{v}_s}$  as the change of basis matrix from the standard basis  $\mathcal{S}$  to a basis  $\tilde{\mathcal{V}}$ , we conclude that

$$D = T_{\tilde{v}\tilde{v}},$$

That is, the matrix of operator  $T$  is diagonal in the basis  $\tilde{\mathcal{V}}$ . Furthermore, Theorem 9.2.3 can also be expressed in terms of similarity transformations between matrices as follows: *A square matrix has a complete set of eigenvectors iff the matrix is similar to a diagonal matrix.* This active point of view is common in the literature.

**EXAMPLE 9.2.3:** Show that the linear operator  $T: \mathbb{R}^2 \rightarrow \mathbb{R}^2$  with matrix  $T = \begin{bmatrix} 1 & 2 \\ 3 & 6 \end{bmatrix}$  in the standard basis of  $\mathbb{R}^2$  is diagonalizable. Find a similarity transformation that converts matrix  $T$  into a diagonal matrix.

**SOLUTION:** To find out whether  $T$  is diagonalizable or not we need to compute its eigenvectors, so we start with its eigenvalues. The characteristic polynomial is

$$p(\lambda) = \begin{vmatrix} 1-\lambda & 2 \\ 3 & 6-\lambda \end{vmatrix} = \lambda(\lambda-7) = 0 \quad \Rightarrow \quad \begin{cases} \lambda_1 = 0, \\ \lambda_2 = 7. \end{cases}$$

Since the eigenvalues are different, we know that the corresponding eigenvectors form a linearly independent set (by Theorem 9.1.7), and so matrix  $T$  has a complete set of eigenvectors. So  $A$  is diagonalizable. The corresponding eigenvectors are the non-zero vectors in the null spaces  $N(T)$  and  $N(T - 7I_2)$ , which are computed as follows:

$$\begin{aligned} \begin{bmatrix} 1 & 2 \\ 3 & 6 \end{bmatrix} &\rightarrow \begin{bmatrix} 1 & 2 \\ 0 & 0 \end{bmatrix} &\Rightarrow & x_1 = -2x_2 &\Rightarrow & x_1 = \begin{bmatrix} -2 \\ 1 \end{bmatrix}; \\ \begin{bmatrix} -6 & 2 \\ 3 & -1 \end{bmatrix} &\rightarrow \begin{bmatrix} 3 & -1 \\ 0 & 0 \end{bmatrix} &\Rightarrow & 3x_1 = x_2 &\Rightarrow & x_2 = \begin{bmatrix} 1 \\ 3 \end{bmatrix}. \end{aligned}$$

Since the set  $\mathcal{V} = \{x_1, x_2\}$  is a complete set of eigenvectors of  $T$ , we conclude that  $T$  is diagonalizable. Proposition 9.2.3 says that

$$D = P^{-1}TP, \quad \text{where } D = \begin{bmatrix} 0 & 0 \\ 0 & 7 \end{bmatrix}, \quad P = \begin{bmatrix} -2 & 1 \\ 1 & 3 \end{bmatrix}.$$

We finally verify this the equation above is correct:

$$PDP^{-1} = \begin{bmatrix} -2 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & 7 \end{bmatrix} \frac{1}{7} \begin{bmatrix} -3 & 1 \\ 1 & 2 \end{bmatrix} = \begin{bmatrix} -2 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 1 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 3 & 6 \end{bmatrix} = T.$$

◁

**9.2.2. Functions of diagonalizable operators.** We have seen in Section 5.3 that the set of all linear operators  $L(V)$  on a vector space  $V$  is itself an algebra, since both the linear combination of linear operators in  $L(V)$  and the composition of linear operators in  $L(V)$  are again linear operators in  $L(V)$ . We have also seen that this algebra structure on  $L(V)$  makes possible to introduce polynomial functions of linear operators. Indeed, given scalars  $a_0, \dots, a_n \in \mathbb{F}$ , an operator-valued polynomial of degree  $n$  is a function  $p: L(V) \rightarrow L(V)$ ,

$$p(T) = a_0 I_V + a_1 T + \dots + a_n T^n.$$

In this Section we introduce, on the one hand, functions of operators more general than polynomial functions, but on the other hand, we define these functions only on diagonalizable operators. The reason of this restriction is that functions of diagonalizable operators are simple to define. It is possible to define general functions of non-diagonalizable operators, but they are more involved, and will not be studied in these notes.

Before computing a general function of a diagonalizable operator, we start finding simple expressions for the power function and a polynomial function of a diagonalizable operator. These formulas simplify the results we have found in Section 5.3.

**Theorem 9.2.4.** *If  $T \in L(V)$  is a diagonalizable linear operator with matrix  $T = P D P^{-1}$  in the standard ordered basis of an  $n$ -dimensional vector space  $V$ , where  $D, P \in \mathbb{F}^{n,n}$  are a diagonal and an invertible matrix, respectively, then, for every  $k \in \mathbb{N}$  holds*

$$T^k = P D^k P^{-1}.$$

**Proof of Theorem 9.2.4:** Consider the case  $k = 2$ . A simple calculation shows

$$T^2 = T T = P D P^{-1} P D P^{-1} = P D D P^{-1} = P D^2 P^{-1}.$$

Suppose now that the formula holds for  $k \in \mathbb{N}$ , that is,  $T^k = P D^k P^{-1}$ , and let us show that it also holds for  $k + 1$ . Indeed,

$$T^{(k+1)} = T T^k = P D P^{-1} P D^k P^{-1} = P D D^k P^{-1} = P D^{(k+1)} P^{-1}.$$

This establishes the Theorem. □

**EXAMPLE 9.2.4:** Given the matrix  $T \in \mathbb{R}^{3,3}$ , compute  $T^k$ , where  $k \in \mathbb{N}$  and  $T = \begin{bmatrix} 1 & 2 \\ 3 & 6 \end{bmatrix}$ .

**SOLUTION:** From Example 9.2.3 we know that  $T$  is diagonalizable, and that

$$D = \begin{bmatrix} 0 & 0 \\ 0 & 7 \end{bmatrix}, \quad P = \begin{bmatrix} -2 & 1 \\ 1 & 3 \end{bmatrix} \quad \Rightarrow \quad P^{-1} = \frac{1}{7} \begin{bmatrix} -3 & 1 \\ 1 & 2 \end{bmatrix}.$$

Therefore,

$$T^k = P D^k P^{-1} = \begin{bmatrix} -2 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & 7^k \end{bmatrix} \frac{1}{7} \begin{bmatrix} -3 & 1 \\ 1 & 2 \end{bmatrix} = 7^{(k-1)} \begin{bmatrix} -2 & 1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & 7 \end{bmatrix} \frac{1}{7} \begin{bmatrix} -3 & 1 \\ 1 & 2 \end{bmatrix}.$$

The final result is

$$T^k = 7^{(k-1)} T.$$

◁

It is simple to generalize Theorem 9.2.4 to polynomial functions.

**Theorem 9.2.5.** *Assume the hypotheses given in Theorem 9.2.4 and denote the diagonal matrix as  $D = \text{diag}[\lambda_1, \dots, \lambda_n]$ . If  $p : \mathbb{F} \rightarrow \mathbb{F}$  is a polynomial of degree  $k \in \mathbb{N}$ , then holds*

$$p(T) = P p(D) P^{-1}, \quad \text{where } p(D) = \text{diag}[p(\lambda_1), \dots, p(\lambda_n)].$$

**Proof of Theorem 9.2.5:** Given scalars  $a_0, \dots, a_k \in \mathbb{F}$ , denote the polynomial  $p$  by

$$p(x) = a_0 + a_1 x + \dots + a_k x^k.$$

Then,

$$\begin{aligned} p(T) &= a_0 I_n + a_1 T + \dots + a_k T^k \\ &= a_0 P P^{-1} + a_1 P D P^{-1} + \dots + a_k P D^k P^{-1} \\ &= P (a_0 I_n + a_1 D + \dots + a_k D^k) P^{-1}. \end{aligned}$$

Noticing that

$$p(D) = a_0 I_n + a_1 D + \dots + a_k D^k = \text{diag}[p(\lambda_1), \dots, p(\lambda_n)],$$

we conclude that

$$p(T) = P p(D) P^{-1}.$$

This establishes the Theorem. □

**9.2.3. The exponential of diagonalizable operators.** Functions that admit a convergent power series expansion can be defined on diagonalizable operators in the same way as polynomial functions in Theorem 9.2.5. We consider first an important particular case, the exponential function  $f(x) = e^x$ . The exponential function is usually defined as the inverse of the natural logarithm function  $g(x) = \ln(x)$ , which in turn is defined as the area under the graph of the function  $h(x) = 1/x$  from 1 to  $x$ , that is,

$$\ln(x) = \int_1^x \frac{1}{y} dy \quad x \in (0, \infty).$$

It is not clear how to use this definition of the exponential function on real numbers to extend it to operators. However, one shows that the exponential function on real numbers has several properties, among them that it can be expressed as a convergent infinite power series,

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!} = 1 + x + \frac{x^2}{2!} + \cdots + \frac{x^k}{k!} + \cdots.$$

It can be proven that defining the exponential function on real numbers as the convergent power series above is equivalent to the definition given earlier as the inverse of the natural logarithm. However, only the power series expression provides the path to generalize the exponential function to a diagonalizable linear operator.

**Definition 9.2.6.** The *exponential* of a linear operator  $\mathbf{T} \in L(V)$  on a finite dimensional vector space, denoted as  $e^{\mathbf{T}}$ , is given by the infinite sum

$$e^{\mathbf{T}} = \sum_{k=0}^{\infty} \frac{\mathbf{T}^k}{k!}.$$

The following result says that the infinite sum of operators given in the Definition 9.2.6 converges in the case that the operator is diagonal or the operator is diagonalizable.

**Theorem 9.2.7.** If a linear operator  $\mathbf{T} \in L(V)$  on a finite dimensional vector space is diagonal, that is, there exists an ordered basis such that  $\mathbf{T} = \text{diag}[\lambda_1, \dots, \lambda_n]$ , then

$$e^{\mathbf{T}} = \text{diag}[e^{\lambda_1}, \dots, e^{\lambda_n}].$$

If a linear operator  $\mathbf{T} \in L(V)$  on a finite dimensional vector space is diagonalizable, that is, there exists an ordered basis such that  $\mathbf{T} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$ , where  $\mathbf{D} = \text{diag}[\lambda_1, \dots, \lambda_n]$ , then

$$e^{\mathbf{T}} = \mathbf{P}e^{\mathbf{D}}\mathbf{P}^{-1}.$$

**Proof of Theorem 9.2.7:** Let  $\mathbf{T}$  be the matrix of the operator  $\mathbf{T} \in L(V)$  in an ordered basis  $\mathcal{V} \subset V$ . For every  $N \in \mathbb{N}$  introduce the partial sum  $S_N(\mathbf{T})$  as follows,

$$S_N(\mathbf{T}) = \sum_{k=0}^N \frac{\mathbf{T}^k}{k!},$$

which is a well defined polynomial in  $\mathbf{T}$ . Assume now that  $\mathbf{T}$  is diagonalizable, so there exist an invertible matrix  $\mathbf{P}$  and a diagonal matrix  $\mathbf{D}$  such that  $\mathbf{T} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$ . Then, a straightforward computation shows that

$$S_N(\mathbf{T}) = \sum_{k=0}^N \frac{\mathbf{P}\mathbf{D}^k\mathbf{P}^{-1}}{k!} = \mathbf{P} \left( \sum_{k=0}^N \frac{\mathbf{D}^k}{k!} \right) \mathbf{P}^{-1}. \quad (9.8)$$

In the particular case that  $\mathbf{T}$  is diagonal, then  $\mathbf{P} = \mathbf{I}_n$  and  $\mathbf{T} = \text{diag}[\lambda_1, \dots, \lambda_n]$ , so the equation above reduced to

$$S_N(\mathbf{T}) = \sum_{k=0}^N \frac{\mathbf{T}^k}{k!} = \text{diag}\left[\sum_{k=0}^N \frac{\lambda_1^k}{k!}, \dots, \sum_{k=0}^N \frac{\lambda_n^k}{k!}\right].$$

In the expression above we can compute the limit as  $N \rightarrow \infty$ , obtaining

$$e^{\mathbf{T}} = \text{diag}[e^{\lambda_1}, \dots, e^{\lambda_n}].$$

This result establishes the first part in Theorem 9.2.7. In the case that  $\mathbf{T}$  is diagonalizable, we go back to Eq. (9.8), where we now denote  $\mathbf{D} = \text{diag}[\lambda_1, \dots, \lambda_n]$ . Since  $\mathbf{D}$  is a diagonal matrix, we can compute the limit  $N \rightarrow \infty$  in Eq. (9.8), that is,

$$e^{\mathbf{T}} = S_\infty(\mathbf{T}) = \mathbf{P} e^{\mathbf{D}} \mathbf{P}^{-1},$$

where we have denoted  $e^{\mathbf{D}} = \text{diag}[e^{\lambda_1}, \dots, e^{\lambda_n}]$ . This establishes the Theorem.  $\square$

**EXAMPLE 9.2.5:** For every  $t \in \mathbb{R}$  find the value of the exponential function  $e^{\mathbf{A}t}$ , where

$$\mathbf{A} = \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix}.$$

**SOLUTION:** From Example 9.2.2 we know that  $\mathbf{A}$  is diagonalizable, and that  $\mathbf{A} = \mathbf{P} \mathbf{D} \mathbf{P}^{-1}$ , where

$$\mathbf{D} = \begin{bmatrix} 4 & 0 \\ 0 & -2 \end{bmatrix}, \quad \mathbf{P} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \Rightarrow \mathbf{P}^{-1} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

Therefore,

$$\mathbf{A}t = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 4t & 0 \\ 0 & -2t \end{bmatrix} \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

and Theorem 9.2.7 imply that

$$e^{\mathbf{A}t} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} e^{4t} & 0 \\ 0 & e^{-2t} \end{bmatrix} \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

$\triangleleft$

The function introduced in Example 9.2.5 above can be seen as an operator-valued function  $f: \mathbb{R} \rightarrow \mathbb{R}^{2,2}$  given by

$$f(t) = e^{\mathbf{A}t}, \quad \mathbf{A} \in \mathbb{R}^{2,2}.$$

It can be shown that this function is actually differentiable, and that

$$\frac{df}{dt}(t) = \mathbf{A} e^{\mathbf{A}t}.$$

A more precise statement is the following.

**Theorem 9.2.8.** *If  $\mathbf{T} \in L(V)$  is a diagonalizable operator on a finite dimensional vector space over  $\mathbb{R}$ , then the operator-valued function  $\mathbf{F}: \mathbb{R} \rightarrow L(V)$ , defined as  $\mathbf{F}(x) = e^{\mathbf{T}x}$  for all  $x \in \mathbb{R}$ , is differentiable and*

$$\frac{d\mathbf{F}}{dx}(x) = \mathbf{T} e^{\mathbf{T}x} = e^{\mathbf{T}x} \mathbf{T}.$$

**Proof of Theorem 9.2.8:** Fix an ordered basis  $\mathcal{V} \subset V$  and denote by  $\mathbf{T}$  the matrix of  $\mathbf{T}$  in that basis. Since  $\mathbf{T}$  is diagonalizable, then  $\mathbf{T} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$ , with  $\mathbf{D} = \text{diag}[\lambda_1, \dots, \lambda_n]$ . Denoting by  $\mathbf{F}$  the matrix of  $\mathbf{F}$  in the basis  $\mathcal{V}$ , it is simple to see that

$$\frac{d\mathbf{F}}{dx}(x) = \frac{d}{dx}(\mathbf{P} e^{\mathbf{D}x} \mathbf{P}^{-1}) = \mathbf{P} \left( \frac{d}{dx} e^{\mathbf{D}x} \right) \mathbf{P}^{-1}.$$

It is not difficult to see that the expression of the far right in equation above is given by

$$\frac{d}{dx} e^{\mathbf{D}x} = \text{diag} \left[ \frac{d}{dx} (e^{\lambda_1 x}), \dots, \frac{d}{dx} (e^{\lambda_n x}) \right] = \text{diag} [\lambda_1 e^{\lambda_1 x}, \dots, \lambda_n e^{\lambda_n x}] = \mathbf{D} e^{\mathbf{D}x} = e^{\mathbf{D}x} \mathbf{D},$$

where we used the expression  $\mathbf{D} = \text{diag}[\lambda_1, \dots, \lambda_n]$ . Recalling that

$$\mathbf{P} \mathbf{D} e^{\mathbf{D}x} \mathbf{P}^{-1} = \mathbf{P} \mathbf{D} \mathbf{P}^{-1} \mathbf{P} e^{\mathbf{D}x} \mathbf{P}^{-1} = \mathbf{T} e^{\mathbf{T}x},$$

we conclude that

$$\frac{d}{dx} e^{\mathbf{T}x} = \mathbf{T} e^{\mathbf{T}x} = e^{\mathbf{T}x} \mathbf{T}.$$

This establishes the Theorem. □

**EXAMPLE 9.2.6:** Find the derivative of  $f(t) = e^{\mathbf{A}t}$ , where  $\mathbf{A} = \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix}$ .

**SOLUTION:** From Example 9.2.5 we know that

$$e^{\mathbf{A}t} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} e^{4t} & 0 \\ 0 & e^{-2t} \end{bmatrix} \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

Then, Theorem 9.2.8 implies that

$$\frac{d}{dt} e^{\mathbf{A}t} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} 4e^{4t} & 0 \\ 0 & -2e^{-2t} \end{bmatrix} \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}.$$

◁

We end this Section presenting a result without proof, that says that given any scalar-valued function with a convergent power series, that function can be extended into an operator-valued function in the case that the operator is diagonalizable.

**Theorem 9.2.9.** Let  $f : \mathbb{F} \rightarrow \mathbb{F}$  be a function given by a power series

$$f(z) = \sum_{k=0}^{\infty} c_k (z - z_0)^k,$$

which converges for  $|z - z_0| < r$ , for some positive real number  $r$ . If  $L_D(V)$  is the subspace of diagonalizable operators on the finite dimensional vector space  $V$ , then the operator-valued function  $\mathbf{F} : L_D(V) \rightarrow L_D(V)$  given by

$$\mathbf{F}(\mathbf{T}) = \sum_{k=0}^{\infty} c_k (\mathbf{T} - z_0 \mathbf{I}_V)^k$$

converges iff every eigenvalue  $\lambda_i$  of  $\mathbf{T}$  satisfies that  $|\lambda_i - z_0| < r$  for all  $i = 1, \dots, n$ .

## 9.2.4. Exercises.

9.2.1.- Which of the following matrices cannot be diagonalized?

$$A = \begin{bmatrix} 2 & -2 \\ 2 & -2 \end{bmatrix},$$

$$B = \begin{bmatrix} 2 & 0 \\ 2 & -2 \end{bmatrix},$$

$$C = \begin{bmatrix} 2 & 0 \\ 2 & 2 \end{bmatrix}.$$

9.2.2.- Verify that the matrix

$$A = \begin{bmatrix} 7/5 & 1/5 \\ -1 & 1/2 \end{bmatrix}$$

has eigenvalues  $\lambda_1 = 1$  and  $\lambda_2 = 9/10$  and associated eigenvectors

$$x_1 = \begin{bmatrix} -1 \\ 2 \end{bmatrix}, \quad x_2 = \begin{bmatrix} -2 \\ 5 \end{bmatrix}.$$

Use this information to compute

$$\lim_{k \rightarrow \infty} A^k.$$

9.2.3.- Given the matrix and vector ,

$$A = \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix}, \quad x_0 = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

compute the function  $x : \mathbb{R} \rightarrow \mathbb{R}^2$

$$x(t) = e^{At} x_0.$$

Verify that this function is solution of the differential equation

$$\frac{d}{dt}x(t) = Ax(t)$$

and satisfies that  $x(t = 0) = x_0$ .

9.2.4.- Let  $A \in \mathbb{R}^{3,3}$  be a matrix with eigenvalues 2, -1 and 3. Find the determinant of A.

9.2.5.- Let  $A \in \mathbb{R}^{4,4}$  be a matrix that can be decomposed as  $A = PDP^{-1}$ , with matrix P an invertible matrix and the matrix

$$D = \text{diag}(2, \frac{1}{4}, 2, 3).$$

Knowing only this information about the matrix A, is it possible to compute the  $\det(A)$ ? If your answer is **no**, explain why not; if your answer is **yes**, compute  $\det(A)$  and show your work.

9.2.6.- Let  $A \in \mathbb{R}^{4,4}$  be a matrix that can be decomposed as  $A = PDP^{-1}$ , with matrix P an invertible matrix and the matrix

$$D = \text{diag}(2, 0, 2, 5).$$

Knowing only this information about the matrix A, is it possible to whether A invertible? Is it possible to know  $\text{tr}(A)$ ? If your answer is **no**, explain why not; if your answer is **yes**, compute  $\text{tr}(A)$  and show your work.



## 9.3. DIFFERENTIAL EQUATIONS

Eigenvalues and eigenvectors of a matrix are useful to find solutions to systems of differential equations. In this Section we first recall what is a system of first order, linear, homogeneous, differential equations with constant coefficients. Then we use the eigenvalues and eigenvectors of the coefficient matrix to obtain solutions to such differential equations.

In order to introduce a linear, first order system of differential equations we need some notation. Let  $A : \mathbb{R} \rightarrow \mathbb{R}^{n,n}$  be a real matrix-valued function,  $\mathbf{x}, \mathbf{b} : \mathbb{R} \rightarrow \mathbb{R}^n$  be real vector-valued functions, with values  $A(t)$ ,  $\mathbf{x}(t)$ , and  $\mathbf{b}(t)$  given by

$$A(t) = \begin{bmatrix} A_{11}(t) & \cdots & A_{1n}(t) \\ \vdots & & \vdots \\ A_{n1}(t) & \cdots & A_{nn}(t) \end{bmatrix}, \quad \mathbf{x}(t) = \begin{bmatrix} x_1(t) \\ \vdots \\ x_n(t) \end{bmatrix}, \quad \mathbf{b}(t) = \begin{bmatrix} b_1(t) \\ \vdots \\ b_n(t) \end{bmatrix}.$$

So,  $A(t)$  is an  $n \times n$  matrix for each value of  $t \in \mathbb{R}$ . An example in the case  $n = 2$  is the matrix-valued function

$$A(t) = \begin{bmatrix} \cos(2\pi t) & -\sin(2\pi t) \\ \sin(2\pi t) & \cos(2\pi t) \end{bmatrix}.$$

The values of this function are rotation matrices on  $\mathbb{R}^2$ , counterclockwise by an angle  $2\pi t$ . So the bigger the parameter  $t$  the bigger the rotation. Derivatives of matrix- and vector-valued functions are computed component-wise, and we use the notation  $\dot{\phantom{x}} = \frac{d}{dt}$ ; for example

$$\frac{d\mathbf{x}}{dt}(t) = \begin{bmatrix} \frac{dx_1}{dt}(t) \\ \vdots \\ \frac{dx_n}{dt}(t) \end{bmatrix} \quad \text{is denoted as} \quad \dot{\mathbf{x}}(t) = \begin{bmatrix} \dot{x}_1(t) \\ \vdots \\ \dot{x}_n(t) \end{bmatrix}.$$

We are now ready to introduce the main definitions.

**Definition 9.3.1.** A system of first order linear differential equations on  $n$  unknowns, with  $n \geq 1$ , is the following: Given a real matrix-valued function  $A : \mathbb{R} \rightarrow \mathbb{R}^{n,n}$ , and a real vector-valued function  $\mathbf{b} : \mathbb{R} \rightarrow \mathbb{R}^n$ , find a vector-valued function  $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^n$  solution of

$$\dot{\mathbf{x}}(t) = A(t)\mathbf{x}(t) + \mathbf{b}(t). \quad (9.9)$$

The system in (9.9) is called **homogeneous** iff  $\mathbf{b}(t) = \mathbf{0}$  for all  $t \in \mathbb{R}$ . The system in (9.9) is called of **constant coefficients** iff the matrix- and vector-valued functions are constant, that is,  $A(t) = A_0$  and  $\mathbf{b}(t) = \mathbf{b}_0$  for all  $t \in \mathbb{R}$ .

The differential equation in (9.9) is called first order because it contains only first derivatives of the the unknown vector-valued function  $\mathbf{x}$ ; it is called linear because the unknown  $\mathbf{x}$  appears linearly in the equation. In this Section we are interested in finding solutions to an initial value problem involving a constant coefficient, homogeneous differential system.

**Definition 9.3.2.** An initial value problem (IVP) for an homogeneous constant coefficients linear differential equation is the following: Given a matrix  $A \in \mathbb{R}^{n,n}$  and a vector  $\mathbf{x}_0 \in \mathbb{R}^n$ , find a vector-valued function  $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^n$  solution of the differential equation and initial condition

$$\dot{\mathbf{x}}(t) = A\mathbf{x}(t), \quad \mathbf{x}(0) = \mathbf{x}_0.$$

In this Section we only consider the case where the coefficient matrix  $A \in \mathbb{R}^{n,n}$  has a complete set of eigenvectors, that is, matrix  $A$  is diagonalizable. In this case it is possible to find all solutions of the initial value problem for an homogeneous and constant coefficient differential equation. These solutions are linear combination of vectors proportional to the

eigenvectors of matrix  $A$ , where the scalars involved in the linear combination depend on the eigenvalues of matrix  $A$ . The explicit form of the solution depends on these eigenvalues and can be classified in three groups: Non-repeated real eigenvalues, non-repeated complex eigenvalues, and repeated eigenvalues. We consider in these notes only the first two cases of non-repeated eigenvalues. In this case, the main result is the following.

**Theorem 9.3.3.** *Assume that matrix  $A \in \mathbb{R}^{n,n}$  has a complete set of eigenvectors denoted as  $\mathcal{V} = \{v_1, \dots, v_n\}$  with corresponding eigenvalues  $\{\lambda_1, \dots, \lambda_n\}$ , all different, and fix  $x_0 \in \mathbb{R}^n$ . Then the initial value problem*

$$\dot{x}(t) = Ax(t), \quad x(0) = x_0, \quad (9.10)$$

has a unique solution given by

$$x(t) = e^{At} x_0, \quad (9.11)$$

where

$$e^{At} = P e^{Dt} P^{-1}, \quad P = [v_1, \dots, v_n], \quad e^{Dt} = \text{diag}[e^{\lambda_1 t}, \dots, e^{\lambda_n t}].$$

The solution given in Eq. (9.11) is often written in an equivalent way, as follows:

$$x(t) = P e^{Dt} P^{-1} x_0 = (P e^{Dt}) (P^{-1} x_0),$$

then, introducing the notation

$$P e^{Dt} = [v_1 e^{\lambda_1 t}, \dots, v_n e^{\lambda_n t}], \quad c = P^{-1} x_0, \quad c = \begin{bmatrix} c_1 \\ \vdots \\ c_n \end{bmatrix},$$

we write the solution  $x(t)$  in the form

$$x(t) = c_1 v_1 e^{\lambda_1 t} + \dots + c_n v_n e^{\lambda_n t}, \quad P c = x_0.$$

This latter notation is common in the literature on ordinary differential equations. The solution  $x(t)$  is expressed as a linear combination of the eigenvectors  $v_i$  of the coefficient matrix  $A$ , where the components are functions of the variable  $t$  given by  $c_i e^{\lambda_i t}$ , for  $i = 1, \dots, n$ . So the eigenvalues and eigenvectors of matrix  $A$  are the crucial information to find the solution  $x(t)$  of the initial value problem in Eq. (9.3.2), as can be seen from the following calculation: Consider the function

$$y_i(t) = e^{\lambda_i t} v_i, \quad i = 1, \dots, n.$$

This function is solution of the differential equation above, since

$$\dot{y}_i(t) = \frac{d}{dt}(e^{\lambda_i t}) v_i = \lambda_i e^{\lambda_i t} v_i = e^{\lambda_i t} (\lambda_i v_i) = e^{\lambda_i t} A v_i = A (e^{\lambda_i t} v_i) = A y_i(t),$$

hence,  $\dot{y}_i(t) = A y_i(t)$ . This calculation is the essential part in the proof of Theorem 9.3.3.

**Proof of Theorem 9.3.3:** Since matrix  $A$  has a complete set of eigenvectors  $\mathcal{V}$ , then for every value of  $t \in \mathbb{R}$  there exist  $t$  dependent scalars  $c_1(t), \dots, c_n(t)$  such that

$$x(t) = c_1(t) v_1 + \dots + c_n(t) v_n. \quad (9.12)$$

The  $t$ -derivative of the expression above is

$$\dot{x}(t) = \dot{c}_1(t) v_1 + \dots + \dot{c}_n(t) v_n.$$

Recalling that  $A v_i = \lambda_i v_i$ , the action of matrix  $A$  in Eq. (9.12) is

$$A x(t) = c_1(t) \lambda_1 v_1 + \dots + c_n(t) \lambda_n v_n.$$

The vector  $x(t)$  is solution of the differential equation in (9.3.2) iff  $\dot{x}(t) = A x(t)$ , that is,

$$\dot{c}_1(t) v_1 + \dots + \dot{c}_n(t) v_n = c_1(t) \lambda_1 v_1 + \dots + c_n(t) \lambda_n v_n,$$

which is equivalent to

$$[\dot{c}_1(t) - \lambda_1 c_1(t)] \mathbf{v}_1 + \cdots + [\dot{c}_n(t) - \lambda_n c_n(t)] \mathbf{v}_n = \mathbf{0}.$$

Since the set  $\mathcal{V}$  is a basis of  $\mathbb{R}^n$ , each term above must vanish, that is, for all  $i = 1, \dots, n$  holds

$$\dot{c}_i(t) = \lambda_i c_i(t) \quad \Rightarrow \quad c_i(t) = c_i e^{\lambda_i t}.$$

So we have obtained the general solution

$$\mathbf{x}(t) = c_1 e^{\lambda_1 t} \mathbf{v}_1 + \cdots + c_n e^{\lambda_n t} \mathbf{v}_n.$$

The initial condition  $\mathbf{x}(0) = \mathbf{x}_0$  fixes a unique set of constants  $c_i$  as follows

$$\mathbf{x}(0) = c_1 \mathbf{v}_1 + \cdots + c_n \mathbf{v}_n = \mathbf{x}_0,$$

since the set  $\mathcal{V}$  is a basis of  $\mathbb{R}^n$ . This expression of the solution can be rewritten as follows: Using the matrix notation  $\mathbf{P} = [\mathbf{v}_1, \dots, \mathbf{v}_n]$ , we see that the vector  $\mathbf{c}$  satisfies the equation  $\mathbf{P} \mathbf{c} = \mathbf{x}_0$ . Also notice that

$$\mathbf{P} e^{\mathbf{D}t} = [\mathbf{v}_1 e^{\lambda_1 t}, \dots, \mathbf{v}_n e^{\lambda_n t}],$$

therefore, the solution  $\mathbf{x}(t)$  can be written as

$$\mathbf{x}(t) = (\mathbf{P} e^{\mathbf{D}t}) (\mathbf{P}^{-1} \mathbf{x}_0) = (\mathbf{P} e^{\mathbf{D}t} \mathbf{P}^{-1}) \mathbf{x}_0.$$

Since  $e^{\mathbf{A}t} = \mathbf{P} e^{\mathbf{D}t} \mathbf{P}^{-1}$ , we conclude that  $\mathbf{x}(t) = e^{\mathbf{A}t} \mathbf{x}_0$ . This establishes the Theorem.  $\square$

**EXAMPLE 9.3.1: (Non-repeated, real eigenvalues)** Given the matrix  $\mathbf{A} \in \mathbb{R}^{2,2}$  and vector  $\mathbf{x}_0 \in \mathbb{R}^2$  below, find the function  $\mathbf{x}: \mathbb{R} \rightarrow \mathbb{R}^2$  solution of the initial value problem

$$\dot{\mathbf{x}}(t) = \mathbf{A} \mathbf{x}(t), \quad \mathbf{x}(0) = \mathbf{x}_0,$$

where

$$\mathbf{A} = \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix}, \quad \mathbf{x}_0 = \begin{bmatrix} 6 \\ 4 \end{bmatrix}.$$

**SOLUTION:** Recall that matrix  $\mathbf{A}$  has a complete set of eigenvectors, with

$$\mathcal{V} = \left\{ \mathbf{v}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \mathbf{v}_2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right\}, \quad \{\lambda_1 = 4, \lambda_2 = -2\}.$$

Then, Theorem 9.3.3 says that the general solution of the differential equation above is

$$\mathbf{x}(t) = c_1 e^{\lambda_1 t} \mathbf{v}_1 + c_2 e^{\lambda_2 t} \mathbf{v}_2 \quad \Leftrightarrow \quad \mathbf{x}(t) = c_1 e^{4t} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + c_2 e^{-2t} \begin{bmatrix} -1 \\ 1 \end{bmatrix}.$$

The constants  $c_1$  and  $c_2$  are obtained from the initial data  $\mathbf{x}_0$  as follows

$$\mathbf{x}(0) = c_1 \begin{bmatrix} 1 \\ 1 \end{bmatrix} + c_2 \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \mathbf{x}_0 = \begin{bmatrix} 6 \\ 4 \end{bmatrix} \quad \Rightarrow \quad \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 4 \end{bmatrix}.$$

The solution of this linear system is

$$\begin{bmatrix} c_1 \\ c_2 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 6 \\ 4 \end{bmatrix} = \begin{bmatrix} 5 \\ -1 \end{bmatrix}.$$

Therefore, the solution  $\mathbf{x}(t)$  of the initial value problem above is

$$\mathbf{x}(t) = 5 e^{4t} \begin{bmatrix} 1 \\ 1 \end{bmatrix} - e^{-2t} \begin{bmatrix} -1 \\ 1 \end{bmatrix}.$$

$\triangleleft$

**9.3.1. Non-repeated real eigenvalues.** We present a qualitative description of the solutions to Eq. (9.10) in the particular case of  $2 \times 2$  linear ordinary differential systems with matrix  $A$  having two real and different eigenvalues. The main tool will be the sketch of **phase diagrams**, also called phase portraits. The solution at a particular value  $t$  is given by a vector

$$\mathbf{x}(t) = \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix}$$

so it can be represented by a point on a plane, while the solution function for all  $t \in \mathbb{R}$  corresponds to a curve on that plane. In the case that the solution vector  $\mathbf{x}(t)$  represents a position function of a particle moving on the plane at the time  $t$ , the curve given in the phase diagram is the trajectory of the particle. Arrows are added to this trajectory to indicate the motion of the particle as time increases.

Since the eigenvalues of the coefficient matrix  $A$  are different, Theorem 9.3.3 says that there always exist two linearly independent eigenvectors  $\mathbf{v}_1$  and  $\mathbf{v}_2$  associated with the eigenvalues  $\lambda_1$  and  $\lambda_2$ , respectively. The general solution to the Eq. (9.10) is then given by

$$\mathbf{x}(t) = c_1 \mathbf{v}_1 e^{\lambda_1 t} + c_2 \mathbf{v}_2 e^{\lambda_2 t}.$$

A phase diagram contains several curves associated with several solutions, that correspond to different values of the free constants  $c_1$  and  $c_2$ . In the case that the eigenvalues are non-zero, the phase diagrams can be classified into three main classes according to the relative signs of the eigenvalues  $\lambda_1 \neq \lambda_2$  of the coefficient matrix  $A$ , as follows:

- (i)  $0 < \lambda_2 < \lambda_1$ , that is, both eigenvalues positive;
- (ii)  $\lambda_2 < 0 < \lambda_1$ , that is, one eigenvalue negative and the other positive;
- (iii)  $\lambda_2 < \lambda_1 < 0$ , that is, both eigenvalues negative.

The study of the cases where one of the eigenvalues vanishes is simpler and is left as an exercise. We now find the phase diagrams for three examples, one for each of the classes presented above. These examples summarize the behavior of the solutions to  $2 \times 2$  linear differential systems with coefficient matrix having two real, different and non-zero eigenvalues  $\lambda_2 < \lambda_1$ . The phase diagrams can be sketched following these steps: First, plot the eigenvectors  $\mathbf{v}_2$  and  $\mathbf{v}_1$  corresponding to the eigenvalues  $\lambda_2$  and  $\lambda_1$ , respectively. Second, draw the whole lines parallel to these vectors and passing through the origin. These straight lines correspond to solutions with one of the coefficients  $c_1$  or  $c_2$  vanishing. Arrows on these lines indicate how the solution changes as the variable  $t$  grows. If  $t$  is interpreted as time, the arrows indicate how the solution changes into the future. The arrows point towards the origin if the corresponding eigenvalue  $\lambda$  is negative, and they point away from the origin if the eigenvalue is positive. Finally, find the non-straight curves correspond to solutions with both coefficient  $c_1$  and  $c_2$  non-zero. Again, arrows on these curves indicate the how the solution moves into the future.

**EXAMPLE 9.3.2: (Case  $0 < \lambda_2 < \lambda_1$ .)** Sketch the phase diagram of the solutions to the differential equation

$$\dot{\mathbf{x}} = A\mathbf{x}, \quad A = \frac{1}{4} \begin{bmatrix} 11 & 3 \\ 1 & 9 \end{bmatrix}. \quad (9.13)$$

**SOLUTION:** The characteristic equation for matrix  $A$  is given by

$$\det(A - \lambda I_2) = \lambda^2 - 5\lambda + 6 = 0 \quad \Rightarrow \quad \begin{cases} \lambda_1 = 3, \\ \lambda_2 = 2. \end{cases}$$

One can show that the corresponding eigenvectors are given by

$$\mathbf{v}_1 = \begin{bmatrix} 3 \\ 1 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} -2 \\ 2 \end{bmatrix}.$$

So the general solution to the differential equation above is given by

$$x(t) = c_1 v_1 e^{\lambda_1 t} + c_2 v_2 e^{\lambda_2 t} \Leftrightarrow x(t) = c_1 \begin{bmatrix} 3 \\ 1 \end{bmatrix} e^{3t} + c_2 \begin{bmatrix} -2 \\ 2 \end{bmatrix} e^{2t}.$$

In Fig. 58 we have sketched four curves, each representing a solution  $x(t)$  corresponding to a particular choice of the constants  $c_1$  and  $c_2$ . These curves actually represent eight different solutions, for eight different choices of the constants  $c_1$  and  $c_2$ , as is described below. The arrows on these curves represent the change in the solution as the variable  $t$  grows. Since both eigenvalues are positive, the length of the solution vector always increases as  $t$  grows. The straight lines correspond to the following four solutions:

- $c_1 = 1, \quad c_2 = 0,$  Line on the first quadrant, starting at the origin, parallel to  $v_1$ ;
- $c_1 = 0, \quad c_2 = 1,$  Line on the second quadrant, starting at the origin, parallel to  $v_2$ ;
- $c_1 = -1, \quad c_2 = 0,$  Line on the third quadrant, starting at the origin, parallel to  $-v_1$ ;
- $c_1 = 0, \quad c_2 = -1,$  Line on the fourth quadrant, starting at the origin, parallel to  $-v_2$ .

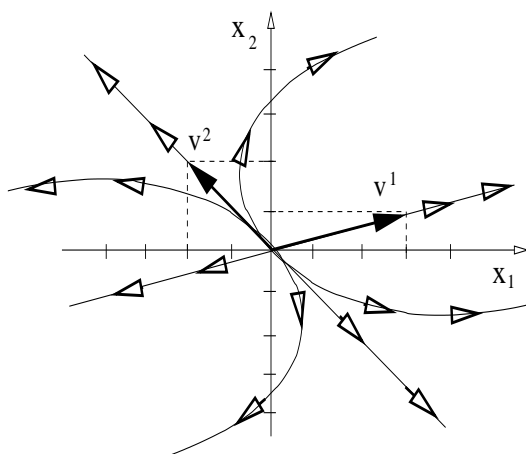


FIGURE 58. The graph of several solutions to Eq. (9.13) corresponding to the case  $0 < \lambda_2 < \lambda_1$ , for different values of the constants  $c_1$  and  $c_2$ . The trivial solution  $x = 0$  is called an unstable point.

Finally, the curved lines on each quadrant start at the origin, and they correspond to the following choices of the constants:

- $c_1 > 0, \quad c_2 > 0,$  Line starting on the second to the first quadrant;
- $c_1 < 0, \quad c_2 > 0,$  Line starting on the second to the third quadrant;
- $c_1 < 0, \quad c_2 < 0,$  Line starting on the fourth to the third quadrant,
- $c_1 > 0, \quad c_2 < 0,$  Line starting on the fourth to the first quadrant.

◁

**EXAMPLE 9.3.3:** (Case  $\lambda_2 < 0 < \lambda_1$ .) Sketch the phase diagram of the solutions to the differential equation

$$\dot{x} = Ax, \quad A = \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix}. \quad (9.14)$$

**SOLUTION:** We known from the calculations performed in Example 9.3.1 that the general solution to the differential equation above is given by

$$\mathbf{x}(t) = c_1 \mathbf{v}_1 e^{\lambda_1 t} + c_2 \mathbf{v}_2 e^{\lambda_2 t} \Leftrightarrow \mathbf{x}(t) = c_1 \begin{bmatrix} 1 \\ 1 \end{bmatrix} e^{4t} + c_2 \begin{bmatrix} -1 \\ 1 \end{bmatrix} e^{-2t},$$

where we have introduced the eigenvalues and eigenvectors

$$\lambda_1 = 4, \quad \mathbf{v}_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \text{and} \quad \lambda_2 = -2, \quad \mathbf{v}_2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}.$$

In Fig. 59 we have sketched four curves, each representing a solution  $\mathbf{x}(t)$  corresponding to a particular choice of the constants  $c_1$  and  $c_2$ . These curves actually represent eight different solutions, for eight different choices of the constants  $c_1$  and  $c_2$ , as is described below. The arrows on these curves represent the change in the solution as the variable  $t$  grows. The part of the solution with positive eigenvalue increases exponentially when  $t$  grows, while the part of the solution with negative eigenvalue decreases exponentially when  $t$  grows. The straight lines correspond to the following four solutions:

- $c_1 = 1, \quad c_2 = 0,$  Line on the first quadrant, starting at the origin, parallel to  $\mathbf{v}_1$ ;
- $c_1 = 0, \quad c_2 = 1,$  Line on the second quadrant, ending at the origin, parallel to  $\mathbf{v}_2$ ;
- $c_1 = -1, \quad c_2 = 0,$  Line on the third quadrant, starting at the origin, parallel to  $-\mathbf{v}_1$ ;
- $c_1 = 0, \quad c_2 = -1,$  Line on the fourth quadrant, ending at the origin, parallel to  $-\mathbf{v}_2$ .

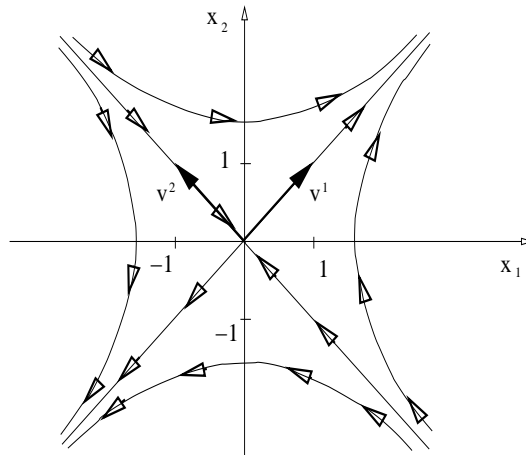


FIGURE 59. The graph of several solutions to Eq. (9.14) corresponding to the case  $\lambda_2 < 0 < \lambda_1$ , for different values of the constants  $c_1$  and  $c_2$ . The trivial solution  $\mathbf{x} = 0$  is called a saddle point.

Finally, the curved lines on each quadrant correspond to the following choices of the constants:

- $c_1 > 0, \quad c_2 > 0,$  Line from the second to the first quadrant,
- $c_1 < 0, \quad c_2 > 0,$  Line from the second to the third quadrant,
- $c_1 < 0, \quad c_2 < 0,$  Line from the fourth to the third quadrant,
- $c_1 > 0, \quad c_2 < 0,$  Line from the fourth to the first quadrant.

◁

**EXAMPLE 9.3.4:** (Case  $\lambda_2 < \lambda_1 < 0$ .) Sketch the phase diagram of the solutions to the differential equation

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x}, \quad \mathbf{A} = \frac{1}{4} \begin{bmatrix} -9 & 3 \\ 1 & -11 \end{bmatrix}. \quad (9.15)$$

**SOLUTION:** The characteristic equation for this matrix  $\mathbf{A}$  is given by

$$\det(\mathbf{A} - \lambda \mathbf{I}) = \lambda^2 + 5\lambda + 6 = 0 \quad \Rightarrow \quad \begin{cases} \lambda_1 = -2, \\ \lambda_2 = -3. \end{cases}$$

One can show that the corresponding eigenvectors are given by

$$\mathbf{v}_1 = \begin{bmatrix} 3 \\ 1 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} -2 \\ 2 \end{bmatrix}.$$

So the general solution to the differential equation above is given by

$$\mathbf{x}(t) = c_1 \mathbf{v}_1 e^{\lambda_1 t} + c_2 \mathbf{v}_2 e^{\lambda_2 t} \quad \Leftrightarrow \quad \mathbf{x}(t) = c_1 \begin{bmatrix} 3 \\ 1 \end{bmatrix} e^{-2t} + c_2 \begin{bmatrix} -2 \\ 2 \end{bmatrix} e^{-3t}.$$

In Fig. 60 we have sketched four curves, each representing a solution  $\mathbf{x}(t)$  corresponding to a particular choice of the constants  $c_1$  and  $c_2$ . These curves actually represent eight different solutions, for eight different choices of the constants  $c_1$  and  $c_2$ , as is described below. The arrows on these curves represent the change in the solution as the variable  $t$  grows. Since both eigenvalues are negative, the length of the solution vector always decreases as  $t$  grows and the solution vector always approaches zero. The straight lines correspond to the following four solutions:

- $c_1 = 1, \quad c_2 = 0,$  Line on the first quadrant, ending at the origin, parallel to  $\mathbf{v}_1$ ;
- $c_1 = 0, \quad c_2 = 1,$  Line on the second quadrant, ending at the origin, parallel to  $\mathbf{v}_2$ ;
- $c_1 = -1, \quad c_2 = 0,$  Line on the third quadrant, ending at the origin, parallel to  $-\mathbf{v}_1$ ;
- $c_1 = 0, \quad c_2 = -1,$  Line on the fourth quadrant, ending at the origin, parallel to  $-\mathbf{v}_2$ .

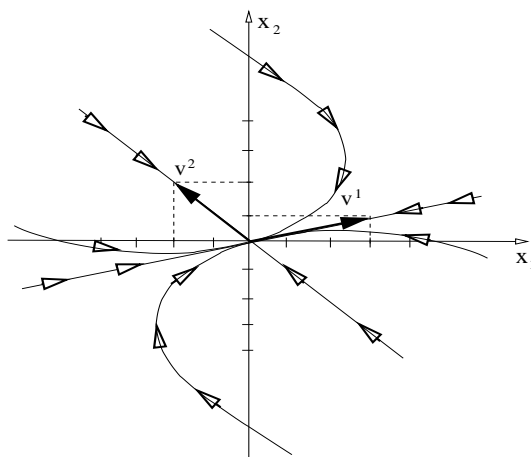


FIGURE 60. The graph of several solutions to Eq. (9.15) corresponding to the case  $\lambda_2 < \lambda_1 < 0$ , for different values of the constants  $c_1$  and  $c_2$ . The trivial solution  $\mathbf{x} = \mathbf{0}$  is called a stable point.

Finally, the curved lines on each quadrant start at the origin, and they correspond to the following choices of the constants:

$c_1 > 0,$	$c_2 > 0,$	Line entering the first from the second quadrant;
$c_1 < 0,$	$c_2 > 0,$	Line entering the third from the second quadrant;
$c_1 < 0,$	$c_2 < 0,$	Line entering the third from the fourth quadrant,
$c_1 > 0,$	$c_2 < 0,$	Line entering the first from the fourth quadrant.

◀

**9.3.2. Non-repeated complex eigenvalues.** The complex eigenvalues of a real valued matrix  $A \in \mathbb{R}^{n,n}$  are always complex conjugate pairs, as it is shown below.

**Lemma 9.3.4. (Conjugate pairs)** *If a real valued matrix  $A \in \mathbb{R}^n$ , has a complex eigenvalue  $\lambda$  with eigenvector  $\mathbf{v}$ , then  $\bar{\lambda}$  and  $\bar{\mathbf{v}}$  are also an eigenvalue and eigenvector of matrix  $A$ .*

**Proof of Lemma 9.3.4:** Complex conjugate the eigenvalue-eigenvector equation for  $\lambda$  and  $\mathbf{v}$  and recalling that  $A = \bar{A}$ , we obtain

$$A\mathbf{v} = \lambda\mathbf{v} \quad \Leftrightarrow \quad A\bar{\mathbf{v}} = \bar{\lambda}\bar{\mathbf{v}}.$$

□

Since the complex eigenvalues of a matrix with real coefficients are always complex conjugate pairs, there is an even number of complex eigenvalues. Denoting the eigenvalue pair by  $\lambda_{\pm}$  and the corresponding eigenvector pair by  $\mathbf{v}_{\pm}$ , it holds that  $\lambda_+ = \bar{\lambda}_-$  and  $\mathbf{v}_+ = \bar{\mathbf{v}}_-$ . Hence, an eigenvalue and eigenvector pairs have the form

$$\lambda_{\pm} = \alpha \pm i\beta, \quad \mathbf{v}_{\pm} = \mathbf{a} \pm i\mathbf{b}, \quad (9.16)$$

where  $\alpha, \beta \in \mathbb{R}$  and  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ . It is simple to obtain two linearly independent solutions to the differential equation in Eq. (9.10) in the case that matrix  $A$  has a complex conjugate pair of eigenvalues and eigenvectors. These solutions can be expressed both as complex-valued or as real-valued functions.

**Theorem 9.3.5. (Conjugate pairs)** *Let  $\lambda_{\pm} = \alpha \pm i\beta$  be eigenvalues of a matrix  $A \in \mathbb{R}^{n,n}$  with respective eigenvectors  $\mathbf{v}_{\pm} = \mathbf{a} \pm i\mathbf{b}$ , where  $\alpha, \beta \in \mathbb{R}$ , while  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ , and  $n \geq 2$ . Then a linearly independent set of complex valued solutions to the differential equation in (9.10) is formed by the functions*

$$\mathbf{x}_+ = \mathbf{v}_+ e^{\lambda_+ t}, \quad \mathbf{x}_- = \mathbf{v}_- e^{\lambda_- t}, \quad (9.17)$$

while a linearly independent set of real valued solutions to Eq. (9.10) is given by the functions

$$\mathbf{x}_1 = [\mathbf{a} \cos(\beta t) - \mathbf{b} \sin(\beta t)] e^{\alpha t}, \quad \mathbf{x}_2 = [\mathbf{a} \sin(\beta t) + \mathbf{b} \cos(\beta t)] e^{\alpha t}. \quad (9.18)$$

**Proof of Theorem 9.3.5:** We know from Theorem 9.3.3 that two linearly independent solutions to Eq. (9.10) are given by Eq. (9.17). The new information in Theorem 9.3.5 above is the real-valued solutions in Eq. (9.18). They can be obtained from Eq. (9.17) as follows:

$$\begin{aligned} \mathbf{x}_{\pm} &= (\mathbf{a} \pm i\mathbf{b}) e^{(\alpha \pm i\beta)t} \\ &= e^{\alpha t} (\mathbf{a} \pm i\mathbf{b}) e^{\pm i\beta t} \\ &= e^{\alpha t} (\mathbf{a} \pm i\mathbf{b}) [\cos(\beta t) \pm i \sin(\beta t)] \\ &= e^{\alpha t} [\mathbf{a} \cos(\beta t) - \mathbf{b} \sin(\beta t)] \pm i e^{\alpha t} [\mathbf{a} \sin(\beta t) + \mathbf{b} \cos(\beta t)]. \end{aligned}$$



Since the differential equation in (9.10) is linear, the functions below are also solutions,

$$\begin{aligned} \mathbf{x}_1 &= \frac{1}{2}(\mathbf{x}_+ + \mathbf{x}_-) = e^{\alpha t} [\mathbf{a} \cos(\beta t) - \mathbf{b} \sin(\beta t)], \\ \mathbf{x}_2 &= \frac{1}{2i}(\mathbf{x}_+ - \mathbf{x}_-) = e^{\alpha t} [\mathbf{a} \sin(\beta t) + \mathbf{b} \cos(\beta t)]. \end{aligned}$$

This establishes the Theorem.  $\square$

**EXAMPLE 9.3.5:** Find a real-valued set of fundamental solutions to the differential equation

$$\dot{\mathbf{x}} = \mathbf{A} \mathbf{x}, \quad \mathbf{A} = \begin{bmatrix} 2 & 3 \\ -3 & 2 \end{bmatrix}, \quad (9.19)$$

and then sketch a phase diagram for the solutions of this equation.

**SOLUTION:** First find the eigenvalues of matrix  $\mathbf{A}$  above,

$$0 = \begin{vmatrix} (2 - \lambda) & 3 \\ -3 & (2 - \lambda) \end{vmatrix} = (\lambda - 2)^2 + 9 \Rightarrow \lambda_{\pm} = 2 \pm 3i.$$

We then find the respective eigenvectors. The one corresponding to  $\lambda_+$  is the solution of the homogeneous linear system with coefficients given by

$$\begin{bmatrix} 2 - (2 + 3i) & 3 \\ -3 & 2 - (2 + 3i) \end{bmatrix} = \begin{bmatrix} -3i & 3 \\ -3 & -3i \end{bmatrix} \rightarrow \begin{bmatrix} -i & 1 \\ -1 & -i \end{bmatrix} \rightarrow \begin{bmatrix} 1 & i \\ -1 & -i \end{bmatrix} \rightarrow \begin{bmatrix} 1 & i \\ 0 & 0 \end{bmatrix}.$$

Therefore the eigenvector  $\mathbf{v}_+ = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$  is given by

$$v_1 = -iv_2 \Rightarrow v_2 = 1, \quad v_1 = -i, \quad \Rightarrow \mathbf{v}_+ = \begin{bmatrix} -i \\ 1 \end{bmatrix}, \quad \lambda_+ = 2 + 3i.$$

The second eigenvector is the complex conjugate of the eigenvector found above, that is,

$$\mathbf{v}_- = \begin{bmatrix} i \\ 1 \end{bmatrix}, \quad \lambda_- = 2 - 3i.$$

Notice that

$$\mathbf{v}_{\pm} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \pm \begin{bmatrix} -1 \\ 0 \end{bmatrix} i.$$

Hence, the real and imaginary parts of the eigenvalues and of the eigenvectors are given by

$$\alpha = 2, \quad \beta = 3, \quad \mathbf{a} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}.$$

So a real-valued expression for a fundamental set of solutions is given by

$$\begin{aligned} \mathbf{x}_1 &= \left( \begin{bmatrix} 0 \\ 1 \end{bmatrix} \cos(3t) - \begin{bmatrix} -1 \\ 0 \end{bmatrix} \sin(3t) \right) e^{2t} \Rightarrow \mathbf{x}_1 = \begin{bmatrix} \sin(3t) \\ \cos(3t) \end{bmatrix} e^{2t}, \\ \mathbf{x}_2 &= \left( \begin{bmatrix} 0 \\ 1 \end{bmatrix} \sin(3t) + \begin{bmatrix} -1 \\ 0 \end{bmatrix} \cos(3t) \right) e^{2t} \Rightarrow \mathbf{x}_2 = \begin{bmatrix} -\cos(3t) \\ \sin(3t) \end{bmatrix} e^{2t}. \end{aligned}$$

The phase diagram of these two fundamental solutions is given in Fig. 61 below. There is also a circle given in that diagram, corresponding to the trajectory of the vectors

$$\tilde{\mathbf{x}}_1 = \begin{bmatrix} \sin(3t) \\ \cos(3t) \end{bmatrix}, \quad \tilde{\mathbf{x}}_2 = \begin{bmatrix} -\cos(3t) \\ \sin(3t) \end{bmatrix}.$$

The trajectory of these vectors is a circle since their length is constant equal to one.

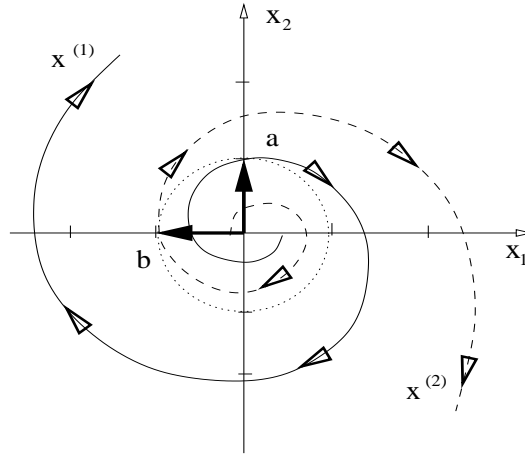


FIGURE 61. The graph of the fundamental solutions  $x_1$  and  $x_2$  of the Eq. (9.19).

In the particular case that the matrix  $A$  in Eq. (9.10) is  $2 \times 2$ , then any solutions is a linear combination of the solutions given in Eq. (9.18). That is, the general solution of given by

$$x(t) = c_1 x_1(t) + c_2 x_2(t),$$

where

$$x_1 = [a \cos(\beta t) - b \sin(\beta t)] e^{\alpha t}, \quad x_2 = [a \sin(\beta t) + b \cos(\beta t)] e^{\alpha t}.$$

We now do a qualitative study of the phase diagrams of the solutions for this case. We first fix the vectors  $a$  and  $b$ , and the plot phase diagrams for solutions having  $\alpha > 0$ ,  $\alpha = 0$ , and  $\alpha < 0$ . These diagrams are given in Fig. 62. One can see that for  $\alpha > 0$  the solutions spiral outward as  $t$  increases, and for  $\alpha < 0$  the solutions spiral inwards to the origin as  $t$  increases..

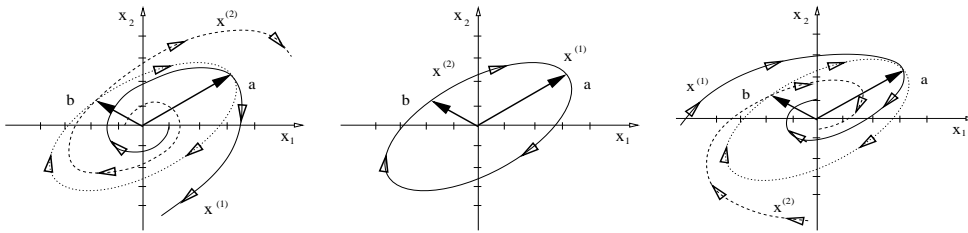


FIGURE 62. The graph of the fundamental solutions  $x_1$  and  $x_2$  (dashed line) of the Eq. (9.18) in the case of  $\alpha > 0$ ,  $\alpha = 0$ , and  $\alpha < 0$ , respectively.

Finally, let us study the following cases: We fix  $\alpha > 0$  and the vector  $a$ , and we plot the phase diagrams for solutions with two choices of the vector  $b$ , as shown in Fig. 63. It can then be seen that the relative directions of the vectors  $a$  and  $b$  determines the rotation direction of the solutions as  $t$  increases.

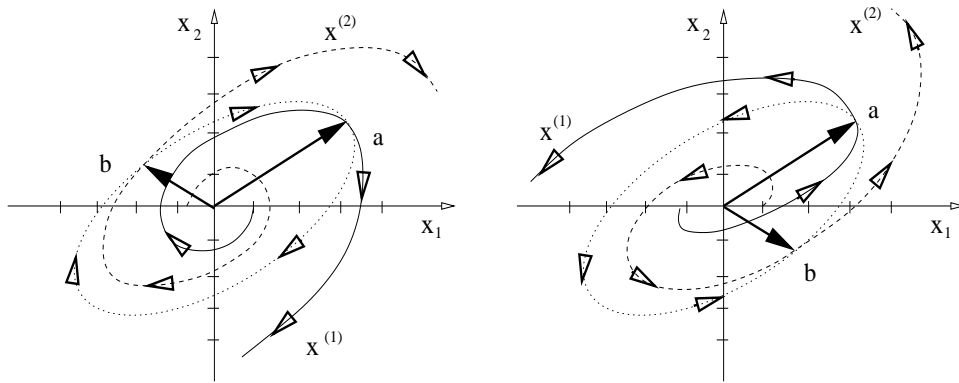


FIGURE 63. The graph of the fundamental solutions  $x_1$  and  $x_2$  (dashed line) of the Eq. (9.18) in the case of  $\alpha > 0$ , for a given vector  $a$  and for two choices of the vector  $b$ . The relative positions of the vectors  $a$  and  $b$  determines the rotation direction.

## 9.3.3. Exercises.

**9.3.1.-** Given the matrix  $A \in \mathbb{R}^{2,2}$  and vector  $x_0 \in \mathbb{R}^2$  below, find the function  $x : \mathbb{R} \rightarrow \mathbb{R}^2$  solution of the initial value problem

$$\dot{x}(t) = Ax(t), \quad x(0) = x_0,$$

where

$$A = \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix}, \quad x_0 = \begin{bmatrix} 3 \\ 4 \end{bmatrix}.$$

**9.3.2.-** Given the matrix  $A$  and the vector  $x_0$  in Exercise **9.3.1** above, compute the operator valued function  $e^{At}$  and verify that the solution  $x$  of the initial value problem given in Exercise **9.3.1** can be written as

$$x(t) = e^{At} x_0.$$

**9.3.3.-** Given the matrix  $A \in \mathbb{R}^{2,2}$  below, find all functions  $x : \mathbb{R} \rightarrow \mathbb{R}^2$  solutions of the differential equation

$$\dot{x}(t) = Ax(t),$$

where

$$A = \begin{bmatrix} 1 & -1 \\ 5 & -3 \end{bmatrix}.$$

Since this matrix has complex eigenvalues, express the solutions as linear combination of real vector valued functions.

**9.3.4.-** Given the matrix  $A \in \mathbb{R}^{3,3}$  and vector  $x_0 \in \mathbb{R}^3$  below, find the function  $x : \mathbb{R} \rightarrow \mathbb{R}^3$  solution of the initial value problem

$$\dot{x}(t) = Ax(t), \quad x(0) = x_0,$$

where

$$A = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 2 & 2 \\ 0 & 0 & 1 \end{bmatrix}, \quad x_0 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}.$$

## 9.4. NORMAL OPERATORS

In this Section we introduce a particular type of linear operators called normal operators, which are defined on vector spaces having an inner product. Particular cases are Hermitian operators and unitary operators, which are widely used in physics. Rotations in space are examples of unitary operators on real vector spaces, while physical observables in quantum mechanics are examples of Hermitian operators. In this Section we restrict our description to finite dimensional inner product spaces. These definitions generalize the notions of unitary and Hermitian matrices already introduced in Chapter 2. We first describe the Riesz Representation Theorem, needed to verify that the notion of adjoint of a linear operator is well-defined. After reviewing the commutator of two operators we then introduce normal operators and discuss the particular cases of unitary and Hermitian operators. Finally we comment on the relations between these notions and the unitary and Hermitian matrices already introduced in Chapter 2.

The Riesz Representation Theorem is a statement concerning linear functionals on an inner product space. Given a vector space  $V$  over the scalar field  $\mathbb{F}$ , a **linear functional** is a scalar-valued linear function  $f : V \rightarrow \mathbb{F}$ , that is, for all  $\mathbf{x}, \mathbf{y} \in V$  and all  $a, b \in \mathbb{F}$  holds  $f(a\mathbf{x} + b\mathbf{y}) = af(\mathbf{x}) + bf(\mathbf{y}) \in \mathbb{F}$ . An example of a linear functional on  $\mathbb{R}^3$  is the function

$$\mathbb{R}^3 \ni \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \mapsto f(\mathbf{x}) = 3x_1 + 2x_2 + x_3 \in \mathbb{R}.$$

This function can be expressed in terms of the dot product in  $\mathbb{R}^3$  as follows

$$f(\mathbf{x}) = \mathbf{u} \cdot \mathbf{x}, \quad \mathbf{u} = \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix}.$$

The Riesz Representation Theorem says that what we did in this example can be done in the general case. In an inner product space  $(V, \langle \cdot, \cdot \rangle)$  every linear functional  $f$  can be expressed in terms of the inner product.

**Theorem 9.4.1.** *Consider a finite dimensional inner product space  $(V, \langle \cdot, \cdot \rangle)$  over the scalar field  $\mathbb{F}$ . For every linear functional  $f : V \rightarrow \mathbb{F}$  there exists a unique vector  $\mathbf{u}_f \in V$  such that holds*

$$f(\mathbf{v}) = \langle \mathbf{u}_f, \mathbf{v} \rangle \quad \forall \mathbf{v} \in V.$$

**Proof of Theorem 9.4.1:** Introduce the set

$$N = \{ \mathbf{v} \in V : f(\mathbf{v}) = 0 \} \subset V.$$

This set is the analogous to linear functionals of the null space of linear operators. Since  $f$  is a linear function the set  $N$  is a subspace of  $V$ . (Proof: Given two elements  $\mathbf{v}_1, \mathbf{v}_2 \in N$  and two scalars  $a, b \in \mathbb{F}$ , holds that  $f(a\mathbf{v}_1 + b\mathbf{v}_2) = af(\mathbf{v}_1) + bf(\mathbf{v}_2) = 0 + 0$ , so  $(a\mathbf{v}_1 + b\mathbf{v}_2) \in N$ .) Introduce the orthogonal complement of  $N$ , that is,

$$N^\perp = \{ \mathbf{w} \in V : \langle \mathbf{w}, \mathbf{v} \rangle = 0 \forall \mathbf{v} \in N \},$$

which is also a subspace of  $V$ . If  $N^\perp = \{\mathbf{0}\}$ , then  $N = (N^\perp)^\perp = (\{\mathbf{0}\})^\perp = V$ . Since the null space of  $f$  is the whole vector space, the functional  $f$  is identically zero, so only for the choice  $\mathbf{u}_f = \mathbf{0}$  holds  $f(\mathbf{v}) = \langle \mathbf{0}, \mathbf{v} \rangle$  for all  $\mathbf{v} \in V$ .

In the case that  $N^\perp \neq \{\mathbf{0}\}$  we now show that this space cannot be very big, in fact it has dimension one, as the following argument shows. Choose  $\tilde{\mathbf{u}} \in N^\perp$  such that  $f(\tilde{\mathbf{u}}) = 1$ .

Then notice that for every  $\mathbf{w} \in N^\perp$  the vector  $\mathbf{w} - f(\mathbf{w})\tilde{\mathbf{u}}$  is trivially in  $N^\perp$  but it is also in  $N$ , since

$$f(\mathbf{w} - f(\mathbf{w})\tilde{\mathbf{u}}) = f(\mathbf{w}) - f(\mathbf{w})f(\tilde{\mathbf{u}}) = f(\mathbf{w}) - f(\mathbf{w}) = 0.$$

A vector both in  $N$  and  $N^\perp$  must vanish, so  $\mathbf{w} = f(\mathbf{w})\tilde{\mathbf{u}}$ . Then every vector in  $N^\perp$  is proportional to  $\tilde{\mathbf{u}}$ , so  $\dim N^\perp = 1$ . This information is used to split any vector  $\mathbf{v} \in V$  as follows  $\mathbf{v} = a\tilde{\mathbf{u}} + \mathbf{x}$  where  $\mathbf{x} \in V$  and  $a \in \mathbb{F}$ . It is clear that

$$f(\mathbf{v}) = f(a\tilde{\mathbf{u}} + \mathbf{x}) = af(\tilde{\mathbf{u}}) + f(\mathbf{x}) = af(\tilde{\mathbf{u}}) = a.$$

However, the function with values  $g(\mathbf{v}) = \left\langle \frac{\tilde{\mathbf{u}}}{\|\tilde{\mathbf{u}}\|^2}, \mathbf{v} \right\rangle$  has precisely the same values as  $f$ , since for all  $\mathbf{v} \in V$  holds

$$g(\mathbf{v}) = \left\langle \frac{\tilde{\mathbf{u}}}{\|\tilde{\mathbf{u}}\|^2}, \mathbf{v} \right\rangle = \left\langle \frac{\tilde{\mathbf{u}}}{\|\tilde{\mathbf{u}}\|^2}, (a\tilde{\mathbf{u}} + \mathbf{x}) \right\rangle = \frac{a}{\|\tilde{\mathbf{u}}\|^2} \langle \tilde{\mathbf{u}}, \tilde{\mathbf{u}} \rangle + \frac{1}{\|\tilde{\mathbf{u}}\|^2} \langle \tilde{\mathbf{u}}, \mathbf{x} \rangle = a.$$

Therefore, choosing  $\mathbf{u}_f = \tilde{\mathbf{u}}/\|\tilde{\mathbf{u}}\|^2$ , holds that

$$f(\mathbf{v}) = \langle \mathbf{u}_f, \mathbf{v} \rangle \quad \forall \mathbf{v} \in V.$$

Since  $\dim N^\perp = 1$ , the choice of  $\mathbf{u}_f$  is unique. This establishes the Theorem.  $\square$

Given a linear operator defined on an inner product space, a new linear operator can be defined through an equation involving the inner product.

**Proposition 9.4.2.** *Let  $\mathbf{T} \in L(V)$  be a linear operator on a finite-dimensional inner product space  $(V, \langle \cdot, \cdot \rangle)$ . There exists one and only one linear operator  $\mathbf{T}^* \in L(V)$  such that*

$$\langle \mathbf{v}, \mathbf{T}^*(\mathbf{u}) \rangle = \langle \mathbf{T}(\mathbf{v}), \mathbf{u} \rangle$$

holds for all vectors  $\mathbf{u}, \mathbf{v} \in V$ .

Given any linear operator  $\mathbf{T}$  on a finite-dimensional inner product space, the operator  $\mathbf{T}^*$  whose existence is guaranteed in Proposition 9.4.2 is called the **adjoint** of  $\mathbf{T}$ .

**Proof of Proposition 9.4.2:** We first establish the following statement: For every vector  $\mathbf{u} \in V$  there exists a unique vector  $\mathbf{w} \in V$  such that

$$\langle \mathbf{T}(\mathbf{v}), \mathbf{u} \rangle = \langle \mathbf{v}, \mathbf{w} \rangle \quad \forall \mathbf{v} \in V. \quad (9.20)$$

The proof starts noticing that for a fixed  $\mathbf{u} \in V$  the scalar-valued function  $f_u : V \rightarrow \mathbb{F}$  given by  $f_u(\mathbf{v}) = \langle \mathbf{u}, \mathbf{T}(\mathbf{v}) \rangle$  is a linear functional. Therefore, the Riesz Representation Theorem 9.4.1 implies that there exists a unique vector  $\mathbf{w} \in V$  such that  $f_u(\mathbf{v}) = \langle \mathbf{w}, \mathbf{v} \rangle$ . This establishes that for every vector  $\mathbf{u} \in V$  there exists a unique vector  $\mathbf{w} \in V$  such that Eq. (9.20) holds. Now that this statement is proven we can define a map, that we choose to denote as  $\mathbf{T}^* : V \rightarrow V$ , given by  $\mathbf{u} \mapsto \mathbf{T}^*(\mathbf{u}) = \mathbf{w}$ . We now show that this map  $\mathbf{T}^*$  is linear. Indeed, for all  $\mathbf{u}_1, \mathbf{u}_2 \in V$  and all  $a, b \in \mathbb{F}$  holds

$$\begin{aligned} \langle \mathbf{v}, \mathbf{T}^*(a\mathbf{u}_1 + b\mathbf{u}_2) \rangle &= \langle \mathbf{T}(\mathbf{v}), (a\mathbf{u}_1 + b\mathbf{u}_2) \rangle \quad \forall \mathbf{v} \in V, \\ &= a \langle \mathbf{T}(\mathbf{v}), \mathbf{u}_1 \rangle + b \langle \mathbf{T}(\mathbf{v}), \mathbf{u}_2 \rangle \\ &= a \langle \mathbf{v}, \mathbf{T}^*(\mathbf{u}_1) \rangle + b \langle \mathbf{v}, \mathbf{T}^*(\mathbf{u}_2) \rangle \\ &= \langle \mathbf{v}, [a\mathbf{T}^*(\mathbf{u}_1) + b\mathbf{T}^*(\mathbf{u}_2)] \rangle \quad \forall \mathbf{v} \in V, \end{aligned}$$

hence  $\mathbf{T}^*(a\mathbf{u}_1 + b\mathbf{u}_2) = a\mathbf{T}^*(\mathbf{u}_1) + b\mathbf{T}^*(\mathbf{u}_2)$ . This establishes the Proposition.  $\square$

The next result relates the adjoint of a linear operator with the concept of the adjoint of a square matrix introduced in Sect. 2.2. Recall that given a basis in the vector space, every linear operator has associated a unique square matrix. Let us use the notation  $[\mathbf{T}]$  and  $[\mathbf{T}^*]$  for the matrices on a given basis of the operators  $\mathbf{T}$  and  $\mathbf{T}^*$ , respectively.

**Proposition 9.4.3.** *Let  $(V, \langle \cdot, \cdot \rangle)$  be a finite-dimensional vector space, let  $\mathcal{V}$  be an orthonormal basis of  $V$ , and let  $[\mathbf{T}]$  be the matrix of the linear operator  $\mathbf{T} \in L(V)$  in the basis  $\mathcal{V}$ . Then, the matrix of the adjoint operator  $\mathbf{T}^*$  in the basis  $\mathcal{V}$  is given by  $[\mathbf{T}^*] = [\mathbf{T}]^*$ .*

Proposition 9.4.3 says that the matrix of the adjoint operator is the adjoint of the matrix of the operator, however this is true only in the case that the basis used to compute the respective matrices is orthonormal. If the basis is not orthonormal, the relation between the matrices  $[\mathbf{T}]$  and  $[\mathbf{T}^*]$  is more involved.

**Proof of Proposition 9.4.3:** Let  $\mathcal{V} = \{\mathbf{e}_1, \dots, \mathbf{e}_n\}$  be an orthonormal basis of  $V$ , that is,

$$\langle \mathbf{e}_i, \mathbf{e}_j \rangle = \begin{cases} 0 & \text{if } i \neq j, \\ 1 & \text{if } i = j. \end{cases}$$

The components of two arbitrary vectors  $\mathbf{u}, \mathbf{v} \in V$  in the basis  $\mathcal{V}$  is denoted as follows

$$\mathbf{u} = \sum_i u_i \mathbf{e}_i, \quad \mathbf{v} = \sum_i v_i \mathbf{e}_i.$$

The action of the operator  $\mathbf{T}$  can also be decomposed in the basis  $\mathcal{V}$  as follows

$$\mathbf{T}(\mathbf{e}_j) = \sum_i [\mathbf{T}]_{ij} \mathbf{e}_i, \quad [\mathbf{T}]_{ij} = [\mathbf{T}(\mathbf{e}_j)]_i.$$

We use the same notation for the adjoint operator, that is,

$$\mathbf{T}^*(\mathbf{e}_j) = \sum_i [\mathbf{T}^*]_{ij} \mathbf{e}_i, \quad [\mathbf{T}^*]_{ij} = [\mathbf{T}^*(\mathbf{e}_j)]_i.$$

The adjoint operator is defined such that the equation  $\langle \mathbf{v}, \mathbf{T}^*(\mathbf{u}) \rangle = \langle \mathbf{T}(\mathbf{v}), \mathbf{u} \rangle$  holds for all  $\mathbf{u}, \mathbf{v} \in V$ . This equation can be expressed in terms of components in the basis  $\mathcal{V}$  as follows

$$\sum_{ijk} \langle v_i \mathbf{e}_i, u_j [\mathbf{T}^*(\mathbf{e}_j)]_k \mathbf{e}_k \rangle = \sum_{ijk} \langle v_i [\mathbf{T}(\mathbf{e}_i)]_k \mathbf{e}_k, u_j \mathbf{e}_j \rangle,$$

that is,

$$\sum_{ijk} \bar{v}_i u_j [\mathbf{T}^*]_{kj} \langle \mathbf{e}_i, \mathbf{e}_k \rangle = \sum_{ijk} \bar{v}_i [\mathbf{T}]_{ki} u_j \langle \mathbf{e}_k, \mathbf{e}_j \rangle.$$

Since the basis  $\mathcal{V}$  is orthonormal we obtain the equation

$$\sum_{ij} \bar{v}_i u_j [\mathbf{T}^*]_{ij} = \sum_{ijk} \bar{v}_i [\mathbf{T}]_{ji} u_j,$$

which holds for all vectors  $\mathbf{u}, \mathbf{v} \in V$ , so we conclude

$$[\mathbf{T}^*]_{ij} = \overline{[\mathbf{T}]_{ji}} \Leftrightarrow [\mathbf{T}^*] = \overline{[\mathbf{T}]^T} \Leftrightarrow [\mathbf{T}^*] = [\mathbf{T}]^*.$$

This establishes the Proposition.  $\square$

**EXAMPLE 9.4.1:** Consider the inner product space  $(\mathbb{C}^3, \cdot)$ . Find the adjoint of the linear operator  $\mathbf{T}$  with matrix in the standard basis of  $\mathbb{C}^3$  given by

$$[\mathbf{T}(\mathbf{x})] = \begin{bmatrix} x_1 + 2ix_2 + ix_3 \\ ix_1 - x_3 \\ x_1 - x_2 + ix_3 \end{bmatrix}, \quad [\mathbf{x}] = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}.$$

**SOLUTION:** The matrix of this operator in the standard basis of  $\mathbb{C}^3$  is given by

$$[\mathbf{T}] = \begin{bmatrix} 1 & 2i & i \\ i & 0 & -1 \\ 1 & -1 & i \end{bmatrix}.$$

Since the standard basis is an orthonormal basis with respect to the dot product, Proposition 9.4.3 implies that

$$[\mathbf{T}^*] = [\mathbf{T}]^* = \begin{bmatrix} 1 & 2i & i \\ i & 0 & -1 \\ 1 & -1 & i \end{bmatrix}^* = \begin{bmatrix} 1 & -i & 1 \\ -2i & 0 & -1 \\ -i & -1 & -i \end{bmatrix} \Rightarrow [\mathbf{T}^*(\mathbf{x})] = \begin{bmatrix} x_1 - ix_2 + x_3 \\ -2ix_1 - x_3 \\ -ix_1 - x_2 - ix_3 \end{bmatrix}.$$

◁

Recall now that the commutator of two linear operators  $\mathbf{T}, \mathbf{S} \in L(V)$  is the linear operator  $[\mathbf{T}, \mathbf{S}] \in L(V)$  given by

$$[\mathbf{T}, \mathbf{S}](\mathbf{u}) = \mathbf{T}(\mathbf{S}(\mathbf{u})) - \mathbf{S}(\mathbf{T}(\mathbf{u})) \quad \forall \mathbf{u} \in V.$$

Two operators  $\mathbf{T}, \mathbf{S} \in L(V)$  are said to commute iff their commutator vanishes, that is,  $[\mathbf{T}, \mathbf{S}] = \mathbf{0}$ . Examples of operators that commute are two rotations on the plane. Examples of operators that do not commute are two arbitrary rotations in space.

**Definition 9.4.4.** A linear operator  $\mathbf{T}$  defined on a finite-dimensional inner product space  $(V, \langle \cdot, \cdot \rangle)$  is called a **normal operator** iff holds  $[\mathbf{T}, \mathbf{T}^*] = \mathbf{0}$ , that is, the operator commutes with its adjoint.

An interesting characterization of normal operators is the following: A linear operator  $\mathbf{T}$  on an inner product space is normal iff  $\|\mathbf{T}(\mathbf{u})\| = \|\mathbf{T}^*(\mathbf{u})\|$  holds for all  $\mathbf{u} \in V$ . Normal operators are particularly important because for these operators hold the Spectral Theorem, which we study in Chapter 9.

Two particular cases of normal operators are often used in physics. A linear operator  $\mathbf{T}$  on an inner product space is called a **unitary operator** iff  $\mathbf{T}^* = \mathbf{T}^{-1}$ , that is, the adjoint is the inverse operator. Unitary operators are normal operators, since

$$\mathbf{T}^* = \mathbf{T}^{-1} \Rightarrow \begin{cases} \mathbf{T}\mathbf{T}^* = \mathbf{I}, \\ \mathbf{T}^*\mathbf{T} = \mathbf{I}, \end{cases} \Rightarrow [\mathbf{T}, \mathbf{T}^*] = \mathbf{0}.$$

Unitary operators preserve the length of a vector, since

$$\|\mathbf{v}\|^2 = \langle \mathbf{v}, \mathbf{v} \rangle = \langle \mathbf{v}, \mathbf{T}^{-1}(\mathbf{T}(\mathbf{v})) \rangle = \langle \mathbf{v}, \mathbf{T}^*(\mathbf{T}(\mathbf{v})) \rangle = \langle \mathbf{T}(\mathbf{v}), \mathbf{T}(\mathbf{v}) \rangle = \|\mathbf{T}(\mathbf{v})\|^2.$$

Unitary operators defined on a complex inner product space are particularly important in quantum mechanics. The particular case of unitary operators on a real inner product space are called **orthogonal operators**. So, orthogonal operators do not change the length of a vector. Examples of orthogonal operators are rotations in  $\mathbb{R}^3$ .

A linear operator  $\mathbf{T}$  on an inner product space is called an **Hermitian operator** iff  $\mathbf{T}^* = \mathbf{T}$ , that is, the adjoint is the original operator. This definition agrees with the definition of Hermitian matrices given in Chapter 2.

**EXAMPLE 9.4.2:** Consider the inner product space  $(\mathbb{C}^3, \cdot)$  and the linear operator  $\mathbf{T}$  with matrix in the standard basis of  $\mathbb{C}^3$  given by

$$[\mathbf{T}(\mathbf{x})] = \begin{bmatrix} x_1 - ix_2 + x_3 \\ ix_1 - x_3 \\ x_1 - x_2 + x_3 \end{bmatrix}, \quad [\mathbf{x}] = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}.$$

Show that  $\mathbf{T}$  is Hermitian.

**SOLUTION:** We need to compute the adjoint of  $\mathbf{T}$ . The matrix of this operator in the standard basis of  $\mathbb{C}^3$  is given by

$$[\mathbf{T}] = \begin{bmatrix} 1 & -i & 1 \\ i & 0 & -1 \\ 1 & -1 & 1 \end{bmatrix}.$$



Since the standard basis is an orthonormal basis with respect to the dot product, Proposition 9.4.3 implies that

$$[\mathbf{T}^*] = [\mathbf{T}]^* = \begin{bmatrix} 1 & -i & 1 \\ i & 0 & -1 \\ 1 & -1 & 1 \end{bmatrix}^* = \begin{bmatrix} 1 & -i & 1 \\ i & 0 & -1 \\ 1 & -1 & 1 \end{bmatrix} = [\mathbf{T}].$$

Therefore,  $\mathbf{T}^* = \mathbf{T}$ .

◁

**9.4.1. Exercises.****9.4.1.-** .**9.4.2.-** .

## CHAPTER 10. APPENDIX

## 10.1. REVIEW EXERCISES

## Chapter 1: Linear systems

- 1.- Consider the linear system

$$\begin{aligned} 2x_1 + 3x_2 - x_3 &= 6 \\ -x_1 - x_2 + 2x_3 &= -2 \\ x_1 + 2x_3 &= 2 \end{aligned}$$

- (a) Use Gauss operations to find the reduced echelon form of the augmented matrix for this system.
- (b) Is this system consistent? If “yes,” find all the solutions.

- 2.- Find all the solutions
- $\mathbf{x}$
- to the linear system
- $A\mathbf{x} = \mathbf{b}$
- and express them in vector form, where

$$A = \begin{bmatrix} 1 & -2 & -1 \\ 2 & 1 & 8 \\ 1 & -1 & 1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}.$$

- 3.- Consider the matrix and the vector

$$A = \begin{bmatrix} 1 & -2 & 7 \\ 1 & 1 & 1 \\ 2 & 2 & 2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ 1 \\ 3 \end{bmatrix}.$$

Is the vector  $\mathbf{b}$  a linear combination of the column vectors of  $A$ ?

- 4.- Let
- $s$
- be a real number, and consider the system

$$\begin{aligned} sx_1 - 2sx_2 &= -1, \\ 3x_1 + 6sx_2 &= 3. \end{aligned}$$

- (a) Determine the values of the parameter  $s$  for which the system above has a unique solution.
- (b) For all the values of  $s$  such that the system above has a unique solution, find that solution.

- 5.- Find the values of
- $k$
- such that the system below has no solution; has one solution; has infinitely many solutions;

$$\begin{aligned} kx_1 + x_2 &= 1 \\ x_1 + kx_2 &= 1. \end{aligned}$$

- 6.- Find a condition on the components of vector
- $\mathbf{b}$
- such that the system
- $A\mathbf{x} = \mathbf{b}$
- is consistent, where

$$A = \begin{bmatrix} 1 & 1 & -1 \\ 2 & 0 & -6 \\ 3 & 1 & -7 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}.$$

- 7.- Find the general solution to the homogeneous linear system with coefficient matrix

$$A = \begin{bmatrix} 1 & 3 & -1 & 5 \\ 2 & 1 & 3 & 0 \\ 3 & 2 & 4 & 1 \end{bmatrix},$$

and write this general solution in vector form.

- 8.- (a) Find a value of the constants
- $h$
- and
- $k$
- such that the non-homogeneous linear system below is consistent and has one free variable.

$$\begin{aligned} x_1 + hx_2 + 5x_3 &= 1, \\ x_2 - 2x_3 &= k, \\ x_1 + 3x_2 - 3x_3 &= 5. \end{aligned}$$

- (b) Using the value of the constants  $h$  and  $k$  found in part (a), find the general solution to the system given in part (a).

- 9.- (a) Find the general solution to the system below and write it in vector form,

$$\begin{aligned} x_1 + 2x_2 - x_3 &= 2, \\ 3x_1 + 7x_2 - 3x_3 &= 7, \\ x_1 + 4x_2 - x_3 &= 4. \end{aligned}$$

- (b) Sketch a graph on  $\mathbb{R}^3$  of the general solution found in part (a).

## Chapter 2: Matrix algebra

- 1.- Consider the vectors

$$\mathbf{u} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad \mathbf{v} = \begin{bmatrix} 1 \\ -1 \end{bmatrix},$$

and the linear function  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  such that

$$T(\mathbf{u}) = \begin{bmatrix} 1 \\ 3 \end{bmatrix}, \quad T(\mathbf{v}) = \begin{bmatrix} 3 \\ 1 \end{bmatrix}.$$

Find the matrix  $\mathbf{A} = [T(\mathbf{e}_1), T(\mathbf{e}_2)]$  of the linear transformation, where

$$\mathbf{e}_1 = \frac{1}{2}(\mathbf{u} + \mathbf{v}), \quad \mathbf{e}_2 = \frac{1}{2}(\mathbf{u} - \mathbf{v}).$$

Show your work.

- 2.- Find the matrix for the linear transformation  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  representing a reflection on the plane along the vertical axis followed by a rotation by  $\theta = \pi/3$  counterclockwise.
- 3.- Which of the following matrices below is equal to  $(\mathbf{A} + \mathbf{B})^2$  for every square matrices  $\mathbf{A}$  and  $\mathbf{B}$ ?

$$\begin{aligned} &(\mathbf{B} + \mathbf{A})^2, \\ &\mathbf{A}^2 + 2\mathbf{A}\mathbf{B} + \mathbf{B}^2, \\ &(\mathbf{A} + \mathbf{B})(\mathbf{B} + \mathbf{A}), \\ &\mathbf{A}^2 + \mathbf{A}\mathbf{B} + \mathbf{B}\mathbf{A} + \mathbf{B}^2, \\ &\mathbf{A}(\mathbf{A} + \mathbf{B}) + (\mathbf{A} + \mathbf{B})\mathbf{B}. \end{aligned}$$

- 4.- Find a matrix  $\mathbf{A}$  solution of the matrix equation

$$\mathbf{A}\mathbf{B} + 2\mathbf{I}_2 = \begin{bmatrix} 5 & 4 \\ -2 & 3 \end{bmatrix},$$

where

$$\mathbf{B} = \begin{bmatrix} 7 & 3 \\ 2 & 1 \end{bmatrix}.$$

- 5.- Consider the matrix

$$\mathbf{A} = \begin{bmatrix} 5 & 3 & s \\ 1 & 2 & -1 \\ 2 & 1 & 1 \end{bmatrix}.$$

Find the value(s) of the constant  $s$  such that the matrix  $\mathbf{A}$  is invertible.

- 6.- Let  $\mathbf{A}$  be an  $n \times n$  matrix,  $\mathbf{D}$  be and  $m \times m$  matrix, and  $\mathbf{C}$  be and  $m \times n$  matrix. Assume that both  $\mathbf{A}$  and  $\mathbf{D}$  are invertible matrices and denote by  $\mathbf{A}^{-1}$ ,  $\mathbf{D}^{-1}$  their respective inverse matrices. Let  $\mathbf{M}$  be the  $(n + m) \times (n + m)$  matrix

$$\mathbf{M} = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}.$$

Find an  $m \times n$  matrix  $\mathbf{X}$  (in terms of any of the matrices  $\mathbf{A}$ ,  $\mathbf{D}$ ,  $\mathbf{A}^{-1}$ ,  $\mathbf{D}^{-1}$ , and  $\mathbf{C}$ ) such that  $\mathbf{M}$  is invertible and the inverse is given by

$$\mathbf{M}^{-1} = \begin{bmatrix} \mathbf{A}^{-1} & \mathbf{0} \\ \mathbf{X} & \mathbf{D}^{-1} \end{bmatrix}.$$

- 7.- Consider the matrix and the vector

$$\mathbf{A} = \begin{bmatrix} 1 & -2 & 7 \\ 1 & 1 & 1 \\ 2 & 2 & 2 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ 1 \\ 3 \end{bmatrix}.$$

- (a) Does vector  $\mathbf{b}$  belong to the  $R(\mathbf{A})$ ?  
 (b) Does vector  $\mathbf{b}$  belong to the  $N(\mathbf{A})$ ?

- 8.- Consider the matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 6 & -7 \\ -2 & 3 & 2 & 0 \\ 0 & -1 & -2 & 2 \end{bmatrix}.$$

Find  $N(\mathbf{A})$  and the  $R(\mathbf{A}^T)$ .

- 9.- Consider the matrix

$$\mathbf{A} = \begin{bmatrix} 2 & 1 & 3 \\ 4 & 5 & 7 \\ 6 & 9 & 12 \end{bmatrix}.$$

- (a) Find the LU-factorization of  $\mathbf{A}$ .  
 (b) Use the LU-factorization above to find the solutions  $x_1$  and  $x_2$  of the systems  $\mathbf{A}\mathbf{x} = \mathbf{e}_1$  and  $\mathbf{A}\mathbf{x}_2 = \mathbf{e}_2$ , where

$$\mathbf{e}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{e}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}.$$

- (c) Use the LU-factorization above to find  $\mathbf{A}^{-1}$ .

**Chapter 3: Determinants**

- 1.- Find the determinant of the matrices

$$A = \begin{bmatrix} 1+i & -3i \\ -4i & 1-2i \end{bmatrix},$$

$$B = \begin{bmatrix} 0 & -1 & 0 & 0 \\ 0 & -2 & 3 & -2 \\ 0 & 0 & -1 & -3 \\ -2 & 0 & 0 & 1 \end{bmatrix},$$

$$C = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}.$$

- 2.- Given matrix A below, find the cofactors matrix C, and explicitly show that
- $AC^T = \det(A)I_3$
- , where

$$A = \begin{bmatrix} 1 & 3 & -1 \\ 4 & 0 & 1 \\ 2 & 1 & 3 \end{bmatrix}.$$

- 3.- Given matrix A below, find the coefficients
- $(A^{-1})_{13}$
- and
- $(A^{-1})_{23}$
- of the inverse matrix
- $A^{-1}$
- , where

$$A = \begin{bmatrix} 5 & 3 & 1 \\ 1 & 2 & -1 \\ 2 & 1 & 1 \end{bmatrix}.$$

- 4.- Find the change in the area of the parallelogram formed by the vectors

$$\mathbf{u} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad \mathbf{v} = \begin{bmatrix} 2 \\ 1 \end{bmatrix},$$

when this parallelogram is transformed under the following linear transformation,  $A: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ ,

$$A = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}.$$

- 5.- Find the volume of the parallelepiped formed by the vectors

$$\mathbf{v}_1 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix}, \quad \mathbf{v}_3 = \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix}.$$

- 6.- Consider the matrix

$$A = \begin{bmatrix} 4 & 2 \\ 1 & 3 \end{bmatrix}.$$

Find the values of the scalar  $\lambda$  such that the matrix  $(A - \lambda I_2)$  is not invertible.

- 7.- Prove the following: If there exists an integer
- $k \geq 1$
- such that
- $A \in \mathbb{F}^{n,n}$
- satisfies
- $A^k = 0$
- , then
- $\det(A) = 0$
- .

- 8.- Assume that matrix
- $A \in \mathbb{F}^{n,n}$
- satisfies the equation
- $A^2 = I_n$
- . Find all possible values of
- $\det(A)$
- .

- 9.- Use Cramer's rule to find the solution of the linear system
- $A\mathbf{x} = \mathbf{b}$
- , where

$$A = \begin{bmatrix} 1 & 4 & -1 \\ 1 & 1 & 1 \\ 2 & 0 & 3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}.$$

### Chapter 4: Vector spaces

- 1.- Determine which of the following subsets of  $\mathbb{R}^{3,3}$  are subspaces:
- The symmetric matrices.
  - The skew-symmetric matrices.
  - The matrices  $A$  with  $A^2 = A$ .
  - The matrices  $A$  with  $\text{tr}(A) = 0$ .
  - The matrices  $A$  with  $\det(A) = 0$ .

In the case that the set is a subspace, find a basis of this subspace.

- 2.- Find the dimension and give a basis of the subspace  $W \subset \mathbb{R}^3$  given by

$$\left\{ \begin{bmatrix} -a + b + c - 3d \\ b + 3c - d \\ a + 2b + 8c \end{bmatrix} \text{ with } a, b, c, d \in \mathbb{R} \right\}.$$

- 3.- Find the dimension of both the null space and the range space of the matrix

$$A = \begin{bmatrix} 1 & 1 & 1 & 5 & 2 \\ 2 & 1 & 4 & 7 & 3 \\ 0 & -1 & 2 & -3 & -1 \end{bmatrix}.$$

- 4.- Consider the matrix

$$A = \begin{bmatrix} 1 & -1 & 5 \\ 0 & 1 & -2 \\ 1 & 3 & -3 \end{bmatrix}.$$

- Find a basis for the null space of  $A$ .
- Find a basis for the subspace in  $\mathbb{R}^3$  consisting of all vectors  $\mathbf{b} \in \mathbb{R}^3$  such that the linear system  $A\mathbf{x} = \mathbf{b}$  is consistent.

- 5.- Show whether the following statement is true or false: Given a vector space  $V$ , if the set

$$\{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3\} \subset V$$

is linearly independent, then so is the set

$$\{\mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3\},$$

where

$$\mathbf{w}_1 = (\mathbf{v}_1 + \mathbf{v}_2),$$

$$\mathbf{w}_2 = (\mathbf{v}_1 + \mathbf{v}_3),$$

$$\mathbf{w}_3 = (\mathbf{v}_2 + \mathbf{v}_3).$$

- 6.- Show that the set  $U \subset \mathbb{P}_3$  given by all polynomials satisfying the condition

$$\int_0^1 \mathbf{p}(x) dx = 0$$

is a subspace of  $\mathbb{P}_3$ . Find a basis for  $U$ .

- 7.- Determine whether the set  $U \subset \mathbb{P}_2$  of all polynomials of the form

$$\mathbf{p}(x) = a + ax + ax^2$$

with  $a \in \mathbb{F}$ , is a subspace of  $\mathbb{P}_2$ .

## Chapter 5: Linear transformations

1.- Consider the matrix

$$A = \begin{bmatrix} 1 & 2 & 6 & -7 \\ -2 & 3 & 2 & 0 \\ 0 & -1 & -2 & 2 \end{bmatrix}.$$

- (1) Find a basis for  $R(A)$ .
- (2) Find a basis for  $N(A)$ .
- (3) Consider the linear transformation  $A : \mathbb{R}^4 \rightarrow \mathbb{R}^3$  determined by  $A$ . Is it injective? Is it surjective? Justify your answers.

2.- Let  $T : \mathbb{R}^3 \rightarrow \mathbb{R}^2$  be the linear transformation given by

$$[T(\mathbf{x}_s)]_{\tilde{s}} = \begin{bmatrix} 2x_1 + 6x_2 - 2x_3 \\ 3x_1 + 8x_2 + 2x_3 \end{bmatrix}_s,$$

$$[\mathbf{x}]_s = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix},$$

where  $\mathcal{S}$  and  $\tilde{\mathcal{S}}$  are the standard bases in  $\mathbb{R}^3$  and  $\mathbb{R}^2$ , respectively.

- (a) Is  $T$  injective? Is  $T$  surjective?
- (b) Find all solutions of the linear system  $T(\mathbf{x}_s) = \mathbf{b}_{\tilde{s}}$ , where

$$\mathbf{b}_{\tilde{s}} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}.$$

- (c) Is the set of all solutions found in part (b) a subspace of  $\mathbb{R}^3$ ?

3.- Let  $A : \mathbb{C}^3 \rightarrow \mathbb{C}^4$  be the linear transformation

$$A = \begin{bmatrix} 1 & 0 & i \\ 0 & i & 1 \\ 1-i & 0 & 1+i \\ i & 1 & 0 \end{bmatrix}.$$

Find a basis for  $N(A)$  and  $R(A)$ .

4.- Let  $D : \mathbb{P}_3 \rightarrow \mathbb{P}_3$  be the differentiation operator,

$$D(\mathbf{p})(x) = \frac{d\mathbf{p}}{dx}(x),$$

and  $I : \mathbb{P}_3 \rightarrow \mathbb{P}_3$  the identity operator. Let  $\mathcal{S} = (1, x, x^2, x^3)$  be the standard ordered basis of  $\mathbb{P}_3$ . Show that the matrix of the operator

$$(I - D^2) : \mathbb{P}_3 \rightarrow \mathbb{P}_3$$

in the basis  $\mathcal{S}$  is invertible.

**Chapter 6: Inner product spaces**

1.- Let  $(V, \langle \cdot, \cdot \rangle)$  be a real inner product space. Show that  $(\mathbf{x} - \mathbf{y}) \perp (\mathbf{x} + \mathbf{y})$  iff  $\|\mathbf{x}\| = \|\mathbf{y}\|$ .

2.- Consider the inner product space given by  $(\mathbb{F}^n, \cdot)$ . Prove that for every matrix  $A \in \mathbb{F}^{n,n}$  holds

$$\mathbf{x} \cdot (A\mathbf{y}) = (A^* \mathbf{x}) \cdot \mathbf{y}.$$

3.- A matrix  $A \in \mathbb{F}^{n,n}$  is called **unitary** iff

$$AA^* = A^*A = I_n.$$

Show that a unitary matrix does not change the norm of a vector in the inner product space  $(\mathbb{F}^n, \cdot)$ , that is, for all  $\mathbf{x} \in \mathbb{F}^n$  and all unitary matrix  $A \in \mathbb{F}^{n,n}$  holds

$$\|A\mathbf{x}\| = \|\mathbf{x}\|.$$

4.- Find all vectors in the inner product space  $(\mathbb{R}^4, \cdot)$  perpendicular to both

$$\mathbf{v} = \begin{bmatrix} 1 \\ 4 \\ 4 \\ 1 \end{bmatrix} \quad \text{and} \quad \mathbf{u} = \begin{bmatrix} 2 \\ 9 \\ 8 \\ 2 \end{bmatrix}.$$

5.- Consider the inner product space given by  $(\mathbb{C}^3, \cdot)$  and the subspace  $W \subset \mathbb{C}^3$  spanned by the vectors

$$\mathbf{u} = \begin{bmatrix} 1+i \\ 1 \\ i \end{bmatrix}, \quad \mathbf{v} = \begin{bmatrix} -1 \\ 0 \\ 2-i \end{bmatrix}.$$

- (a) Use the Gram-Schmidt method to find an orthonormal basis for  $W$ .
- (b) Extend the orthonormal basis of  $W$  into an orthonormal basis for  $\mathbb{C}^3$ .



**Chapter 7: Approximation methods**

1.- .

2.- .

**Chapter 8: Normed spaces**

1.- .

2.- .

**Chapter 9: Spectral decomposition**

1.- .

2.- .

## 10.2. PRACTICE EXAMS

**Instructions to use the Practice Exams.** The idea of these lecture notes is to help anyone interested in learning linear algebra. More often than not such person is a college student. An unfortunate aspect of our education system is that students must pass an exam to prove they understood the ideas of a given subject. To prepare the student for such exam is the purpose of the practice exams in this Section.

These practice exams once were actual exams taken by student in previous courses. These exams can be useful to other students if they are used properly; one way is the following: Study the course material first, do all the exercises at the end of every Section; do the review problems in Section 10.1; then and only then take the first practice exam and do it. Think of it as an actual exam. Do not look at the notes or any other literature. Do the whole exam. Watch your time. You have only two hours to do it. After you finish, you can grade yourself. You have the solutions to the exam at the end of the Chapter. Never, ever look at the solutions before you finish the exam. If you do it, the practice exam is worthless. Really, worthless; you will not have the solutions when you do the actual exam. The story does not finish here. Pay close attention at the exercises you did wrong, if any. Go back to the class material and do extra problems on those subjects. Review all subjects you think you are not well prepared. Then take the second practice exam. Follow a similar procedure. Review the subject related to the practice exam exercises you had difficulties to solve. Then, do the next practice exam. You have three of them. After doing the last one, you should have a good idea of what your actual grade in the class should be.

**Practice Exam 1.** (Two hours.)

1. Consider the matrix  $A = \begin{bmatrix} 5 & 3 & 1 \\ 1 & 2 & -1 \\ 2 & 1 & 1 \end{bmatrix}$ . Find the coefficients  $(A^{-1})_{21}$  and  $(A^{-1})_{32}$  of the matrix  $A^{-1}$ , that is, of the inverse matrix of  $A$ . Show your work.

2. (a) Find  $k \in \mathbb{R}$  such that the volume of the parallelepiped formed by the vectors below is equal to 4, where

$$\mathbf{v}_1 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix}, \quad \mathbf{v}_3 = \begin{bmatrix} k \\ 1 \\ 1 \end{bmatrix}$$

- (b) Set  $k = 1$  and define the matrix  $A = [\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3]$ . Matrix  $A$  determines the linear transformation  $A: \mathbb{R}^3 \rightarrow \mathbb{R}^3$ . Is this linear transformation injective (one-to-one)? Is it surjective (onto)?

3. Determine whether the subset  $V \subset \mathbb{R}^3$  is a subspace, where

$$V = \left\{ \begin{bmatrix} -a + b \\ a - 2b \\ a - 7b \end{bmatrix} \text{ with } a, b \in \mathbb{R} \right\}.$$

If the set is a subspace, find an orthogonal basis in the inner product space  $(\mathbb{R}^3, \cdot)$ .

4. True or False: (Justify your answers.)

- (a) If the set of columns of  $A \in \mathbb{F}^{m,n}$  is a linearly independent set, then  $A\mathbf{x} = \mathbf{b}$  has exactly one solution for every  $\mathbf{b} \in \mathbb{F}^m$ .  
 (b) The set of column vectors of an  $5 \times 7$  is never linearly independent.

5. Consider the vector space  $\mathbb{R}^2$  with the standard basis  $\mathcal{S}$  and let  $T: \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be the linear transformation

$$[T]_{ss} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix}_{ss}.$$

Find  $[T]_{bb}$ , the matrix of  $T$  in the basis  $\mathcal{B}$ , where  $\mathcal{B} = \left\{ [\mathbf{b}_1]_s = \begin{bmatrix} 1 \\ 2 \end{bmatrix}_s, [\mathbf{b}_2]_s = \begin{bmatrix} 1 \\ -2 \end{bmatrix}_s \right\}$ .

6. Consider the linear transformations  $T: \mathbb{R}^3 \rightarrow \mathbb{R}^2$  and  $S: \mathbb{R}^3 \rightarrow \mathbb{R}^3$  given by

$$\left[ T \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{s_3} \right) \right]_{s_2} = \begin{bmatrix} x_1 - x_2 + x_3 \\ -x_1 + 2x_2 + x_3 \end{bmatrix}_{s_2}, \quad \left[ S \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{s_3} \right) \right]_{s_3} = \begin{bmatrix} 3x_3 \\ 2x_2 \\ x_1 \end{bmatrix}_{s_3},$$

where  $\mathcal{S}_3$  and  $\mathcal{S}_2$  are the standard basis of  $\mathbb{R}^3$  and  $\mathbb{R}^2$ , respectively.

- (a) Find a matrix  $[T]_{s_3 s_2}$  and the matrix  $[S]_{s_3 s_3}$ . Show your work.  
 (b) Find the matrix of the composition  $T \circ S: \mathbb{R}^3 \rightarrow \mathbb{R}^2$  in the standard basis, that is, find  $[T \circ S]_{s_3 s_2}$ .  
 (c) Is  $T \circ S$  injective (one-to-one)? Is  $T \circ S$  surjective (onto)? Justify your answer.

7. Consider the matrix  $A = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 2 & 1 \end{bmatrix}$  and the vector  $\mathbf{b} = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}$ .

- (a) Find the least-squares solution  $\hat{\mathbf{x}}$  to the matrix equation  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .  
 (b) Verify whether the vector  $\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}$  belong to the space  $R(\mathbf{A})^\perp$ ? Justify your answers.

8. Consider the matrix  $\mathbf{A} = \begin{bmatrix} -1/2 & -3 \\ 1/2 & 2 \end{bmatrix}$ .

- (a) Show that matrix  $\mathbf{A}$  is diagonalizable.  
 (b) Using that  $\mathbf{A}$  is diagonalizable, find the  $\lim_{k \rightarrow \infty} \mathbf{A}^k$ .

9. Let  $(V, \langle \cdot, \cdot \rangle)$  be an inner product space with inner product norm  $\|\cdot\|$ . Let  $\mathbf{T}: V \rightarrow V$  be a linear transformation and  $\mathbf{x}, \mathbf{y} \in V$  be vectors satisfying the following conditions:

$$\mathbf{T}(\mathbf{x}) = 2\mathbf{x}, \quad \mathbf{T}(\mathbf{y}) = -3\mathbf{y}, \quad \|\mathbf{x}\| = 1/3, \quad \|\mathbf{y}\| = 1, \quad \mathbf{x} \perp \mathbf{y}.$$

- (a) Compute  $\|\mathbf{v}\|$  for the vector  $\mathbf{v} = 3\mathbf{x} - \mathbf{y}$ .  
 (b) Compute  $\|\mathbf{T}(\mathbf{v})\|$  for the vector  $\mathbf{v}$  given above.

10. Consider the matrix  $\mathbf{A} = \begin{bmatrix} 2 & -1 & 2 \\ 0 & 1 & h \\ 0 & 0 & 2 \end{bmatrix}$ .

- (a) Find all eigenvalues of matrix  $\mathbf{A}$  and their corresponding algebraic multiplicities.  
 (b) Find the value(s) of the real number  $h$  such that the matrix  $\mathbf{A}$  above has a two-dimensional eigenspace, and find a basis for this eigenspace.

**Practice Exam 2.** (Two hours.)

1. Consider the matrix  $A = \begin{bmatrix} -2 & 3 & -1 \\ 1 & 2 & -1 \\ -2 & -1 & 1 \end{bmatrix}$ . Find the coefficients  $(A^{-1})_{13}$  and  $(A^{-1})_{21}$  of the inverse matrix of  $A$ . Show your work.

2. Consider the vector space  $\mathbb{P}_3([0, 1])$  with the inner product

$$\langle \mathbf{p}, \mathbf{q} \rangle = \int_0^1 \mathbf{p}(x)\mathbf{q}(x) dx.$$

Given the set  $\mathcal{U} = \{\mathbf{p}_1 = x^2, \mathbf{p}_2 = x^3\}$ , find an orthogonal basis for the subspace  $U = \text{Span}(\mathcal{U})$  using the Gram-Schmidt method on the set  $\mathcal{U}$  starting with the vector  $\mathbf{p}_1$ .

3. Consider the matrix  $A = \begin{bmatrix} 1 & 3 & 1 & 1 \\ 2 & 6 & 3 & 0 \\ 3 & 9 & 5 & -1 \end{bmatrix}$ .

- (a) Verify that the vector  $\mathbf{v} = \begin{bmatrix} 3 \\ 1 \\ -4 \\ -2 \end{bmatrix}$  belongs to the null space of  $A$ .

- (b) Extend the set  $\{\mathbf{v}\}$  into a basis of the null space of  $A$ .

4. Use Cramer's rule to find the solution to the linear system

$$2x_1 + x_2 - x_3 = 0$$

$$x_1 + x_3 = 1$$

$$x_1 + 2x_2 + 3x_3 = 0.$$

5. Let  $\mathcal{S}_3$  and  $\mathcal{S}_2$  be standard bases of  $\mathbb{R}^3$  and  $\mathbb{R}^2$ , respectively, and consider the linear transformation  $\mathbf{T}: \mathbb{R}^3 \rightarrow \mathbb{R}^2$  given by

$$\left[ \mathbf{T} \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) \right]_{\mathcal{S}_2} = \begin{bmatrix} -x_1 + 2x_2 - x_3 \\ x_1 + x_3 \end{bmatrix}_{\mathcal{S}_2},$$

and introduce the bases  $\mathcal{U} \subset \mathbb{R}^3$  and  $\mathcal{V} \subset \mathbb{R}^2$  given by

$$\mathcal{U} = \left\{ [\mathbf{u}_1]_{\mathcal{S}_3} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}_{\mathcal{S}_3}, [\mathbf{u}_2]_{\mathcal{S}_3} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}_{\mathcal{S}_3}, [\mathbf{u}_3]_{\mathcal{S}_3} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}_{\mathcal{S}_3} \right\},$$

$$\mathcal{V} = \left\{ [\mathbf{v}_1]_{\mathcal{S}_2} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}_{\mathcal{S}_2}, [\mathbf{v}_2]_{\mathcal{S}_2} = \begin{bmatrix} -3 \\ 2 \end{bmatrix}_{\mathcal{S}_2} \right\}.$$

Find the matrices  $[\mathbf{T}]_{\mathcal{S}_3\mathcal{S}_2}$  and  $[\mathbf{T}]_{\mathcal{U}\mathcal{V}}$ . Show your work.

6. Consider the inner product space  $(\mathbb{R}^{2,2}, \langle \cdot, \cdot \rangle_F)$  and the subspace

$$W = \text{Span} \left\{ \mathbf{E}_1 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \mathbf{E}_2 = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \right\} \subset \mathbb{R}^{2,2}.$$

Find a basis for  $W^\perp$ , the orthogonal complement of  $W$ .

7. Consider the matrix  $\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 1 & -1 \\ -2 & 1 \end{bmatrix}$  and the vector  $\mathbf{b} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ .
- Find the least-squares solution  $\hat{\mathbf{x}}$  to the matrix equation  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .
  - Verify that the vector  $\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}$ , where  $\hat{\mathbf{x}}$  is the least-squares solution found in part (7a), belongs to the space  $R(\mathbf{A})^\perp$ , the orthogonal complement of  $R(\mathbf{A})$ .
8. Suppose that a matrix  $\mathbf{A} \in \mathbb{R}^{3,3}$  has eigenvalues  $\lambda_1 = 1$ ,  $\lambda_2 = 2$ , and  $\lambda_3 = 4$ .
- Find the trace of  $\mathbf{A}$ , find the trace of  $\mathbf{A}^2$ , and find the determinant of  $\mathbf{A}$ .
  - Is matrix  $\mathbf{A}$  invertible? If your answer is “yes”, then prove it and find  $\det(\mathbf{A}^{-1})$ ; if your answer is “no”, then prove it.
9. Consider the matrix  $\mathbf{A} = \begin{bmatrix} 7 & 5 \\ 3 & -7 \end{bmatrix}$ .
- Find the eigenvalues and eigenvectors of  $\mathbf{A}$ .
  - Compute the matrix  $e^{\mathbf{A}}$ .
10. Find the function  $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^2$  solution of the initial value problem
- $$\frac{d}{dt}\mathbf{x}(t) = \mathbf{A}\mathbf{x}(t), \quad \mathbf{x}(0) = \mathbf{x}_0,$$
- where the matrix  $\mathbf{A} = \begin{bmatrix} -5 & 2 \\ -12 & 5 \end{bmatrix}$  and the vector  $\mathbf{x}_0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ .



**Practice Exam 3.** (Two hours.)

1. (a) Find the LU-factorization of the matrix  $A = \begin{bmatrix} 1 & 2 & 3 \\ -1 & 2 & 2 \\ 2 & -8 & -3 \end{bmatrix}$ .
- (b) Use the LU-factorization above to find the solution of the linear system  $A\mathbf{x} = \mathbf{b}$ , where  $\mathbf{b} = \begin{bmatrix} 1 \\ -2 \\ -1 \end{bmatrix}$ .

2. Determine whether the following sets  $W_1$  and  $W_2$  are subspaces of the vector space  $\mathbb{P}_2([0, 1])$ . If your answer is “yes,” find a basis of the subspace.

(a)  $W_1 = \{\mathbf{p} \in \mathbb{P}_2([0, 1]) : \int_0^1 \mathbf{p}(x) dx \leq 1\}$ ;

(b)  $W_2 = \{\mathbf{p} \in \mathbb{P}_2([0, 1]) : \int_0^1 x \mathbf{p}(x) dx = 0\}$ .

3. Find a basis of  $\mathbb{R}^4$  containing a basis of the null space of matrix

$$A = \begin{bmatrix} 1 & 2 & -4 & 3 \\ 2 & 1 & 1 & 3 \\ 1 & 1 & -1 & 2 \\ 3 & 2 & 0 & 5 \end{bmatrix}.$$

4. Given a matrix  $A \in \mathbb{F}^{n,n}$ , introduce its characteristic polynomial  $p_A(\lambda) = \det(A - \lambda I_n)$ . This polynomial has the form  $p_A(\lambda) = a_0 + a_1\lambda + \cdots + a_n\lambda^n$  for appropriate scalars  $a_0, \dots, a_n$ . Now introduce a matrix-valued function  $P_A : \mathbb{F}^{n,n} \rightarrow \mathbb{F}^{n,n}$  as follows

$$P_A(\mathbf{X}) = a_0 I_n + a_1 \mathbf{X} + \cdots + a_n \mathbf{X}^n.$$

Determine whether the following statement is true or false and justify your answer: If matrix  $A \in \mathbb{F}^{n,n}$  is diagonalizable, then  $P_A(A) = 0$ . (That is,  $P_A(A)$  is the zero matrix.)

5. Consider the vector space  $\mathbb{R}^2$  with ordered bases

$$\mathcal{S} = \left( \mathbf{e}_{1s} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}_s, \mathbf{e}_{2s} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}_s \right), \quad \mathcal{U} = \left( \mathbf{u}_{1s} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}_s, \mathbf{u}_{2s} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}_s \right).$$

Let  $\mathbf{T} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be a linear transformation given by

$$[\mathbf{T}(\mathbf{u}_1)]_s = \begin{bmatrix} -3 \\ 1 \end{bmatrix}_s, \quad [\mathbf{T}(\mathbf{u}_2)]_s = \begin{bmatrix} 1 \\ 3 \end{bmatrix}_s.$$

- (a) Find the matrix  $\mathbf{T}_{ss}$ .
- (b) Find the matrix  $\mathbf{T}_{uu}$ .

6. Consider the matrix  $A = \begin{bmatrix} 1 & 2 \\ 2 & -1 \\ 1 & 1 \end{bmatrix}$  and the vector  $\mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$ .

- (a) Find the least-squares solution  $\hat{\mathbf{x}}$  to the matrix equation  $A\mathbf{x} = \mathbf{b}$ .
- (b) Find the vector on  $R(A)$  that is the closest to the vector  $\mathbf{b}$  in  $\mathbb{R}^3$ .

7. Consider the inner product space  $(\mathbb{C}^3, \cdot)$  and the subspace  $W \subset \mathbb{C}^3$  given by

$$W = \text{Span} \left( \left\{ \begin{bmatrix} i \\ 1 \\ i \end{bmatrix} \right\} \right).$$

- (a) Find a basis for  $W^\perp$ , the orthogonal complement of  $W$ .
- (b) Find an orthonormal basis for  $W^\perp$ .

8. (a) Find the eigenvalues and eigenvectors of matrix  $A = \begin{bmatrix} 5 & 4 \\ 4 & 5 \end{bmatrix}$ , and show that  $A$  is diagonalizable.
- (b) Knowing that matrix  $A$  above is diagonalizable, explicitly find a square root of matrix  $A$ , that is, find a matrix  $X$  such that  $X^2 = A$ . How many square roots does matrix  $A$  have?
9. Consider the matrix  $A = \begin{bmatrix} 8 & -18 \\ 3 & -7 \end{bmatrix}$ .
- (a) Find the eigenvalues and eigenvectors of  $A$ .
- (b) Compute explicitly the matrix-valued function  $e^{At}$  for  $t \in \mathbb{R}$ .
10. Find the function  $x : \mathbb{R} \rightarrow \mathbb{R}^2$  solution of the initial value problem

$$\frac{d}{dt}x(t) = Ax(t), \quad x(0) = x_0,$$

where the matrix  $A = \begin{bmatrix} -7 & 12 \\ -4 & 7 \end{bmatrix}$  and the vector  $x_0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ .

## 10.3. ANSWERS TO EXERCISES

## Chapter 1: Linear systems

## Section 1.1: Row and Column Pictures

1.1.1.-  $x = 2, y = 3$ .

1.1.2.- The system is not consistent.

1.1.3.- The nonlinear system has two solutions,

$$(x = \sqrt{2}, y = \sqrt{2}),$$

$$(x = -\sqrt{2}, y = -\sqrt{2}).$$

1.1.4.- The system is inconsistent.

1.1.5.- Subtract the first equation from the second. To the resulting equation subtract the third equation. One gets  $0 = 1$ . The system is inconsistent.

1.1.6.-  $k = \pm 2$ .

1.1.7.- The system is consistent. The solution is

$$x_1 = -1, \quad x_2 = 1.$$

1.1.8.-

(a) The vectors  $A_1$  and  $A_2$  are collinear, while vector  $b$  is not parallel to  $A_1$  and  $A_2$ .

(b) The system is not consistent.

(c) The system is consistent and the solution is not unique. The solutions are  $x_1 = 3 + x_2/2$  and  $x_2$  arbitrary.

1.1.9.-  $h = 1$ .

1.1.10.-  $A_3 = -A_1 + 2A_2$ . The system

$$A_1x_1 + A_2x_2 + A_3x_3 = 0$$

is consistent with infinitely many solutions given by

$$x_1 = x_3, \quad x_2 = -2x_3,$$

and  $x_3$  arbitrary.

## Section 1.2: Gauss-Jordan method

1.2.1.-  $x_1 = 1, \quad x_2 = 0, \quad x_3 = 0$ .

1.2.2.-  $x_1 = 3, \quad x_2 = 1, \quad x_3 = 0$ .

1.2.3.-  $x_1 = 1, \quad x_2 = 0, \quad x_3 = -1$ .

1.2.4.-  $x_1 = 2, \quad x_2 = 4, \quad x_3 = 5$ .

$$x_1 = -4, \quad x_2 = -7, \quad x_3 = -8.$$

1.2.5.-  $x_1 = 1, \quad x_2 = 1, \quad x_3 = 1$ .

$$x_1 = 1, \quad x_2 = 2, \quad x_3 = 2.$$

$$x_1 = 1, \quad x_2 = 2, \quad x_3 = 3.$$

## Section 1.3: Echelon forms

1.3.1.-  $\text{rank}(A) = 2$ , pivot columns are the first and fourth.

1.3.2.-  $x_1 = -x_4, x_2 = 0, x_3 = -2x_4$  and  $x_4$  is a free variable.

1.3.3.-  $x_1 = -1 + 2x_3, x_2 = 2 - 3x_3$  and  $x_3$  is a free variable.

1.3.4.- For example:

$$A = \begin{bmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \quad \mathbf{c} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.$$

1.3.5.- Matrix  $A$  has a reduced echelon form  $E_A$  with a pivot on every column, therefore  $E_A$  has no row with all coefficients zero. Hence, any system with such coefficient matrix  $A$  is always consistent.

1.3.6.-

(a)  $k = -1$ .

(b)  $\mathbf{x} = \begin{bmatrix} -1/2 \\ 1 \end{bmatrix}$ .

## Section 1.4: Non-homogeneous equations

$$1.4.1.- \quad x = \begin{bmatrix} -2 \\ 1 \\ 0 \\ 0 \end{bmatrix} x_2 + \begin{bmatrix} -1 \\ 0 \\ -1 \\ 1 \end{bmatrix} x_4.$$

$$1.4.2.- \quad x = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} -3 \\ -2 \\ 1 \end{bmatrix} x_3.$$

1.4.3.- Hint: Recall that the matrix-vector product is a linear operation.

$$1.4.4.- \quad x = \begin{bmatrix} -2 \\ 1 \\ 0 \\ 0 \end{bmatrix} x_2 + \begin{bmatrix} -1 \\ 0 \\ -1 \\ 1 \end{bmatrix} x_4 + \begin{bmatrix} 1 \\ 0 \\ 2 \\ 0 \end{bmatrix}.$$

1.4.5.-  $x = \begin{bmatrix} 5 \\ -2 \\ 0 \end{bmatrix} + \begin{bmatrix} 4 \\ -7 \\ 1 \end{bmatrix} x_3$ . Writing the solution as  $x = y + zx_3$ , the solution vectors have end points on line passing through the end point of  $y$  and the line is tangent to the line parallel to  $z$ .

1.4.6.-  $x = \begin{bmatrix} 0 \\ 8 \\ 2 \\ 0 \end{bmatrix} + \begin{bmatrix} 3 \\ 4 \\ -5 \\ 1 \end{bmatrix} x_4$ . Writing the solution as  $x = y + zx_4$ , the solution vectors have end points on line passing through the end point of  $y$  and the line is tangent to the line parallel to  $z$ .

1.4.7.-

- (a) For  $k \neq 3$  there exists a unique solution.  
 (b) If  $k = 3$  the solutions are

$$x = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ -3/2 \\ 1 \end{bmatrix} x_3.$$

## Section 1.5: Floating-point numbers

1.5.1.-

- (a)  $x_1 = 0, x_2 = -1$ .  
 (b) **Hint:** Modify the coefficients of the original system so that the approximate solution is an exact solution of the modified system.

One answer is: We keep the first equation coefficients and we modify the sources:

$$\begin{aligned} 10^{-3}x_1 - x_2 &= 1 \\ x_1 + x_2 &= -1. \end{aligned}$$

- (c)  $x_1 = 1, x_2 = -1$ .  
 (d) One answer is: We keep the first equation coefficients and we modify the sources:

$$\begin{aligned} 10^{-3}x_1 - x_2 &= 1 + 10^{-3} \\ x_1 + x_2 &= 0. \end{aligned}$$

- (e)  $x_1 = \frac{1}{1 + 10^{-3}}, x_2 = \frac{-1}{1 + 10^{-3}}$ .  
 (f)  $x_1 = 0.999, x_2 = -0.999$ .

1.5.2.-

- (a)  $x_1 = 0, x_2 = 1$ .  
 (b)  $x_1 = 2, x_2 = 1$ .  
 (c)  $x_1 = 2, x_2 = 1$ .  
 (d)  $x_1 = 2, x_2 = 1$ .

(e)  $x_1 = \frac{2}{1 + 10^{-3}},$

$$x_2 = \frac{1 + 3 \times 10^{-3}}{1 + 10^{-3}}.$$

1.5.3.- Without partial pivoting:

$$x_1 = 1.01, \quad x_2 = 1.03.$$

With partial pivoting:

$$x_1 = 1, \quad x_2 = 1.$$

Solution in  $\mathbb{R}$ :

$$x_1 = 1, \quad x_2 = 1.$$

**Chapter 2: Matrix algebra****Section 2.1: Linear transformations****2.1.1.-**

- (a)  $T$  is not linear.
- (b)  $T$  is linear.
- (c)  $T$  is not linear.

$$2.1.2.- \quad x = \begin{bmatrix} -1 \\ 3 \end{bmatrix}.$$

$$2.1.3.- \quad x = \begin{bmatrix} 3 \\ 1 \\ 0 \end{bmatrix} + \begin{bmatrix} -3 \\ -2 \\ 1 \end{bmatrix} x_3.$$

$$2.1.4.- \quad T\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} 3x_1 + x_2 \\ x_1 + 3x_2 \end{bmatrix}.$$

**2.1.5.-**

- (a) Rotation by an angle  $\pi$ .
- (b) Expansion by 2.
- (c) Projection onto  $x_2$  axis.
- (d) Reflection along the line  $x_1 = x_2$ .

$$2.1.6.- \quad A = \begin{bmatrix} -2 & 7 \\ 5 & -3 \end{bmatrix}.$$

**Section 2.2: Linear combinations****2.2.1.-**

$$(a) \quad A = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

$$(b) \quad A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 5 & 4 \\ 3 & 4 & 6 \end{bmatrix}.$$

$$(c) \quad A = \begin{bmatrix} 1 & 2+i & 3 \\ 2-i & 5 & 4-i \\ 3 & 4+i & 6 \end{bmatrix}.$$

$$2.2.2.- \quad x = -1/2, y = -6, z = 0.$$

**2.2.3.-** Hint: Compute the transpose of  $A + A^T$ . Do the same with  $A - A^T$ .

**2.2.4.-** Hint: Express  $A = C + D$ , with  $C$  symmetric and  $D$  skew-symmetric. Find the expressions of matrices  $C$  and  $D$  in terms of matrix  $A$ .

**2.2.5.-**

- (a) Hint: Introduce  $i = j$  in the skew-symmetric matrix component condition.
- (b) Hint: Introduce  $i = j$  in the skew-Hermitian matrix component condition.
- (c) Hint: Compute  $B^T$ .

**2.2.6.-** Hint: Use the properties of the adjoint (complex conjugate and transpose) operation.

**2.2.7.-** Hint: See the proof of Theorem 2.2.8.

### Section 2.3: Matrix multiplication

#### 2.3.1.-

(a) BA, CB not possible,

$$AB = \begin{bmatrix} 10 & 15 \\ 12 & 8 \\ 28 & 52 \end{bmatrix}, \quad C^T B = [10 \quad 31].$$

(b)  $B^2$  is not possible, and

$$A^2 = \begin{bmatrix} 13 & -1 & 19 \\ 16 & 13 & 12 \\ 36 & -17 & 64 \end{bmatrix},$$

$$C^T C = [14], \quad CC^T = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{bmatrix}.$$

#### 2.3.2.-

(a)  $[A, B] = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}.$

(b)  $ABC = \begin{bmatrix} 9 & 26 \\ 12 & 33 \end{bmatrix}.$

#### 2.3.3.-

$$A^2 = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad A^3 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

2.3.4.-  $A^n = \begin{bmatrix} 1 & na \\ 0 & 1 \end{bmatrix}.$

2.3.5.- Hint: Expand  $(A + B)^2$  and recall that  $AB = BA$  iff  $[A, B] = 0$ .

2.3.6.- If  $B = a \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ , then

$$A = \begin{bmatrix} I_3 & \vdots & B \\ \dots & \dots & \dots \\ 0 & \vdots & B \end{bmatrix}.$$

Recalling that  $a = 1/3$ , we obtain,

$$A^n = \begin{bmatrix} I_3 & \vdots & 2^{(n-1)}B \\ \dots & \dots & \dots \\ 0 & \vdots & B \end{bmatrix}.$$

2.3.7.- Hint: Write down in components each side of the equation.

2.3.8.- Hint: Express  $\text{tr}(A^T A)$  in components, and show that it is the sum of squares of every component of matrix A.

2.3.9.- Hint: Recall that AB is symmetric iff  $(AB)^T = AB$ .

2.3.10.- Hint: Take the trace on each side of  $[A, X] = I$  and use the properties of the trace to get a contradiction.

**Section 2.4: Inverse matrix****2.4.1.-**

(a)  $A^{-1} = \begin{bmatrix} 3 & -2 \\ -1 & 1 \end{bmatrix}.$

(b)  $A^{-1} = \begin{bmatrix} -2 & 3 & -10 \\ 0 & -1 & 3 \\ -1 & 1 & -4 \end{bmatrix}.$

(c) A is not invertible.

**2.4.2.-**  $k = 2$  or  $k = -3$ .**2.4.3.-**

(a)  $A^{-1} = \begin{bmatrix} 1/2 & 1/2 & 1/2 \\ -1 & 0 & 0 \\ -1 & 0 & -1 \end{bmatrix}.$

(b)  $x = \begin{bmatrix} 3/2 \\ -1 \\ -4 \end{bmatrix}.$

**2.4.4.-** Hint: Just write down the commutator of two matrices:

$$[A, B] = AB - BA,$$

with  $B = A^{-1}$ .

**2.4.5.-**  $X = \begin{bmatrix} 2 & 4 \\ -1 & -2 \\ 3 & 3 \end{bmatrix}.$

**2.4.6.-** Hint: Use that

$$(A^{-1})^T = (A^T)^{-1}.$$

**2.4.7.-** Hint: Generalize for matrices the identity for real numbers,

$$(1 - a)(1 + a) = 1 - a^2.$$

**2.4.8.-** Hint: Generalize for matrices the identity for real numbers,

$$(1 - a)(1 + a + a^2) = 1 - a^3.$$

**2.4.9.-**

(a) Hint: Use Theorem 2.4.3.

(b) Hint: Use Theorem 2.4.3.

(c) Hint:  $\text{tr}(AB) = \text{tr}(BA)$ .**2.4.10.-** Hint: Multiply the matrices using block multiplication.

## Section 2.5: Null and Range spaces

2.5.1.-

$$R(\mathbf{A}) = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}\right\}\right),$$

$$R(\mathbf{A}^T) = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 2 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 2 \\ 4 \\ 1 \\ 3 \end{bmatrix}\right\}\right).$$

2.5.2.-

$$N(\mathbf{A}) = \text{Span}\left(\left\{\begin{bmatrix} -2 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 2 \\ 0 \\ -3 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \\ -4 \\ 0 \\ 1 \end{bmatrix}\right\}\right),$$

$$R(\mathbf{A}) = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix}\right\}\right),$$

$$N(\mathbf{A}^T) = \text{Span}\left(\left\{\begin{bmatrix} -2 \\ -1/2 \\ 1 \end{bmatrix}\right\}\right),$$

$$R(\mathbf{A}^T) = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 2 \\ 1 \\ 1 \\ 5 \end{bmatrix}, \begin{bmatrix} -2 \\ -4 \\ 0 \\ 4 \\ -2 \end{bmatrix}\right\}\right).$$

2.5.3.-

- (a) The system  $\mathbf{Ax} = \mathbf{b}$  is consistent since  $\mathbf{b} \in R(\mathbf{A})$ , because of

$$\mathbf{b} = \begin{bmatrix} 1 \\ 8 \\ 5 \end{bmatrix} = 3 \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} - 2 \begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix} \in R(\mathbf{A}).$$

- (b) Since  $N(\mathbf{A}) \neq \{\mathbf{0}\}$ , given any solution  $\mathbf{x}$  of the system  $\mathbf{Ax} = \mathbf{b}$ , then another solution is

$$\tilde{\mathbf{x}} = \mathbf{x} + c \begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix},$$

for any  $c \in \mathbb{R}$ .

2.5.4.-

- (a) A linear system  $\mathbf{Ax} = \mathbf{b}$  is consistent iff  $\mathbf{b}$  is a linear combination of the columns of  $\mathbf{A}$  iff  $\mathbf{b} \in R(\mathbf{A})$ .  
 (b) Hint:  $\mathbf{A}(\mathbf{x} + \hat{\mathbf{x}}) = \mathbf{Ax}$  for every vector  $\hat{\mathbf{x}} \in N(\mathbf{A})$ .

2.5.5.- Theorem 2.4.3 part (c) says that  $\mathbf{A}$  is invertible iff  $\mathbf{E}_A = \mathbf{I}$ , which is equivalent to say that  $R(\mathbf{A}) = \mathbb{F}^n$ .

2.5.6.-

$$N(\mathbf{A}) = N(\mathbf{A}^T) = \{\mathbf{0}\},$$

$$R(\mathbf{A}) = R(\mathbf{A}^T) = \mathbb{F}^n.$$

2.5.7.-

- (a) No, since  $\mathbf{E}_{A^T} \neq \mathbf{E}_{B^T}$ .  
 (b) Yes, since  $\mathbf{E}_A = \mathbf{E}_B$ .  
 (c) Yes, since  $\mathbf{E}_A = \mathbf{E}_B$ .  
 (d) No, since  $\mathbf{E}_{A^T} \neq \mathbf{E}_{B^T}$ .



**Section 2.6: LU-factorization****2.6.1.-**  $A = LU$  where

$$L = \begin{bmatrix} 1 & 0 \\ -3 & 1 \end{bmatrix}, U = \begin{bmatrix} 5 & 2 \\ 0 & 3 \end{bmatrix}.$$

**2.6.2.-**  $A = LU$  where

$$L = \begin{bmatrix} 1 & 0 \\ 2 & 1 \end{bmatrix}, U = \begin{bmatrix} 2 & 1 & 3 \\ 0 & 4 & 1 \end{bmatrix}.$$

**2.6.3.-**  $A = LU$  where

$$L = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 3 & -2 & 1 \end{bmatrix}, U = \begin{bmatrix} 2 & 1 & 2 \\ 0 & 3 & 1 \\ 0 & 0 & 1 \end{bmatrix}.$$

**2.6.4.-** Matrix  $A$  does not have an LU-factorization.**2.6.5.-**  $T = LU$  where

$$L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ -1/2 & 1 & 0 & 0 \\ 0 & -2/3 & 1 & 0 \\ 0 & 0 & -3/4 & 1 \end{bmatrix},$$

$$U = \begin{bmatrix} 2 & -1 & 0 & 0 \\ 0 & 3/2 & -1 & 0 \\ 0 & 0 & 4/3 & -1 \\ 0 & 0 & 0 & 1/4 \end{bmatrix}.$$

**2.6.6.-**  $A = LU$  where

$$L = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 3 & 4 & 1 \end{bmatrix}, U = \begin{bmatrix} 2 & 2 & 2 \\ 0 & 3 & 3 \\ 0 & 0 & 4 \end{bmatrix}.$$

Then,

$$y = \begin{bmatrix} 12 \\ 0 \\ -24 \end{bmatrix}, x = \begin{bmatrix} 6 \\ 6 \\ -6 \end{bmatrix}.$$

**2.6.7.-**  $c = 0$  or  $c = \pm\sqrt{2}$ .

## Chapter 3: Determinants

### Section 3.1: Definitions and properties

**3.1.1.-**  $\det(A) = 8$  and  $\det(B) = -8$ .

**3.1.2.-**  $V = 14$ .

**3.1.3.-**

$$\det(U) = (1)(4)(6) = 24,$$

$$\det(L) = (1)(3)(6) = 18.$$

**3.1.4.-**  $\det(A) = -8$ .

**3.1.5.-** For example:

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, B = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}.$$

On the one hand  $\det(A + B) = 1$ , while on the other hand  $\det(A) = \det(B) = 0$ .

**3.1.6.-**

$$\det(2A) = 2^n \det(A)$$

$$\det(-A) = (-1)^n \det(A)$$

$$\det(A^2) = [\det(A)]^2.$$

**3.1.7.-**

$$\begin{aligned} \det(B) &= \det(P^{-1}AP) \\ &= \det(P^{-1}) \det(A) \det(P) \\ &= \frac{1}{\det(P)} \det(A) \det(P) \\ &= \det(A). \end{aligned}$$

**3.1.8.-**

$$\begin{aligned} \det(A^*) &= \det(\overline{A^T}) \\ &= \det(\overline{A}) \\ &= \overline{\det(A)}. \end{aligned}$$

**3.1.9.-**

$$\begin{aligned} \det(A^*A) &= \det(A^*) \det(A) \\ &= \overline{\det(A)} \det(A) \\ &= |\det(A)|^2 \geq 0. \end{aligned}$$

**3.1.10.-**

$$\begin{aligned} \det(kA) &= |kA_1, \dots, kA_n| \\ &= k^n |A_1, \dots, A_n| \\ &= k^n \det(A). \end{aligned}$$

**3.1.11.-**

$$\begin{aligned} \det(A) &= \det(A^T) \\ &= \det(-A) \\ &= (-1)^n \det(A) \\ &= -\det(A) \Rightarrow 2 \det(A) = 0. \end{aligned}$$

Therefore,  $\det(A) = 0$ .

**3.1.12.-**

$$\begin{aligned} 1 &= \det(I_n) \\ &= \det(A^T A) \\ &= \det(A^T) \det(A) \\ &= [\det(A)]^2 \Rightarrow \det(A) = \pm 1. \end{aligned}$$

### Section 3.2: Applications

**3.2.1.-**  $A^{-1} = \frac{1}{8} \begin{bmatrix} 0 & 1 & -1 \\ -8 & 4 & 4 \\ 16 & -6 & -2 \end{bmatrix}.$

**3.2.2.-**

$$(A^{-1})_{12} = \frac{8}{41}, \quad (A^{-1})_{32} = -\frac{19}{41}.$$

**3.2.3.-**  $\det(A) = 10$  and  $\det(B) = 0$ .

**3.2.4.-** Hint:

$$\begin{vmatrix} 1 & a & a^2 \\ 1 & b & b^2 \\ 1 & c & c^2 \end{vmatrix} = \begin{vmatrix} 1 & a & a^2 \\ 0 & (b-a) & (b^2-a^2) \\ 0 & 0 & (c-a)(c-b) \end{vmatrix}.$$

**3.2.5.-**  $k \neq \pm 1$ .

**3.2.6.-**  $x = \frac{1}{(ad-bc)} \begin{bmatrix} d \\ -c \end{bmatrix}.$

**3.2.7.-**  $x = \begin{bmatrix} -2 \\ 4 \\ -1 \end{bmatrix}.$

**Chapter 4: Vector Spaces****Section 4.1: Spaces and subspaces****4.1.1.-**

- (a) No.
- (b) Yes.
- (c) No.
- (d) Yes.
- (e) No.
- (f) No.

**4.1.2.-**

- (a) Yes.
- (b) No.
- (c) No.
- (d) Yes.
- (e) No.
- (f) Yes.
- (g) Yes.

**4.1.3.-** The result is  $W_1 + W_2 = \mathbb{R}^3$ . An example is  $W_1 = \text{Span}(\{\mathbf{e}_1, \mathbf{e}_2\})$ , and  $W_2 = \text{Span}(\{\mathbf{e}_3\})$ . Then,

$$W_1 + W_2 = \text{Span}(\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3\}) = \mathbb{R}^3.$$

**4.1.4.-**

- (a) The line  $\text{Span}\left(\left\{\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}\right\}\right)$ .
- (b) The plane  $\text{Span}\left(\left\{\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}\right\}\right)$ .
- (c) The space  $\mathbb{R}^3$ .

**4.1.5.-** If  $S_1 = \{\mathbf{u}_1, \dots, \mathbf{u}_k\}$  and  $S_2 = \{\mathbf{v}_1, \dots, \mathbf{v}_l\}$ , then

$$\text{Span}(S_1 \cup S_2) =$$

$\{(a_1 \mathbf{u}_1 + \dots + a_k \mathbf{u}_k) + (b_1 \mathbf{v}_1 + \dots + b_l \mathbf{v}_l)\}$  for every set of scalars  $a_1, \dots, a_k \in \mathbb{F}$ ,  $b_1, \dots, b_l \in \mathbb{F}$ . From the definition of sum of subspaces we see that this last subspace can be rewritten as

$$\text{Span}(S_1 \cup S_2) =$$

$$\{a_1 \mathbf{u}_1 + \dots + a_k \mathbf{u}_k\} + \{b_1 \mathbf{v}_1 + \dots + b_l \mathbf{v}_l\} \\ = \text{Span}(S_1) + \text{Span}(S_2).$$

**4.1.6.-** Find a vector not in  $\text{Span}(W_1)$ .

Denoting  $\mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix}$ , we compute

$$\begin{bmatrix} 1 & 1 & u_1 \\ 2 & 0 & u_2 \\ 3 & 1 & u_3 \end{bmatrix} \rightarrow$$

$$\rightarrow \begin{bmatrix} 1 & 1 & u_1 \\ 0 & -2 & u_2 - 2u_1 \\ 0 & 0 & u_3 - u_1 - u_2 \end{bmatrix}.$$

Choose any solution of  $u_3 - u_1 - u_2 \neq 0$ .

For example,  $\mathbf{u} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ , and

$$W_2 = \text{Span}(\{\mathbf{u}\}).$$

**Section 4.2: Linear dependence**

**4.2.1.-**

- (a) Linearly dependent.  
 (b) Linearly dependent. (A four vector set in  $\mathbb{R}^3$  is always linearly dependent.)  
 (c) Linearly independent.

**4.2.2.-**

- (a) One possible solution is:

$$\left\{ \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 0 \\ 2 \\ 3 \end{bmatrix} \right\}.$$

- (b) 11 sets.

**4.2.3.-** If  $S = \{\mathbf{0}, \mathbf{v}_1, \dots, \mathbf{v}_k\}$ , then

$$1\mathbf{0} + 0\mathbf{v}_1 + \dots + 0\mathbf{v}_k = \mathbf{0}.$$

**4.2.4.-** Hint: Show that the constants  $c_1, c_2$  satisfying

$$c_1(\mathbf{v} + \mathbf{w}) + c_2(\mathbf{v} - \mathbf{w}) = \mathbf{0}$$

are both zero iff the constants  $\tilde{c}_1, \tilde{c}_2$  satisfying

$$\tilde{c}_1 \mathbf{v} + \tilde{c}_2 \mathbf{w} = \mathbf{0}$$

are both zero. Once that is proven, and since the latter statement is true, so is the former.

**4.2.5.-** Linearly dependent.

**4.2.6.-** Hint: Use the definition. Prove that there are non-zero constants  $c_1, c_2, c_3, c_4$  solutions of

$$c_1 + c_2c + c_3x^2 + c_4(1 + x + x^2) = 0.$$

**4.2.7.-** Linearly independent.

**Section 4.3: Bases and dimension**

**4.3.1.-** Bases are not unique. Here are possible choices:

$$\mathcal{N}_A = \left\{ \begin{bmatrix} -2 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \\ -1 \\ 1 \end{bmatrix} \right\},$$

$$\mathcal{R}_A = \left\{ \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} \right\},$$

$$\mathcal{N}_{A^T} = \left\{ \begin{bmatrix} 1/3 \\ -5/3 \\ 1 \end{bmatrix} \right\},$$

$$\mathcal{R}_{A^T} = \left\{ \begin{bmatrix} 1 \\ 2 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 2 \\ 4 \\ 1 \\ 3 \end{bmatrix} \right\}.$$

**4.3.2.-** The dimension is 3.

**4.3.3.-**

- (a)  $\dim \mathbb{P}_n = n + 1$ .
- (b)  $\dim \mathbb{F}^{m,n} = mn$ .
- (c)  $\dim(\text{Sym}_n) = n(n + 1)/2$ .
- (d)  $\dim(\text{SkewSym}_n) = n(n - 1)/2$ .

**4.3.4.-** One example is the following:  
 $\mathbf{v}_1 = \mathbf{e}_1$ ,  $\mathbf{v}_2 = \mathbf{e}_2$ , and

$$W = \text{Span}\left(\left\{\begin{bmatrix} 1 \\ 1 \end{bmatrix}\right\}\right).$$

**4.3.5.-** To verify that  $\mathbf{v} \in N(A)$  compute  $A\mathbf{v}$  and check that the result is 0. A basis for  $N(A)$  containing  $\mathbf{v}$  is

$$\mathcal{N} = \left\{ \begin{bmatrix} -8 \\ 1 \\ 3 \\ 3 \\ 0 \end{bmatrix}, \begin{bmatrix} -2 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \\ -2 \\ 0 \\ 1 \end{bmatrix} \right\}.$$

**4.3.6.-** The set  $\mathcal{B}$  is a basis of that subspace.

## Section 4.4: Vector components

4.4.1.-

$$v_u = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}.$$

4.4.2.-

$$v_u = \begin{bmatrix} 9 \\ 2 \\ -3 \end{bmatrix}.$$

4.4.3.-

$$(a) r_s = \begin{bmatrix} 2 \\ 3 \\ -1 \end{bmatrix}.$$

$$(b) r_q = \begin{bmatrix} 6 \\ -4 \\ -1 \end{bmatrix}.$$

4.4.4.-

- (a) Hint: Prove that  $\text{Span}(\mathcal{M}) = \mathbb{R}^{2,2}$ , and that  $\mathcal{M}$  is linearly independent. For the first part, show that for every matrix

$$B = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}$$

there exist constants  $c_1, c_2, c_3, c_4$  solution of

$$c_1M_1 + c_2M_2 + c_3M_3 + c_4M_4 = B.$$

For the second part, choose  $B = 0$  and show that the only constants  $c_1, c_2, c_3, c_4$  solution of

$$c_1M_1 + c_2M_2 + c_3M_3 + c_4M_4 = 0$$

are  $c_1 = c_2 = c_3 = c_4 = 0$ .

(b)

$$A = \frac{5}{2}M_1 + \frac{1}{2}M_2 + \frac{5}{2}M_3 - \frac{3}{2}M_4.$$

Using the notation

$$D_M = \begin{bmatrix} d_1 & d_2 \\ d_3 & d_4 \end{bmatrix} \Leftrightarrow$$

$$D = d_1M_1 + d_2M_2 + d_3M_3 + d_4M_4$$

then we get

$$A_M = \frac{1}{2} \begin{bmatrix} 5 & 1 \\ 5 & -3 \end{bmatrix}.$$

**Chapter 5: Linear transformations****Section 5.1: Linear transformations**

**5.1.1.-** The projection  $A$  neither injective nor surjective.

**5.1.2.-** The rotation  $R(\theta)$  is both injective and surjective for  $\theta \in [0, 2\pi)$ .

**5.1.3.-**

- (a) Linear.
- (b) Linear.
- (c) Non-linear.
- (d) Linear.
- (e) Linear.

**5.1.4.-** Hint: Show that

$$T(ax + by) = aT(x) + bT(y)$$

for all  $x, y \in \mathbb{F}^n$  and all  $a, b \in \mathbb{F}$ .

$T$  is not a linear operator, it is a linear functional.

**5.1.5.-** Hint: Show that

$$\Delta(ap + bq) = a\Delta(p) + b\Delta(q)$$

for all  $p, q \in \mathbb{P}_3$  and all  $a, b \in \mathbb{F}$ .

$\Delta$  is not injective, but it is surjective.

**5.1.6.-** Hint: Recall the definition of a linearly independent set.

**5.1.7.-** Hint: Recall the Nullity-Rank Theorem 5.1.6.

**5.1.8.-** Hint: Recall the Nullity-Rank Theorem 5.1.6.

**Section 5.2: The inverse transformations**

**5.2.1.-** Hint: Show that  $T$  is bijective.

The inverse transformation is

$$T^{-1}\left(\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}\right) = \frac{1}{5} \begin{bmatrix} 2y_1 + y_2 \\ 3y_1 - y_2 \end{bmatrix}.$$

**5.2.2.-** Hint: Show that  $\Delta$  is bijective.

The inverse transformation is

$$\Delta^{-1}(b_0 + b_1x + b_2x^2) = \frac{b_0}{2}x^2 + \frac{b_1}{6}x^3 + \frac{b_2}{12}x^4.$$

**5.2.3.-** Hint: Choose any isomorphism  $T : V \rightarrow W$  and use the Nullity-Rank Theorem.

**5.2.4.-** Hint: Use the coordinate map to find an isomorphism between these spaces.

**5.2.5.-** Hint: Generalize the proof in Exercise 5.2.4.

**5.2.6.-** Hint: Generalize the proof in Example 5.2.8.

**5.2.7.-** Hint: Use the ideas given in Example 5.2.9.

**5.2.8.-** The basis  $\{\phi_1, \phi_2\}$  is given by

$$\phi_1\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = -x_1 + 3x_2,$$

$$\phi_2\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = x_1 - 2x_2.$$

Using row vector notation,

$$[\phi_1] = [-1, 3],$$

$$[\phi_2] = [1, -2].$$

### Section 5.3: The algebra of linear operators

#### 5.3.1.-

$$(a) (3\mathbf{T} - 2\mathbf{S})\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} 3x_2 \\ x_1 \end{bmatrix}.$$

$$(b) (\mathbf{T} \circ \mathbf{S})\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1 \\ 0 \end{bmatrix}.$$

$$(\mathbf{S} \circ \mathbf{T})\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} 0 \\ x_2 \end{bmatrix}.$$

$$(c) \mathbf{T}^2\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}.$$

$$\mathbf{S}^2\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

5.3.2.-  $\dim L(\mathbb{F}^4) = 16$ ,  $\dim L(\mathbb{P}_2) = 9$ ,  
 $\dim L(\mathbb{F}^{3,2}) = 36$ .

#### 5.3.3.-

$$\mathbf{T}^{-2}\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1/4 \\ -x_1 + x_2 \end{bmatrix}.$$

#### 5.3.4.-

$$p(\mathbf{T})\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} 2x_1 + 6x_2 \\ 9x_1 + 11x_2 \end{bmatrix}.$$

5.3.5.- Hint: Use that  $\cos^2(\alpha) + \sin^2(\alpha) = 1$ .

5.3.6.- Hint: Use the definition of the commutator.

### Section 5.4: Transformation components

$$5.4.1.- \mathsf{T}_{uu} = \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}.$$

#### 5.4.2.-

$$(a) \mathsf{T}_{uu} = \frac{1}{2} \begin{bmatrix} 2 & -3 & 1 \\ -2 & 1 & 1 \\ 0 & 1 & -1 \end{bmatrix}.$$

(b) To verify  $[\mathbf{T}(v)]_u = \mathsf{T}_{uu}v_u$  one must compute the left-hand side and the right-hand side, and check one obtains the same result. For the left-hand side:

$$[\mathbf{T}(v_s)]_s = \begin{bmatrix} 0 \\ 0 \\ -1 \end{bmatrix}_s$$

and

$$\begin{bmatrix} 0 \\ 0 \\ -1 \end{bmatrix}_s = \frac{1}{2} \begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix}_u.$$

For the right-hand side:

$$v_u = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}_u,$$

and  $\mathsf{T}_{uu}v_u$  is given by

$$\frac{1}{2} \begin{bmatrix} 2 & -3 & 1 \\ -2 & 1 & 1 \\ 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}_u = \frac{1}{2} \begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix}_u.$$

5.4.3.- Hint: Recall the definition

$$\mathsf{T}_{s\bar{s}} = [[\mathbf{T}(e_1)]_{\bar{s}}, \dots, [\mathbf{T}(e_n)]_{\bar{s}}].$$

If  $\mathbf{A} = [\mathbf{A}_1, \dots, \mathbf{A}_n]$ , then for  $i = 1, \dots, n$  holds  $[\mathbf{T}(e_i)]_{\bar{s}} = \mathbf{A}e_i = \mathbf{A}_i$ . Therefore,

$$\mathsf{T}_{s\bar{s}} = \mathbf{A}.$$

$$5.4.4.- \mathsf{T}_{s\bar{s}} = \begin{bmatrix} 0 & -1 & 2 & 0 \\ 0 & 0 & -2 & 6 \\ 0 & 0 & 0 & -3 \end{bmatrix}.$$

#### 5.4.5.-

$$(\mathbf{D} \circ \mathbf{S})_{s\bar{s}} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

$$(\mathbf{S} \circ \mathbf{D})_{s\bar{s}} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$



**Section 5.5: Change of basis****5.5.1.-**

(a)  $\mathbf{l}_{us} = \begin{bmatrix} 2 & 1 \\ -9 & 8 \end{bmatrix},$

$$\mathbf{l}_{su} = \frac{1}{25} \begin{bmatrix} 8 & -1 \\ 9 & 2 \end{bmatrix}.$$

(b)  $\mathbf{x}_s = \begin{bmatrix} 5 \\ -10 \end{bmatrix}, \mathbf{x}_u = \begin{bmatrix} 2 \\ 1 \end{bmatrix}.$

**5.5.2.-**

(a)  $\mathbf{l}_{bc} = \frac{1}{9} \begin{bmatrix} -1 & -6 & -2 \\ 2 & 3 & 4 \\ -5 & -3 & -1 \end{bmatrix},$

$$\mathbf{l}_{cb} = \begin{bmatrix} 1 & 0 & -2 \\ -2 & -1 & 0 \\ 1 & 3 & 1 \end{bmatrix}.$$

(b)  $\mathbf{x}_b = \begin{bmatrix} -3 \\ 0 \\ -3 \end{bmatrix}, \mathbf{x}_c = \begin{bmatrix} 1 \\ -2 \\ 2 \end{bmatrix}.$

**5.5.3.-**

(a)  $\mathbf{x}_s = \begin{bmatrix} 7 \\ 8 \end{bmatrix}.$

(b)  $\mathbf{e}_{1u} = \frac{1}{3} \begin{bmatrix} -1 \\ 2 \end{bmatrix}, \mathbf{e}_{2u} = \frac{1}{3} \begin{bmatrix} 2 \\ -1 \end{bmatrix}.$

**5.5.4.-**

(a)  $\mathbf{x}_s = \begin{bmatrix} -1 \\ 5 \end{bmatrix}.$

(b)  $\mathbf{x}_b = \frac{1}{4} \begin{bmatrix} 3 \\ -7 \end{bmatrix}.$

**5.5.5.-**  $\mathbf{b}_{1s} = \begin{bmatrix} -3 \\ 2 \end{bmatrix}, \mathbf{b}_{2s} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}.$

**5.5.6.-** Hint: Find an invertible matrix  $\mathbf{P}$  such that  $\mathbf{C} = \mathbf{P}^{-1}\mathbf{A}\mathbf{P}$ .**5.5.7.-**

$$\mathbf{T}_{ss} = \begin{bmatrix} 1 & 2 & -1 \\ 0 & -1 & 0 \\ 1 & 0 & 7 \end{bmatrix},$$

$$\mathbf{T}_{uu} = \begin{bmatrix} 1 & 4 & 3 \\ -1 & -2 & -9 \\ 1 & 1 & 8 \end{bmatrix}.$$

**5.5.8.-**

$$\mathbf{T}_{us} = \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix}, \mathbf{T}_{ss} = \begin{bmatrix} 2 & -1 \\ 2 & 1 \end{bmatrix},$$

$$\mathbf{T}_{uu} = \begin{bmatrix} 2 & 2 \\ -1 & 1 \end{bmatrix}, \mathbf{T}_{su} = \begin{bmatrix} 2 & 0 \\ 0 & -1 \end{bmatrix}.$$

## Chapter 6: Inner product spaces

### Section 6.1: The dot product

**6.1.1.-**  $\|u\| = 5$ ,  $\|v\| = 2$ ,  
 $d(u, v) = \sqrt{31}$ ,  $\theta = \arccos(-1/10)$ .

**6.1.2.-**  
 $u = \frac{1}{\sqrt{13}} \begin{bmatrix} -2 \\ 3 \end{bmatrix}$ ,  $v = \frac{1}{\sqrt{13}} \begin{bmatrix} 2 \\ -3 \end{bmatrix}$ .

**6.1.3.-**  $u = \frac{1}{\sqrt{10}} \begin{bmatrix} 1 + 2i \\ 2 - i \end{bmatrix}$ .

**6.1.4.-**  
 (a)  $\left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right\}$ .

(b)  $\left\{ \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right\}$ .

**6.1.5.-** Hint: Expand the expression

$$\|x - y\|^2 = (x - y) \cdot \overline{(x - y)}.$$

**6.1.6.-**

(a) Hint: For every  $x, y \in \mathbb{R}^n$  holds

$$\operatorname{Re}(x \cdot y) = 0 \Leftrightarrow x \cdot y = 0.$$

(b) Hint: For  $x, y \in \mathbb{C}^n$  satisfying  $\operatorname{Im}(x \cdot y) \neq 0$  holds

$$\operatorname{Re}(x \cdot y) = 0 \not\Leftrightarrow x \cdot y = 0.$$

**6.1.7.-** Hint: Use Problem 6.1.5.

### Section 6.2: Inner product

**6.2.1.-**

(a) No.

(b) No.

(c) Yes.

(d) No.

**6.2.2.-**

(a) Hint: Choose a particular  $x$ .

(b) Hint: Use the properties of the inner product.

**6.2.3.-**

(a) Yes.

(b) No.

(c) No.

**6.2.4.-** Hint:  $N(A) = \{0\}$  is needed to show the positivity property of the inner product.

**6.2.5.-**  $k = -3$ .

**6.2.6.-**  $\|A\|_F = \sqrt{10}$ ,  $\|B\|_F = \sqrt{3}$ .

**6.2.7.-** Hint: Use properties of the trace operation.

**6.2.8.-**  $q = \pm \frac{1}{2}(3 - 5x^2)$ .

**Section 6.3: Orthogonal vectors**

**6.3.1.-** Hint: ( $\Leftrightarrow$ ): choose particular values for  $a, b$ . Choose  $b = 1, a = \langle x, y \rangle$ .

**6.3.2.-**  $x \in \text{Span}\left(\left\{\begin{bmatrix} -2 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} -3 \\ 0 \\ 1 \end{bmatrix}\right\}\right)$ .

**6.3.3.-**  $x \in \text{Span}\left(\left\{\begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix}\right\}\right)$ .

**6.3.4.-** Hint: Compute all inner products  $\langle p_i, p_j \rangle$  for  $i, j = 0, 1, 2, 3$ .

**6.3.5.-**

(a) Show that  $U^T U = I_3$ , where the columns in  $U$  are the basis vectors in  $\mathcal{U}$ .

(b)  $[x]_{\mathcal{U}} = Ux = \frac{1}{\sqrt{6}} \begin{bmatrix} \sqrt{3} \\ -\sqrt{2} \\ -5 \end{bmatrix}$ .

**6.3.6.-**

(a) Hint: Compute all inner products  $\langle E_i, E_j \rangle$  for  $i, j = 1, 2, 3, 4$ .

(b) Hint: Use the definition,  $[A]_i = \langle M_i, A \rangle_F$ .

**6.3.7.-**  $\theta = \pi/2$ .

**6.3.8.-**  $U_{:3} = \frac{1}{\sqrt{42}} \begin{bmatrix} -5 \\ 4 \\ 1 \end{bmatrix}$ .

**Section 6.4: Orthogonal projections**

**6.4.1.-**  $x = \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} + \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$ .

**6.4.2.-**

(a)  $w_2 = \begin{bmatrix} 1/2 \\ 1 \\ 1/2 \end{bmatrix} + \begin{bmatrix} 3/2 \\ -2 \\ 5/2 \end{bmatrix}$ .

(b)  $x = \frac{5}{3} \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} + \frac{1}{3} \begin{bmatrix} 7 \\ -1 \\ -5 \end{bmatrix}$ .

**6.4.3.-**  $x = \frac{1}{3} \begin{bmatrix} 2 \\ 1 \\ -4 \end{bmatrix} + \frac{1}{3} \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}$ .

**6.4.4.-**  $\mathcal{R}^\perp = \left\{ \begin{bmatrix} 2 \\ -4 \\ 1 \end{bmatrix} \right\}$ .

**6.4.5.-**

(a)  $\mathcal{W}^\perp = \left\{ \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} -2 \\ 0 \\ 1 \end{bmatrix} \right\}$ .

(b)  $\tilde{\mathcal{W}}^\perp = \left\{ \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} \right\}$ .

**6.4.6.-**  $\mathcal{W}^\perp = \left\{ \begin{bmatrix} -2 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} -3 \\ 0 \\ 0 \\ 1 \end{bmatrix} \right\}$ .

**6.4.7.-**

(a) Hint: Use the definition of the orthogonal complement of a subspace.

(b) Hint: Show  $(X + Y)^\perp \subset X^\perp \cap Y^\perp$  using (a). Show the other inclusion,  $X^\perp \cap Y^\perp \subset (X + Y)^\perp$  using the definition of orthogonal complement of a subspace.

(c) Hint: Use part (b), that is,  $(\tilde{X} + \tilde{Y})^\perp = \tilde{X}^\perp \cap \tilde{Y}^\perp$ , for  $\tilde{X} = X^\perp$  and  $\tilde{Y} = Y^\perp$ .

**Section 6.5: Gram-Schmidt method**

$$6.5.1.- \left\{ \frac{1}{3} \begin{bmatrix} -2 \\ 2 \\ -1 \end{bmatrix}, \frac{1}{\sqrt{2}} \begin{bmatrix} -1 \\ -1 \\ 0 \end{bmatrix} \right\}.$$

6.5.2.-

$$(a) \left\{ \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \frac{1}{5} \begin{bmatrix} 3 \\ 0 \\ 4 \end{bmatrix} \right\}.$$

$$(b) \mathbf{x} = \frac{1}{5} \begin{bmatrix} 9 \\ 5 \\ 12 \end{bmatrix} + \frac{4}{5} \begin{bmatrix} 4 \\ 0 \\ -3 \end{bmatrix}.$$

6.5.3.-

$$\left\{ \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, \frac{1}{\sqrt{6}} \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix}, \frac{1}{\sqrt{3}} \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix} \right\}.$$

$$6.5.4.- \left\{ \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, \frac{1}{\sqrt{6}} \begin{bmatrix} 1 \\ 2 \\ -1 \end{bmatrix} \right\}.$$

$$6.5.5.- \left\{ \mathbf{q}_0 = 1, \mathbf{q}_1 = -\frac{1}{2} + x, \mathbf{q}_2 = \frac{1}{6} - x + x^2 \right\}.$$

**Section 6.6: The adjoint operator**

6.6.1.-

6.6.2.-

**Chapter 7: Approximation methods**

**Section 7.1: Best approximation**

7.1.1.-

7.1.2.-

**Section 7.2: Least squares**

7.2.1.-

7.2.2.-

**Section 7.3: Finite difference method**

7.3.1.-

7.3.2.-

**Section 7.4: Finite element method**

7.4.1.-

7.4.2.-

## 10.4. SOLUTIONS TO PRACTICE EXAMS

**Solutions to Practice Exam 1.**

1. Consider the matrix  $A = \begin{bmatrix} 5 & 3 & 1 \\ 1 & 2 & -1 \\ 2 & 1 & 1 \end{bmatrix}$ . Find the coefficients  $(A^{-1})_{21}$  and  $(A^{-1})_{32}$  of the matrix  $A^{-1}$ , that is, of the inverse matrix of  $A$ . Show your work.

**SOLUTION:** The formula for the inverse matrix  $A^{-1} = C^T / \det(A)$  implies that

$$(A^{-1})_{21} = \frac{C_{12}}{\det(A)} \quad (A^{-1})_{32} = \frac{C_{23}}{\det(A)}.$$

We start computing the  $\det(A)$ , that is,

$$\begin{vmatrix} 5 & 3 & 1 \\ 1 & 2 & -1 \\ 2 & 1 & 1 \end{vmatrix} = 5 \begin{vmatrix} 2 & -1 \\ 1 & 1 \end{vmatrix} - 3 \begin{vmatrix} 1 & -1 \\ 2 & 1 \end{vmatrix} + \begin{vmatrix} 1 & 2 \\ 2 & 1 \end{vmatrix} = 15 - 9 - 3 \Rightarrow \det(A) = 3.$$

Now the cofactors,

$$C_{12} = (-1)^{(1+2)} \begin{vmatrix} 1 & -1 \\ 2 & 1 \end{vmatrix} = -3, \quad C_{23} = (-1)^{(2+3)} \begin{vmatrix} 5 & 3 \\ 2 & 1 \end{vmatrix} = 1.$$

Therefore, we conclude that

$$\boxed{(A^{-1})_{21} = -1}, \quad \boxed{(A^{-1})_{32} = \frac{1}{3}}.$$

◀

2. (a) Find  $k \in \mathbb{R}$  such that the volume of the parallelepiped formed by the vectors below is equal to 4, where

$$\mathbf{v}_1 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix}, \quad \mathbf{v}_3 = \begin{bmatrix} k \\ 1 \\ 1 \end{bmatrix}$$

- (b) Set  $k = 1$  and define the matrix  $A = [\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3]$ . Matrix  $A$  determines the linear transformation  $A : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ . Is this linear transformation injective (one-to-one)? Is it surjective (onto)?

**SOLUTION:**

(a) The volume of a parallelepiped formed with vectors  $\mathbf{v}_1$ ,  $\mathbf{v}_2$ , and  $\mathbf{v}_3$  is the absolute value of the determinant of a matrix with these vectors as column vectors. In this case that determinant is given by

$$\begin{vmatrix} 1 & 3 & k \\ 2 & 2 & 1 \\ 3 & 1 & 1 \end{vmatrix} = \begin{vmatrix} 2 & 1 \\ 1 & 1 \end{vmatrix} - 3 \begin{vmatrix} 2 & 1 \\ 3 & 1 \end{vmatrix} + k \begin{vmatrix} 2 & 2 \\ 3 & 1 \end{vmatrix} = 4 - 4k.$$

The volume of the parallelepiped is  $V = 4$  iff holds

$$4 = |4 - 4k| \Rightarrow \begin{cases} 4 = 4 - 4k & \Rightarrow \boxed{k = 0}, \\ 4 = -4 + 4k & \Rightarrow \boxed{k = 2}. \end{cases}$$

(b) If  $k = 1$ , then  $\det(A) = 0$ . This implies that the set  $\{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3\}$  is linearly dependent, so  $N(A) \neq \{\mathbf{0}\}$ , which implies that  $A$  is **not injective**. This matrix is also **not surjective**, since the Nullity-Rank Theorem implies

$$3 = \dim N(A) + \dim R(A) \geq 1 + \dim R(A) \Rightarrow \dim R(A) \leq 2.$$

This establishes that  $R(A) \neq \mathbb{R}^3$ .

◀

3. Determine whether the subset  $V \subset \mathbb{R}^3$  is a subspace, where

$$V = \left\{ \begin{bmatrix} -a + b \\ a - 2b \\ a - 7b \end{bmatrix} \text{ with } a, b \in \mathbb{R} \right\}.$$

If the set is a subspace, find an orthogonal basis in the inner product space  $(\mathbb{R}^3, \cdot)$ .

**SOLUTION:** Vectors in  $V$  can be written as follows,

$$\begin{bmatrix} -a + b \\ a - 2b \\ a - 7b \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} a + \begin{bmatrix} 1 \\ -2 \\ -7 \end{bmatrix} b$$

which implies that

$$V = \text{Span} \left( \left\{ \mathbf{v}_1 = \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix}, \mathbf{v}_2 = \begin{bmatrix} 1 \\ -2 \\ -7 \end{bmatrix} \right\} \right) \Rightarrow \boxed{V \text{ is a subspace.}}$$

A basis for this subspace  $V$  is the set  $\{\mathbf{v}_1, \mathbf{v}_2\}$  since these vectors are not collinear. This basis is not orthogonal though, because

$$\mathbf{v}_1 \cdot \mathbf{v}_2 = -10 \neq 0.$$

We construct an orthogonal basis for  $V$  by projecting  $\mathbf{v}_2$  onto  $\mathbf{v}_1$ . That is, let us find the vector

$$\mathbf{v}_{2\perp} = \mathbf{v}_2 - \frac{\mathbf{v}_1 \cdot \mathbf{v}_2}{\|\mathbf{v}_1\|^2} \mathbf{v}_1 = \begin{bmatrix} 1 \\ -2 \\ -7 \end{bmatrix} - \frac{(-10)}{3} \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix} = \frac{1}{3} \left( \begin{bmatrix} 3 \\ -6 \\ -21 \end{bmatrix} + \begin{bmatrix} -10 \\ 10 \\ 10 \end{bmatrix} \right) \Rightarrow \mathbf{v}_{2\perp} = \frac{1}{3} \begin{bmatrix} -7 \\ 4 \\ -11 \end{bmatrix}.$$

Since any non-zero vector proportional  $\mathbf{v}_{2\perp}$  is orthogonal to  $\mathbf{v}_1$ , an orthogonal basis for the subspace  $V$  is

$$\left\{ \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} -7 \\ 4 \\ -11 \end{bmatrix} \right\}.$$

◁

4. True or False: (Justify your answers.)

- (a) If the set of columns of  $A \in \mathbb{F}^{m,n}$  is a linearly independent set, then  $A\mathbf{x} = \mathbf{b}$  has exactly one solution for every  $\mathbf{b} \in \mathbb{F}^m$ .  
 (b) The set of column vectors of an  $5 \times 7$  is never linearly independent.

**SOLUTION:**

(a) **False.** The reason why this statement is false lies in the part of the sentence “for every  $\mathbf{b} \in \mathbb{F}^m$ ”. Do not get confused by the “exactly one solution” part. We now give an example of a system where the set of coefficient column vectors is linearly independent but the system has no solution. Take  $A \in \mathbb{R}^{3,2}$  and  $\mathbf{b} \in \mathbb{R}^3$  given by,

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}.$$

In this case  $A\mathbf{x} = \mathbf{b}$  has no solution.

(b) **True.** The biggest linearly independent set in  $\mathbb{R}^5$  contains five vectors, so any set containing seven vectors is linearly dependent. ◁

5. Consider the vector space  $\mathbb{R}^2$  with the standard basis  $\mathcal{S}$  and let  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be the linear transformation

$$[T]_{ss} = \begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix}_{ss}.$$

Find  $[\mathbf{T}]_{bb}$ , the matrix of  $\mathbf{T}$  in the basis  $\mathcal{B}$ , where  $\mathcal{B} = \left\{ [\mathbf{b}_1]_s = \begin{bmatrix} 1 \\ 2 \end{bmatrix}_s, [\mathbf{b}_2]_s = \begin{bmatrix} 1 \\ -2 \end{bmatrix}_s \right\}$ .

**SOLUTION:** From the change of basis formulas we know that

$$[\mathbf{T}]_{bb} = \mathbf{P}^{-1}[\mathbf{T}]_{ss}\mathbf{P}, \quad \mathbf{P} = \mathbf{I}_{bs} = \begin{bmatrix} 1 & 1 \\ 2 & -2 \end{bmatrix}, \quad \mathbf{P}^{-1} = \frac{1}{4} \begin{bmatrix} 2 & 1 \\ 2 & -1 \end{bmatrix}.$$

Therefore, we obtain

$$[\mathbf{T}]_{bb} = \frac{1}{4} \begin{bmatrix} 2 & 1 \\ 2 & -1 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 2 & 3 \end{bmatrix}_{ss} \begin{bmatrix} 1 & 1 \\ 2 & -2 \end{bmatrix} \Rightarrow \boxed{[\mathbf{T}]_{bb} = \frac{1}{2} \begin{bmatrix} 9 & -5 \\ 1 & -1 \end{bmatrix}}.$$

◀

6. Consider the linear transformations  $\mathbf{T}: \mathbb{R}^3 \rightarrow \mathbb{R}^2$  and  $\mathbf{S}: \mathbb{R}^3 \rightarrow \mathbb{R}^3$  given by

$$\left[ \mathbf{T} \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{s_3} \right) \right]_{s_2} = \begin{bmatrix} x_1 - x_2 + x_3 \\ -x_1 + 2x_2 + x_3 \end{bmatrix}_{s_2}, \quad \left[ \mathbf{S} \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{s_3} \right) \right]_{s_3} = \begin{bmatrix} 3x_3 \\ 2x_2 \\ x_1 \end{bmatrix}_{s_3},$$

where  $\mathcal{S}_3$  and  $\mathcal{S}_2$  are the standard basis of  $\mathbb{R}^3$  and  $\mathbb{R}^2$ , respectively.

- Find a matrix  $[\mathbf{T}]_{s_3s_2}$  and the matrix  $[\mathbf{S}]_{s_3s_3}$ . Show your work.
- Find the matrix of the composition  $\mathbf{T} \circ \mathbf{S}: \mathbb{R}^3 \rightarrow \mathbb{R}^2$  in the standard basis, that is, find  $[\mathbf{T} \circ \mathbf{S}]_{s_3s_2}$ .
- Is  $\mathbf{T} \circ \mathbf{S}$  injective (one-to-one)? Is  $\mathbf{T} \circ \mathbf{S}$  surjective (onto)? Justify your answer.

**SOLUTION:**

- From the definition of the matrix of a linear transformation we get,

$$\boxed{[\mathbf{T}]_{s_3s_2} = \begin{bmatrix} 1 & -1 & 1 \\ -1 & 2 & 1 \end{bmatrix}}, \quad \boxed{[\mathbf{S}]_{s_3s_3} = \begin{bmatrix} 0 & 0 & 3 \\ 0 & 2 & 0 \\ 1 & 0 & 0 \end{bmatrix}}.$$

- Recalling the formula  $[\mathbf{T} \circ \mathbf{S}]_{s_3s_2} = [\mathbf{T}]_{s_3s_2}[\mathbf{S}]_{s_3s_3}$ , we obtain

$$[\mathbf{T} \circ \mathbf{S}]_{s_3s_2} = \begin{bmatrix} 1 & -1 & 1 \\ -1 & 2 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 3 \\ 0 & 2 & 0 \\ 1 & 0 & 0 \end{bmatrix} \Rightarrow \boxed{[\mathbf{T} \circ \mathbf{S}]_{s_3s_2} = \begin{bmatrix} 1 & -2 & 3 \\ 1 & 4 & -3 \end{bmatrix}}.$$

- We now obtain the reduced echelon form of matrix  $[\mathbf{T} \circ \mathbf{S}]_{s_3s_2}$  to find out whether this composition is injective or surjective.

$$[\mathbf{T} \circ \mathbf{S}]_{s_3s_2} = \begin{bmatrix} 1 & -2 & 3 \\ 1 & 4 & -3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -2 & 3 \\ 0 & 6 & -6 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \end{bmatrix}.$$

Since  $\dim N(\mathbf{T} \circ \mathbf{S}) = 3 - \text{rank}(\mathbf{T} \circ \mathbf{S}) = 1$ , we conclude that  $\mathbf{T} \circ \mathbf{S}$  is **not injective**. Since  $\dim R(\mathbf{T} \circ \mathbf{S}) = 2 = \dim \mathbb{R}^2$  and  $R(\mathbf{T} \circ \mathbf{S}) \subset \mathbb{R}^2$ , we conclude that  $R(\mathbf{T} \circ \mathbf{S}) = \mathbb{R}^2$ , which shows that  $\mathbf{T} \circ \mathbf{S}$  is **surjective**. ◀

7. Consider the matrix  $\mathbf{A} = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 2 & 1 \end{bmatrix}$  and the vector  $\mathbf{b} = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}$ .

- Find the least-squares solution  $\hat{\mathbf{x}}$  to the matrix equation  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .
- Verify whether the vector  $\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}$  belong to the space  $R(\mathbf{A})^\perp$ ? Justify your answers.

**SOLUTION:**



(a) We find the normal equation  $A^T A \hat{x} = A^T b$  for  $\hat{x}$  and we solve it as follows.

$$A^T A = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 2 & 1 \end{bmatrix} \Rightarrow A^T A = \begin{bmatrix} 6 & 5 \\ 5 & 6 \end{bmatrix}.$$

$$A^T b = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} \Rightarrow A^T b = 5 \begin{bmatrix} 1 \\ 1 \end{bmatrix}.$$

The normal equation is

$$\begin{bmatrix} 6 & 5 \\ 5 & 6 \end{bmatrix} \hat{x} = 5 \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow \hat{x} = \frac{1}{(36-25)} \begin{bmatrix} 6 & -5 \\ -5 & 6 \end{bmatrix} 5 \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow \boxed{\hat{x} = \frac{5}{11} \begin{bmatrix} 1 \\ 1 \end{bmatrix}}.$$

(b) We first compute the vector

$$(A\hat{x} - b) = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 2 & 1 \end{bmatrix} \frac{5}{11} \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} = \frac{5}{11} \begin{bmatrix} 2 \\ 3 \\ 3 \end{bmatrix} - \frac{11}{11} \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} = \frac{1}{11} \begin{bmatrix} -12 \\ 4 \\ 4 \end{bmatrix}.$$

We conclude that

$$(A\hat{x} - b) = \frac{4}{11} \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}.$$

Now it is simple to verify that  $(A\hat{x} - b) \in R(A)^\perp$ , since

$$\begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix} = 0, \quad \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} = 0 \Rightarrow \boxed{(A\hat{x} - b) \in R(A)^\perp}.$$

◁

8. Consider the matrix  $A = \begin{bmatrix} -1/2 & -3 \\ 1/2 & 2 \end{bmatrix}$ .

(a) Show that matrix  $A$  is diagonalizable.

(b) Using that  $A$  is diagonalizable, find the  $\lim_{k \rightarrow \infty} A^k$ .

**SOLUTION:**

(a) We need to compute the eigenvalues and eigenvectors of matrix  $A$ . The eigenvalues are the solutions of the characteristic equation

$$p(\lambda) = \begin{vmatrix} (-\frac{1}{2} - \lambda) & -3 \\ \frac{1}{2} & (2 - \lambda) \end{vmatrix} = (\lambda - 2)\left(\lambda + \frac{1}{2}\right) + \frac{3}{2} = \lambda^2 - \frac{3}{2}\lambda + \frac{1}{2} = 0.$$

We then obtain,

$$\boxed{\lambda_+ = 1}, \quad \boxed{\lambda_- = \frac{1}{2}}.$$

The eigenvector corresponding to  $\lambda_+ = 1$  is a non-zero solution of  $(A - I_2)v^+ = 0$ . The Gauss method implies,

$$\begin{bmatrix} -\frac{3}{2} & -3 \\ \frac{1}{2} & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 \\ 0 & 0 \end{bmatrix} \Rightarrow v_1^+ = -2v_2^+ \Rightarrow \boxed{v^+ = \begin{bmatrix} -2 \\ 1 \end{bmatrix}}.$$

The eigenvector corresponding to  $\lambda_- = 1/2$  is a non-zero solution of  $(A - I_2)v^- = 0$  and it is computed as follows,

$$\begin{bmatrix} -1 & -3 \\ \frac{1}{2} & \frac{3}{2} \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 3 \\ 0 & 0 \end{bmatrix} \Rightarrow v_1^- = -3v_2^- \Rightarrow \boxed{v^- = \begin{bmatrix} -3 \\ 1 \end{bmatrix}}.$$

WE now know that matrix  $A$  is diagonalizable since  $\{v^+, v^-\}$  is linearly independent. Introducing matrices  $P$  and  $D$  as follows,

$$P = \begin{bmatrix} -2 & -3 \\ 1 & 1 \end{bmatrix}, \quad D = \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{2} \end{bmatrix},$$

we conclude that  $A = P D P^{-1}$ , that is,

$$A = \begin{bmatrix} -2 & -3 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 1 & 3 \\ -1 & -2 \end{bmatrix}.$$

(b) Knowing that  $A$  is diagonalizable it is simple to compute  $\lim_{k \rightarrow \infty} A^k$ . The calculation is

$$\begin{aligned} \lim_{k \rightarrow \infty} A^k &= \lim_{k \rightarrow \infty} \begin{bmatrix} -2 & -3 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \left(\frac{1}{2}\right)^k \end{bmatrix} \begin{bmatrix} 1 & 3 \\ -1 & -2 \end{bmatrix} \\ &= \begin{bmatrix} -2 & -3 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \\ &= \begin{bmatrix} -2 & 0 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ -1 & -2 \end{bmatrix} \Rightarrow \boxed{\lim_{k \rightarrow \infty} A^k = \begin{bmatrix} -2 & -6 \\ 1 & 3 \end{bmatrix}}. \end{aligned}$$

◀

**9.** Let  $(V, \langle \cdot, \cdot \rangle)$  be an inner product space with inner product norm  $\|\cdot\|$ . Let  $T: V \rightarrow V$  be a linear transformation and  $x, y \in V$  be vectors satisfying the following conditions:

$$T(x) = 2x, \quad T(y) = -3y, \quad \|x\| = 1/3, \quad \|y\| = 1, \quad x \perp y.$$

- (a) Compute  $\|v\|$  for the vector  $v = 3x - y$ .  
 (b) Compute  $\|T(v)\|$  for the vector  $v$  given above.

**SOLUTION:**

(a)

$$\begin{aligned} \|v\|^2 &= \|3x - y\|^2 \\ &= \langle (3x - y), (3x - y) \rangle \\ &= 9\|x\|^2 + \|y\|^2 - 3\langle x, y \rangle - 3\langle y, x \rangle \\ &= 9\|x\|^2 + \|y\|^2 \\ &= 9\left(\frac{1}{3}\right)^2 + 1 \Rightarrow \boxed{\|v\| = \sqrt{2}}. \end{aligned}$$

(b) We first compute  $T(v) = T(3x - y) = 3T(x) - T(y)$ . Recall that  $T(x) = 2x$ ,  $T(y) = -3y$ , we conclude that  $T(v) = 6x - 3y$ . Then, it is simple to see that

$$\begin{aligned} \|T(v)\|^2 &= \|6x - 3y\|^2 \\ &= \langle (6x - 3y), (6x - 3y) \rangle \\ &= 36\|x\|^2 + 9\|y\|^2 \\ &= 36\left(\frac{1}{3}\right)^2 + 9 \\ &= 4 + 9 \Rightarrow \boxed{\|T(v)\| = \sqrt{13}}. \end{aligned}$$

◀

**10.** Consider the matrix  $A = \begin{bmatrix} 2 & -1 & 2 \\ 0 & 1 & h \\ 0 & 0 & 2 \end{bmatrix}$ .

- (a) Find all eigenvalues of matrix  $A$  and their corresponding algebraic multiplicities.

- (b) Find the value(s) of the real number  $h$  such that the matrix  $A$  above has a two-dimensional eigenspace, and find a basis for this eigenspace.

**SOLUTION:**

(a) Since matrix  $A$  is upper triangular, its eigenvalues are the diagonal elements. Denoting by  $\lambda$  the eigenvalues and  $r$  their corresponding algebraic multiplicities we obtain,

$$\boxed{\lambda_1 = 2, \quad r_1 = 2}, \quad \text{and} \quad \boxed{\lambda_2 = 1, \quad r_2 = 1}.$$

(b) Since the eigenvalue  $\lambda_2 = 1$  has algebraic multiplicity  $r_2 = 1$ , this eigenvalue cannot have a two-dimensional eigenspace, since  $\dim E_{\lambda_2} \leq r_2 = 1$ . The only candidate is the eigenspace  $E_{\lambda_1}$ . Let us compute the non-zero elements in this eigenspace,

$$A - 2I_3 = \begin{bmatrix} 0 & -1 & 2 \\ 0 & -1 & h \\ 0 & 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 1 & -2 \\ 0 & 0 & h-2 \\ 0 & 0 & 0 \end{bmatrix}.$$

We look for the value of  $h$  such that the eigenspace  $E_{\lambda_1}$  is two-dimensional, which means that the reduced echelon form associated with the matrix above has rank equal to one. This means that this matrix has only one pivot, which in turns is a condition on  $h$ . The condition is

$$\boxed{h = 2}.$$

In this case, the eigenspace is given by the solutions of the homogeneous linear system with the coefficient matrix

$$\begin{bmatrix} 0 & 1 & -2 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \Rightarrow \begin{cases} x_2 = -2x_3 \\ x_1, x_3 \text{ free.} \end{cases} \Rightarrow \mathbf{x} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} x_2 + \begin{bmatrix} 0 \\ -2 \\ 1 \end{bmatrix} x_3.$$

Therefore, a basis for  $E_{\lambda_1}$  is given by

$$\boxed{\left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -2 \\ 1 \end{bmatrix} \right\}}.$$

◁

**Solutions to Practice Exam 2.**

1. Consider the matrix  $A = \begin{bmatrix} -2 & 3 & -1 \\ 1 & 2 & -1 \\ -2 & -1 & 1 \end{bmatrix}$ . Find the coefficients  $(A^{-1})_{13}$  and  $(A^{-1})_{21}$  of the inverse matrix of  $A$ . Show your work.

**SOLUTION:** The formula for the inverse matrix  $A^{-1} = C^T / \det(A)$  implies that

$$(A^{-1})_{13} = \frac{C_{31}}{\det(A)} \quad (A^{-1})_{21} = \frac{C_{12}}{\det(A)}.$$

We start computing the  $\det(A)$ , that is,

$$\begin{vmatrix} -2 & 3 & -1 \\ 1 & 2 & -1 \\ -2 & -1 & 1 \end{vmatrix} = -2 \begin{vmatrix} 2 & -1 \\ -1 & 1 \end{vmatrix} - 3 \begin{vmatrix} 1 & -1 \\ -2 & 1 \end{vmatrix} - \begin{vmatrix} 1 & 2 \\ -2 & -1 \end{vmatrix} = -2 + 3 - 3,$$

which gives us the result  $\det(A) = -2$ . Now the cofactors,

$$C_{31} = (-1)^{(3+1)} \begin{vmatrix} 3 & -1 \\ 2 & -1 \end{vmatrix} = -1, \quad C_{12} = (-1)^{(1+2)} \begin{vmatrix} 1 & -1 \\ -2 & 1 \end{vmatrix} = 1.$$

Therefore, we conclude that

$$\boxed{(A^{-1})_{13} = \frac{1}{2}}, \quad \boxed{(A^{-1})_{12} = -\frac{1}{2}}.$$

◀

2. Consider the vector space  $\mathbb{P}_3([0, 1])$  with the inner product

$$\langle \mathbf{p}, \mathbf{q} \rangle = \int_0^1 \mathbf{p}(x)\mathbf{q}(x) dx.$$

Given the set  $\mathcal{U} = \{\mathbf{p}_1 = x^2, \mathbf{p}_2 = x^3\}$ , find an orthogonal basis for the subspace  $U = \text{Span}(\mathcal{U})$  using the Gram-Schmidt method on the set  $\mathcal{U}$  starting with the vector  $\mathbf{p}_1$ .

**SOLUTION:** We call the orthogonal basis  $\{\mathbf{q}_1, \mathbf{q}_2\}$ . Stating with vector  $\mathbf{p}_1$  means  $\mathbf{q}_1 = \mathbf{p}_1$ , hence

$$\boxed{\mathbf{q}_1(x) = x^2}.$$

The vector  $\mathbf{q}_2$  is found with the formula

$$\mathbf{q}_2 = \mathbf{p}_2 - \frac{\langle \mathbf{q}_1, \mathbf{p}_2 \rangle}{\|\mathbf{q}_1\|^2} \mathbf{q}_1.$$

We need to compute the scalars

$$\|\mathbf{q}_1\|^2 = \int_0^1 x^4 dx = \frac{x^5}{5} \Big|_0^1 \Rightarrow \|\mathbf{q}_1\|^2 = \frac{1}{5}.$$

$$\langle \mathbf{q}_1, \mathbf{p}_2 \rangle = \int_0^1 x^5 dx = \frac{x^6}{6} \Big|_0^1 \Rightarrow \langle \mathbf{q}_1, \mathbf{p}_2 \rangle = \frac{1}{6}.$$

Therefore, we conclude that

$$\boxed{\mathbf{q}_2 = x^3 - \frac{5}{6}x^2}.$$

◀

3. Consider the matrix  $A = \begin{bmatrix} 1 & 3 & 1 & 1 \\ 2 & 6 & 3 & 0 \\ 3 & 9 & 5 & -1 \end{bmatrix}$ .

- (a) Verify that the vector  $\mathbf{v} = \begin{bmatrix} 3 \\ 1 \\ -4 \\ -2 \end{bmatrix}$  belongs to the null space of  $\mathbf{A}$ .
- (b) Extend the set  $\{\mathbf{v}\}$  into a basis of the null space of  $\mathbf{A}$ .

**SOLUTION:**

- (a) The vector  $\mathbf{v}$  belongs to  $N(\mathbf{A})$ , since

$$\mathbf{A}\mathbf{v} = \begin{bmatrix} 1 & 3 & 1 & 1 \\ 2 & 6 & 3 & 0 \\ 3 & 9 & 5 & -1 \end{bmatrix} \begin{bmatrix} 3 \\ 1 \\ -4 \\ -2 \end{bmatrix} = \begin{bmatrix} 3+3-4-2 \\ 6+6-12+0 \\ 9+9-20+2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \Rightarrow \boxed{\mathbf{A}\mathbf{v} = \mathbf{0}}.$$

- (b) We first need to find a basis for the  $N(\mathbf{A})$ . So we find all solutions to the homogeneous equation  $\mathbf{A}\mathbf{x} = \mathbf{0}$ . We perform Gauss operations on matrix  $\mathbf{A}$ ;

$$\mathbf{A} \rightarrow \begin{bmatrix} 1 & 3 & 1 & 1 \\ 0 & 0 & 1 & -2 \\ 0 & 0 & 2 & -4 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 3 & 0 & 3 \\ 0 & 0 & 1 & -2 \\ 0 & 0 & 0 & 0 \end{bmatrix} \stackrel{3R}{\rightarrow} \begin{cases} x_1 = -3x_2 - 3x_4, \\ x_3 = 2x_4. \end{cases}$$

Therefore, the solution  $\mathbf{x}$  has the form

$$\mathbf{x} = \begin{bmatrix} -3 \\ 1 \\ 0 \\ 0 \end{bmatrix} x_2 + \begin{bmatrix} -3 \\ 0 \\ 2 \\ 1 \end{bmatrix} x_4 \Rightarrow \mathcal{U} = \left\{ \begin{bmatrix} -3 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} -3 \\ 0 \\ 2 \\ 1 \end{bmatrix} \right\} \text{ is a basis of } N(\mathbf{A}).$$

The following calculations is useful to find a basis of  $N(\mathbf{A})$  containing  $\mathbf{v}$ :

$$\begin{bmatrix} 3 & -3 & -3 \\ 1 & 1 & 0 \\ -4 & 0 & 2 \\ 2 & 0 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -1 & -1 \\ 1 & 1 & 0 \\ -4 & 0 & 2 \\ 2 & 0 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -1 & -1 \\ 0 & 2 & 1 \\ 0 & -4 & 2 \\ 0 & -2 & 3 \end{bmatrix}.$$

From the last matrix we know that the pivot columns are the first and second columns. That says that the set  $\mathcal{N}$  below form a basis to the  $N(\mathbf{A})$ , where

$$\boxed{\mathcal{N} = \left\{ \begin{bmatrix} 3 \\ 1 \\ -4 \\ 2 \end{bmatrix}, \begin{bmatrix} -3 \\ 1 \\ 0 \\ 0 \end{bmatrix} \right\}}.$$

This set includes vector  $\mathbf{v}$ .

◁

4. Use Cramer's rule to find the solution to the linear system

$$\begin{aligned} 2x_1 + x_2 - x_3 &= 0 \\ x_1 + x_3 &= 1 \\ x_1 + 2x_2 + 3x_3 &= 0. \end{aligned}$$

**SOLUTION:** If we write this system in the form  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , then

$$\mathbf{A} = \begin{bmatrix} 2 & 1 & -1 \\ 1 & 0 & 1 \\ 1 & 2 & 3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}.$$

Following Cramer's rule, the components of the solution vector  $\mathbf{x} = [x_i]$ , for  $i = 1, 2, 3$ , are given by  $x_i = \det(\mathbf{A}_i) / \det(\mathbf{A})$ . Since

$$\det(\mathbf{A}) = \begin{vmatrix} 2 & 1 & -1 \\ 1 & 0 & 1 \\ 1 & 2 & 3 \end{vmatrix} = -(3-1) - 2(2+1) \Rightarrow \det(\mathbf{A}) = -8,$$

we obtain,

$$\begin{aligned} x_1 &= \frac{1}{(-8)} \begin{vmatrix} 2 & 1 & -1 \\ 1 & 0 & 1 \\ 1 & 2 & 3 \end{vmatrix} = \frac{-(3+2)}{(-8)} & \Rightarrow & x_1 = \frac{5}{8}; \\ x_2 &= \frac{1}{(-8)} \begin{vmatrix} 2 & 1 & -1 \\ 1 & 0 & 1 \\ 1 & 2 & 3 \end{vmatrix} = \frac{(6+1)}{(-8)} & \Rightarrow & x_2 = -\frac{7}{8}; \\ x_3 &= \frac{1}{(-8)} \begin{vmatrix} 2 & 1 & -1 \\ 1 & 0 & 1 \\ 1 & 2 & 3 \end{vmatrix} = \frac{-(4-1)}{(-8)} & \Rightarrow & x_3 = \frac{3}{8}. \end{aligned}$$

We conclude that

$$\mathbf{v} = \frac{1}{8} \begin{bmatrix} 5 \\ -7 \\ 3 \end{bmatrix}.$$

◁

5. Let  $\mathcal{S}_3$  and  $\mathcal{S}_2$  be standard bases of  $\mathbb{R}^3$  and  $\mathbb{R}^2$ , respectively, and consider the linear transformation  $\mathbf{T}: \mathbb{R}^3 \rightarrow \mathbb{R}^2$  given by

$$\left[ \mathbf{T} \left( \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right) \right]_{s_2} = \begin{bmatrix} -x_1 + 2x_2 - x_3 \\ x_1 + x_3 \end{bmatrix}_{s_2},$$

and introduce the bases  $\mathcal{U} \subset \mathbb{R}^3$  and  $\mathcal{V} \subset \mathbb{R}^2$  given by

$$\begin{aligned} \mathcal{U} &= \left\{ [\mathbf{u}_1]_{s_3} = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}_{s_3}, [\mathbf{u}_2]_{s_3} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}_{s_3}, [\mathbf{u}_3]_{s_3} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}_{s_3} \right\}, \\ \mathcal{V} &= \left\{ [\mathbf{v}_1]_{s_2} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}_{s_2}, [\mathbf{v}_2]_{s_2} = \begin{bmatrix} -3 \\ 2 \end{bmatrix}_{s_2} \right\}. \end{aligned}$$

Find the matrices  $[\mathbf{T}]_{s_3 s_2}$  and  $[\mathbf{T}]_{uv}$ . Show your work.

**SOLUTION:** Matrix  $[\mathbf{T}]_{s_3 s_2} = [\mathbf{T}(\mathbf{e}_1)]_{s_2}, \mathbf{T}(\mathbf{e}_2)]_{s_2}, \mathbf{T}(\mathbf{e}_3)]_{s_2}$ , that is,

$$[\mathbf{T}]_{s_3 s_2} = \begin{bmatrix} -1 & 2 & -1 \\ 1 & 0 & 1 \end{bmatrix}.$$

The matrix  $[\mathbf{T}]_{uv}$  can be computed with the change of basis formula

$$[\mathbf{T}]_{uv} = \mathbf{Q}^{-1} [\mathbf{T}]_{s_3 s_2} \mathbf{P}$$

where the change of basis matrices  $\mathbf{Q}$  and  $\mathbf{P}$  are given by

$$\mathbf{Q} = I_{vs_2} = \begin{bmatrix} 1 & -3 \\ 2 & 2 \end{bmatrix} \Rightarrow \mathbf{Q}^{-1} = I_{s_2 v} = \frac{1}{8} \begin{bmatrix} 2 & -2 \\ 3 & 1 \end{bmatrix}, \quad \mathbf{P} = I_{us_3} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix}.$$

Therefore, the matrix  $[\mathbf{T}]_{uv}$  is given by

$$\begin{aligned} [\mathbf{T}]_{uv} &= \frac{1}{8} \begin{bmatrix} 2 & -2 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} -1 & 2 & -1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix} \\ &= \frac{1}{8} \begin{bmatrix} 2 & -2 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} 1 & -2 & 1 \\ 1 & 2 & 1 \end{bmatrix} \Rightarrow \boxed{[\mathbf{T}]_{uv} = \frac{1}{8} \begin{bmatrix} 5 & 2 & 5 \\ -1 & 6 & -1 \end{bmatrix}}. \end{aligned}$$

◀

6. Consider the inner product space  $(\mathbb{R}^{2,2}, \langle \cdot, \cdot \rangle_F)$  and the subspace

$$W = \text{Span}\left\{ \mathbf{E}_1 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \mathbf{E}_2 = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \right\} \subset \mathbb{R}^{2,2}.$$

Find a basis for  $W^\perp$ , the orthogonal complement of  $W$ .

**SOLUTION:** We first find the whole subspace  $W^\perp$  and afterwards we find a basis. The matrix  $\mathbf{X} \in W^\perp$  iff  $\langle \mathbf{E}_1, \mathbf{X} \rangle_F = 0$  and  $\langle \mathbf{E}_2, \mathbf{X} \rangle_F = 0$ . If we denote

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix},$$

then the equations above have the form

$$\begin{aligned} 0 = \langle \mathbf{E}_1, \mathbf{X} \rangle_F & & 0 = \langle \mathbf{E}_2, \mathbf{X} \rangle_F \\ = \text{tr} \left( \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix} \right) & & = \text{tr} \left( \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix} \right) \\ = x_{21} + x_{12}; & & = x_{11} - x_{22}. \end{aligned}$$

We can express the solutions in the following way,

$$x_{11} = x_{22}, \quad x_{21} = -x_{12}.$$

Then, a general element  $\mathbf{X} \in W^\perp$  has the form

$$\mathbf{X} = \begin{bmatrix} x_{22} & x_{12} \\ -x_{12} & x_{22} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} x_{22} \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} x_{12}$$

for arbitrary scalars  $x_{22}, x_{12}$ . So a basis of  $W^\perp$  is given by

$$\boxed{\left\{ \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \right\}}.$$

◀

7. Consider the matrix  $\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 1 & -1 \\ -2 & 1 \end{bmatrix}$  and the vector  $\mathbf{b} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ .

- Find the least-squares solution  $\hat{\mathbf{x}}$  to the matrix equation  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .
- Verify that the vector  $\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}$ , where  $\hat{\mathbf{x}}$  is the least-squares solution found in part (7a), belongs to the space  $R(\mathbf{A})^\perp$ , the orthogonal complement of  $R(\mathbf{A})$ .

**SOLUTION:**

(a) We first compute the normal equation  $\mathbf{A}^T \mathbf{A} \hat{\mathbf{x}} = \mathbf{A}^T \mathbf{b}$ . The coefficient matrix is

$$\mathbf{A}^T \mathbf{A} = \begin{bmatrix} 1 & 1 & -2 \\ 2 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 1 & -1 \\ -2 & 1 \end{bmatrix} = \begin{bmatrix} 6 & -1 \\ -1 & 6 \end{bmatrix};$$

while the source vector is

$$\mathbf{A}^T \mathbf{b} = \begin{bmatrix} 1 & 1 & -2 \\ 2 & -1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}.$$

Therefore, the least-squares solution  $\hat{\mathbf{x}}$  is obtained as follows,

$$\begin{bmatrix} 6 & -1 \\ -1 & 6 \end{bmatrix} \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix} \Rightarrow \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \end{bmatrix} = \frac{1}{35} \begin{bmatrix} 6 & 1 \\ 1 & 6 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} \Rightarrow \boxed{\hat{\mathbf{x}} = \frac{1}{7} \begin{bmatrix} 1 \\ -1 \end{bmatrix}}.$$

(b) We first compute the vector  $(\mathbf{A}\hat{\mathbf{x}} - \mathbf{b})$ , which is given by

$$(\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}) = \begin{bmatrix} 1 & 2 \\ 1 & -1 \\ -2 & 1 \end{bmatrix} = \frac{1}{7} \begin{bmatrix} 1 \\ 1 \\ -2 \end{bmatrix} - \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \frac{1}{7} \begin{bmatrix} -1 \\ 2 \\ -3 \end{bmatrix} - \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \Rightarrow (\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}) = \frac{1}{7} \begin{bmatrix} -1 \\ -5 \\ -3 \end{bmatrix}.$$

It is simple to see that

$$\boxed{\begin{bmatrix} 1 \\ 1 \\ -2 \end{bmatrix} \cdot \begin{bmatrix} -1 \\ -5 \\ -3 \end{bmatrix} = -1 - 5 + 6 = 0}, \quad \boxed{\begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} -1 \\ -5 \\ -3 \end{bmatrix} = -2 + 5 - 3 = 0}.$$

◁

8. Suppose that a matrix  $\mathbf{A} \in \mathbb{R}^{3,3}$  has eigenvalues  $\lambda_1 = 1$ ,  $\lambda_2 = 2$ , and  $\lambda_3 = 4$ .
- Find the trace of  $\mathbf{A}$ , find the trace of  $\mathbf{A}^2$ , and find the determinant of  $\mathbf{A}$ .
  - Is matrix  $\mathbf{A}$  invertible? If your answer is “yes”, then prove it and find  $\det(\mathbf{A}^{-1})$ ; if your answer is “no”, then prove it.

**SOLUTION:** Since Matrix  $\mathbf{A}$  is  $3 \times 3$  and has three different eigenvalues, then  $\mathbf{A}$  is diagonalizable. Therefore, there exists an invertible matrix  $\mathbf{P}$  such that  $\mathbf{A} = \mathbf{P} \mathbf{D} \mathbf{P}^{-1}$ , where  $\mathbf{D} = \text{diag}[1, 2, 4]$ .

(a)

$$\text{tr}(\mathbf{A}) = \text{tr}(\mathbf{P} \mathbf{D} \mathbf{P}^{-1}) = \text{tr}(\mathbf{P}^{-1} \mathbf{P} \mathbf{D}) = \text{tr}(\mathbf{D}) = 1 + 2 + 4 \Rightarrow \boxed{\text{tr}(\mathbf{A}) = 7}.$$

$$\text{tr}(\mathbf{A}^2) = \text{tr}(\mathbf{P} \mathbf{D}^2 \mathbf{P}^{-1}) = \text{tr}(\mathbf{P}^{-1} \mathbf{P} \mathbf{D}^2) = \text{tr}(\mathbf{D}^2) = 1 + 4 + 16 \Rightarrow \boxed{\text{tr}(\mathbf{A}^2) = 21}.$$

$$\det(\mathbf{A}) = \det(\mathbf{P} \mathbf{D} \mathbf{P}^{-1}) = \det(\mathbf{P}^{-1}) \det(\mathbf{D}) \det(\mathbf{P}) = \frac{1}{\det(\mathbf{P})} \det(\mathbf{D}) \det(\mathbf{P}) = \det(\mathbf{D}).$$

Since  $\det(\mathbf{D}) = (1)(2)(4)$ , we conclude that

$$\boxed{\det(\mathbf{A}) = 8}.$$

(b) Since matrix  $\mathbf{A}$  has non-zero determinant, then it is invertible. We also know that  $\det(\mathbf{A}^{-1}) = 1/\det(\mathbf{A})$ , therefore

$$\boxed{\det(\mathbf{A}^{-1}) = \frac{1}{8}}.$$

◁

9. Consider the matrix  $\mathbf{A} = \begin{bmatrix} 7 & 5 \\ 3 & -7 \end{bmatrix}$ .

- Find the eigenvalues and eigenvectors of  $\mathbf{A}$ .
- Compute the matrix  $e^{\mathbf{A}}$ .

**SOLUTION:**

(a) The eigenvalues are the roots of the characteristic polynomial,

$$p(\lambda) = \begin{vmatrix} (7 - \lambda) & 5 \\ 3 & (-7 - \lambda) \end{vmatrix} = (\lambda - 7)(\lambda + 7) - 15 = \lambda^2 - 64.$$



Therefore, the eigenvalues are

$$\boxed{\lambda_+ = 8}, \quad \boxed{\lambda_- = -8}.$$

The corresponding eigenvectors are the following: For  $\lambda_+ = 8$  we obtain,

$$\mathbf{A} - 8\mathbf{I}_2 = \begin{bmatrix} -1 & 5 \\ 3 & -15 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -5 \\ 0 & 0 \end{bmatrix} \Rightarrow v_1^+ = 5v_2^+ \Rightarrow \boxed{\mathbf{v}^+ = \begin{bmatrix} 5 \\ 1 \end{bmatrix}}.$$

For  $\lambda_- = -8$  we obtain

$$\mathbf{A} + 8\mathbf{I}_2 = \begin{bmatrix} 15 & 5 \\ 3 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 3 & 1 \\ 0 & 0 \end{bmatrix} \Rightarrow 3v_1^- = -v_2^- \Rightarrow \boxed{\mathbf{v}^- = \begin{bmatrix} 1 \\ -3 \end{bmatrix}}.$$

(b) Since  $\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$ , then  $e^{\mathbf{A}} = \mathbf{P}e^{\mathbf{A}}\mathbf{P}^{-1}$ , where

$$\mathbf{P} = \begin{bmatrix} 5 & 1 \\ 1 & -3 \end{bmatrix} \Rightarrow \mathbf{P}^{-1} = \frac{1}{16} \begin{bmatrix} 3 & 1 \\ 1 & -5 \end{bmatrix}.$$

Therefore, it is simple to compute

$$\begin{aligned} e^{\mathbf{A}} &= \begin{bmatrix} 5 & 1 \\ 1 & -3 \end{bmatrix} \begin{bmatrix} 8 & 0 \\ 0 & -8 \end{bmatrix} \frac{1}{16} \begin{bmatrix} 3 & 1 \\ 1 & -5 \end{bmatrix}, \\ &= \frac{1}{2} \begin{bmatrix} 5 & -1 \\ 1 & 3 \end{bmatrix} \begin{bmatrix} 3 & 1 \\ 1 & -5 \end{bmatrix}, \\ &= \frac{1}{2} \begin{bmatrix} 14 & 10 \\ 6 & -14 \end{bmatrix} \Rightarrow \boxed{e^{\mathbf{A}} = \begin{bmatrix} 7 & 5 \\ 3 & -7 \end{bmatrix}}. \end{aligned}$$

◁

**10.** Find the function  $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^2$  solution of the initial value problem

$$\frac{d}{dt}\mathbf{x}(t) = \mathbf{A}\mathbf{x}(t), \quad \mathbf{x}(0) = \mathbf{x}_0,$$

where the matrix  $\mathbf{A} = \begin{bmatrix} -5 & 2 \\ -12 & 5 \end{bmatrix}$  and the vector  $\mathbf{x}_0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ .

**SOLUTION:** We first need to compute the eigenvalues and eigenvectors of matrix  $\mathbf{A}$ . The eigenvalues are the roots of the characteristic polynomial,

$$p(\lambda) = \begin{vmatrix} (-5 - \lambda) & 2 \\ -12 & (5 - \lambda) \end{vmatrix} = (\lambda + 5)(\lambda - 5) + 24 = \lambda^2 - 1.$$

Therefore, the eigenvalues are

$$\boxed{\lambda_+ = 1}, \quad \boxed{\lambda_- = -1}.$$

The corresponding eigenvectors are the following: For  $\lambda_+ = 1$  we obtain,

$$\mathbf{A} - \mathbf{I}_2 = \begin{bmatrix} -6 & 2 \\ -12 & 4 \end{bmatrix} \rightarrow \begin{bmatrix} -3 & 1 \\ 0 & 0 \end{bmatrix} \Rightarrow 3v_1^+ = v_2^+ \Rightarrow \boxed{\mathbf{v}^+ = \begin{bmatrix} 1 \\ 3 \end{bmatrix}}.$$

For  $\lambda_- = -1$  we obtain

$$\mathbf{A} + \mathbf{I}_2 = \begin{bmatrix} -4 & 2 \\ 12 & 6 \end{bmatrix} \rightarrow \begin{bmatrix} -2 & 1 \\ 0 & 0 \end{bmatrix} \Rightarrow 2v_1^- = v_2^- \Rightarrow \boxed{\mathbf{v}^- = \begin{bmatrix} 1 \\ 2 \end{bmatrix}}.$$

Since the general solution to the differential equation above is

$$\mathbf{x}(t) = c_+\mathbf{v}^+e^{\lambda_+t} + c_-\mathbf{v}^-e^{\lambda_-t},$$

we conclude that

$$\mathbf{x}(t) = c_+ \begin{bmatrix} 1 \\ 3 \end{bmatrix} e^t + c_- \begin{bmatrix} 1 \\ 2 \end{bmatrix} e^{-t}.$$

The initial condition implies that

$$\begin{bmatrix} 1 \\ 1 \end{bmatrix} = \mathbf{x}_0 = \mathbf{x}(0) = c_+ \begin{bmatrix} 1 \\ 3 \end{bmatrix} + c_- \begin{bmatrix} 1 \\ 2 \end{bmatrix}.$$

We need to solve the linear system

$$\begin{bmatrix} 1 & 1 \\ 3 & 2 \end{bmatrix} \begin{bmatrix} c_+ \\ c_- \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow \begin{bmatrix} c_+ \\ c_- \end{bmatrix} = \frac{1}{(-1)} \begin{bmatrix} 2 & -1 \\ -3 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow \begin{bmatrix} c_+ \\ c_- \end{bmatrix} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}.$$

So the solution to the initial value problem above is

$$\mathbf{x}(t) = - \begin{bmatrix} 1 \\ 3 \end{bmatrix} e^t + 2 \begin{bmatrix} 1 \\ 2 \end{bmatrix} e^{-t}.$$

◀

**Solutions to Practice Exam 3.**

1. (a) Find the LU-factorization of the matrix  $A = \begin{bmatrix} 1 & 2 & 3 \\ -1 & 2 & 2 \\ 2 & -8 & -3 \end{bmatrix}$ .
- (b) Use the LU-factorization above to find the solution of the linear system  $Ax = b$ , where  $b = \begin{bmatrix} 1 \\ -2 \\ -1 \end{bmatrix}$ .

**SOLUTION:**

(a) We use Gauss operations to transform matrix  $A$  into upper triangular, as follows,

$$A = \begin{bmatrix} 1 & 2 & 3 \\ -1 & 2 & 2 \\ 2 & -8 & -3 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 0 & 4 & 5 \\ 0 & -12 & -9 \end{bmatrix} \rightarrow \boxed{\begin{bmatrix} 1 & 2 & 3 \\ 0 & 4 & 2 \\ 0 & 0 & 6 \end{bmatrix}} = U.$$

The factors used to construct matrix  $U$  define the lower triangular matrix  $L$  as follows,

$$L = \boxed{\begin{bmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 2 & -3 & 1 \end{bmatrix}}.$$

(b) We find the solution  $x$  of  $Ax = b$  using the LU-factorization  $A = LU$  as follows:

$$LUx = b \Leftrightarrow \begin{cases} Ly = b, \\ Ux = y. \end{cases}$$

We first find vector  $y$  using forward substitution,

$$\begin{bmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 2 & -3 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} 1 \\ -2 \\ -1 \end{bmatrix} \Rightarrow y = \begin{bmatrix} 1 \\ -1 \\ -6 \end{bmatrix}.$$

We now solve for vector  $x$  using back substitution,

$$\begin{bmatrix} 1 & 2 & 3 \\ 0 & 4 & 5 \\ 0 & 0 & 6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \\ -6 \end{bmatrix} \Rightarrow \boxed{x = \begin{bmatrix} 2 \\ 1 \\ -1 \end{bmatrix}}.$$

◀

2. Determine whether the following sets  $W_1$  and  $W_2$  are subspaces of the vector space  $\mathbb{P}_2([0, 1])$ . If your answer is “yes,” find a basis of the subspace.

(a)  $W_1 = \{p \in \mathbb{P}_2([0, 1]) : \int_0^1 p(x) dx \leq 1\};$

(b)  $W_2 = \{p \in \mathbb{P}_2([0, 1]) : \int_0^1 x p(x) dx = 0\}.$

**SOLUTION:**

(a)  $W_1$  is **not** a subspace. Take any element  $p \in W_1$  satisfying  $\int_0^1 p(x) dx < -1$ . Multiply that element by  $(-1)$ . The result is not in  $W_1$ , since  $\int_0^1 -p(x) dx > 1$ .

(b)  $W_2$  is a subspace. Proof: given two elements  $p \in W_2$ ,  $q \in W_2$ , that is,

$$\int_0^1 x p(x) dx = 0, \quad \int_0^1 x q(x) dx = 0,$$

then an arbitrary linear combination ( $a\mathbf{p} + b\mathbf{q}$ ) also belongs to  $W_2$ , since

$$\int_0^1 x (a\mathbf{p}(x) + b\mathbf{q}(x)) dx = a \int_0^1 x \mathbf{p}(x) dx + b \int_0^1 x \mathbf{q}(x) dx = 0.$$

◁

3. Find a basis of  $\mathbb{R}^4$  containing a basis of the null space of matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & -4 & 3 \\ 2 & 1 & 1 & 3 \\ 1 & 1 & -1 & 2 \\ 3 & 2 & 0 & 5 \end{bmatrix}.$$

**SOLUTION:** We first find a basis of  $N(\mathbf{A})$  using the Gauss method,

$$\begin{bmatrix} 1 & 2 & -4 & 3 \\ 2 & 1 & 1 & 3 \\ 1 & 1 & -1 & 2 \\ 3 & 2 & 0 & 5 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & -4 & 3 \\ 0 & -3 & 9 & -3 \\ 0 & -1 & 3 & -1 \\ 0 & -4 & 12 & -4 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & -4 & 3 \\ 0 & 1 & -3 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 2 & 1 \\ 0 & 1 & -3 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Therefore, a vector  $\mathbf{x} = [x_i] \in \mathbb{R}^4$  belongs to  $N(\mathbf{A})$  iff

$$\left. \begin{array}{l} x_1 = -2x_3 - x_4, \\ x_2 = 3x_3 - x_4, \end{array} \right\} \Rightarrow \mathbf{x} = \begin{bmatrix} -2 \\ 3 \\ 1 \\ 0 \end{bmatrix} x_3 + \begin{bmatrix} -1 \\ -1 \\ 0 \\ 1 \end{bmatrix} x_4,$$

so a basis for  $N(\mathbf{A})$  is the set

$$\mathcal{N} = \left\{ \begin{bmatrix} -2 \\ 3 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} -1 \\ -1 \\ 0 \\ 1 \end{bmatrix} \right\}.$$

Now we can find a basis of  $\mathbb{R}^4$  containing  $\mathcal{N}$  as follows. We start with any basis of  $\mathbb{R}^4$ , say the standard basis, and we construct a matrix,  $\mathbf{B}$  as follows:

$$\mathbf{B} = \begin{bmatrix} -2 & -1 & 1 & 0 & 0 & 0 \\ 3 & -1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \end{bmatrix}.$$

We then transform  $\mathbf{B}$  into reduced echelon form using Gauss operations. The pivot positions indicate the vectors forming a linearly independent set:

$$\begin{bmatrix} -2 & -1 & 1 & 0 & 0 & 0 \\ 3 & -1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \\ -2 & -1 & 1 & 0 & 0 & 0 \\ 3 & -1 & 0 & 1 & 0 & 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 2 & 1 \\ 0 & 0 & 0 & 1 & -3 & 1 \end{bmatrix}.$$

Therefore, a basis of  $\mathbb{R}^4$  containing  $\mathcal{N}$  is the following,

$$\left\{ \begin{bmatrix} -2 \\ 3 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} -1 \\ -1 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \right\}.$$

◁

4. Given a matrix  $\mathbf{A} \in \mathbb{F}^{n,n}$ , introduce its characteristic polynomial  $p_A(\lambda) = \det(\mathbf{A} - \lambda \mathbf{I}_n)$ . This polynomial has the form  $p_A(\lambda) = a_0 + a_1\lambda + \cdots + a_n\lambda^n$  for appropriate scalars  $a_0, \dots, a_n$ . Now introduce a matrix-valued function  $P_A: \mathbb{F}^{n,n} \rightarrow \mathbb{F}^{n,n}$  as follows

$$P_A(\mathbf{X}) = a_0 \mathbf{I}_n + a_1 \mathbf{X} + \cdots + a_n \mathbf{X}^n.$$

Determine whether the following statement is true or false and justify your answer: If matrix  $\mathbf{A} \in \mathbb{F}^{n,n}$  is diagonalizable, then  $P_A(\mathbf{A}) = \mathbf{0}$ . (That is,  $P_A(\mathbf{A})$  is the zero matrix.)

**SOLUTION:** We start with a comment:  $P_A(\mathbf{X}) \neq \det(\mathbf{A} - \mathbf{X}\mathbf{I})$ . Notice that the left-hand side is a matrix while the right hand side is a scalar. So an argument saying that the statement is true because  $\det(\mathbf{A} - \mathbf{A}) = 0$  is wrong. Once again,  $P_A(\mathbf{A}) \in \mathbb{F}^{n,n}$  and  $\det(\mathbf{A} - \mathbf{A}) = 0 \in \mathbb{F}$ .

Having clarified that, let us start saying that matrix  $\mathbf{A}$  is diagonalizable, that is, there exists an invertible matrix  $\mathbf{P}$  and a diagonal matrix  $\mathbf{D} = \text{diag}[\lambda_1, \dots, \lambda_n]$  such that  $\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$ , where  $\lambda_1, \dots, \lambda_n$  are eigenvalues of matrix  $\mathbf{A}$ . Therefore,

$$\begin{aligned} P_A(\mathbf{A}) &= P_A(\mathbf{P}\mathbf{D}\mathbf{P}^{-1}) \\ &= a_0 \mathbf{P}\mathbf{P}^{-1} + a_1 \mathbf{P}\mathbf{D}\mathbf{P}^{-1} + \cdots + a_n \mathbf{P}\mathbf{D}^n\mathbf{P}^{-1} \\ &= \mathbf{P}(a_0 \mathbf{I}_n + a_1 \mathbf{D} + \cdots + a_n \mathbf{D}^n)\mathbf{P}^{-1} \\ &= \mathbf{P} \text{diag}[p_A(\lambda_1), \dots, p_A(\lambda_n)] \mathbf{P}^{-1}. \end{aligned}$$

Since  $\lambda_1, \dots, \lambda_n$  are eigenvalues of matrix  $\mathbf{A}$ , we get that

$$p_A(\lambda_1) = 0, \dots, p_A(\lambda_n) = 0.$$

We then conclude that  $P_A(\mathbf{A}) = \mathbf{0}$ . So the statement is **true**. ◀

5. Consider the vector space  $\mathbb{R}^2$  with ordered bases

$$\mathcal{S} = \left( \mathbf{e}_{1s} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}_s, \mathbf{e}_{2s} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}_s \right), \quad \mathcal{U} = \left( \mathbf{u}_{1s} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}_s, \mathbf{u}_{2s} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}_s \right).$$

Let  $\mathbf{T}: \mathbb{R}^2 \rightarrow \mathbb{R}^2$  be a linear transformation given by

$$[\mathbf{T}(\mathbf{u}_1)]_s = \begin{bmatrix} -3 \\ 1 \end{bmatrix}_s, \quad [\mathbf{T}(\mathbf{u}_2)]_s = \begin{bmatrix} 1 \\ 3 \end{bmatrix}_s.$$

- (a) Find the matrix  $\mathbf{T}_{ss}$ .  
 (b) Find the matrix  $\mathbf{T}_{uu}$ .

**SOLUTION:**

- (a) From the data of the problem we get matrix

$$\mathbf{T}_{us} = \begin{bmatrix} -3 & 1 \\ 1 & 3 \end{bmatrix}.$$

We compute matrix  $\mathbf{T}_{ss}$  with the change of basis formula

$$\mathbf{T}_{ss} = \mathbf{Q}^{-1} \mathbf{T}_{us} \mathbf{P}, \quad \mathbf{P} = \mathbf{I}_{su}, \quad \mathbf{Q} = \mathbf{I}_{ss}.$$

From the data of the problem we get

$$\mathbf{I}_{us} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \Rightarrow \mathbf{P} = \frac{1}{3} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}.$$

Therefore,

$$\mathbf{T}_{ss} = \begin{bmatrix} -3 & 1 \\ 1 & 3 \end{bmatrix} \frac{1}{3} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \Rightarrow \boxed{\mathbf{T}_{ss} = \frac{1}{3} \begin{bmatrix} -7 & 5 \\ -1 & 5 \end{bmatrix}}.$$

- (b) We compute matrix  $\mathbf{T}_{uu}$  with the change of basis formula

$$\mathbf{T}_{uu} = \mathbf{Q}^{-1} \mathbf{T}_{us} \mathbf{P}, \quad \mathbf{P} = \mathbf{I}_{uu}, \quad \mathbf{Q} = \mathbf{I}_{us} = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}.$$

Therefore,

$$\mathsf{T}_{uu} = \frac{1}{3} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} -3 & 1 \\ 1 & 3 \end{bmatrix} \Rightarrow \boxed{\mathsf{T}_{uu} = \frac{1}{3} \begin{bmatrix} -7 & -1 \\ 5 & 5 \end{bmatrix}}.$$

◁

6. Consider the matrix  $\mathbf{A} = \begin{bmatrix} 1 & 2 \\ 2 & -1 \\ 1 & 1 \end{bmatrix}$  and the vector  $\mathbf{b} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$ .

- (a) Find the least-squares solution  $\hat{\mathbf{x}}$  to the matrix equation  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .  
 (b) Find the vector on  $R(\mathbf{A})$  that is the closest to the vector  $\mathbf{b}$  in  $\mathbb{R}^3$ .

**SOLUTION:**

(a) We first compute the normal equation  $\mathbf{A}^T\mathbf{A}\hat{\mathbf{x}} = \mathbf{A}^T\mathbf{b}$ . The coefficient matrix is

$$\mathbf{A}^T\mathbf{A} = \begin{bmatrix} 1 & 2 & 1 \\ 2 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 2 & -1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 6 & 1 \\ 1 & 6 \end{bmatrix};$$

while the source vector is

$$\mathbf{A}^T\mathbf{b} = \begin{bmatrix} 1 & 2 & 1 \\ 2 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 4 \\ 2 \end{bmatrix}.$$

Therefore, the least-squares solution  $\hat{\mathbf{x}}$  is obtained as follows,

$$\begin{bmatrix} 6 & 1 \\ 1 & 6 \end{bmatrix} \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \end{bmatrix} = \begin{bmatrix} 4 \\ 2 \end{bmatrix} \Rightarrow \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \end{bmatrix} = \frac{1}{35} \begin{bmatrix} 6 & -1 \\ -1 & 6 \end{bmatrix} \begin{bmatrix} 4 \\ 2 \end{bmatrix} \Rightarrow \boxed{\hat{\mathbf{x}} = \frac{2}{35} \begin{bmatrix} 11 \\ 4 \end{bmatrix}}.$$

(b) The vector in  $R(\mathbf{A})$  that is the closest to  $\mathbf{b}$  is the vector  $\mathbf{b}_n = \mathbf{A}\hat{\mathbf{x}}$ , that is,

$$\mathbf{b}_n = \mathbf{A}\hat{\mathbf{x}} = \begin{bmatrix} 1 & 2 \\ 2 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 11 \\ 4 \end{bmatrix} = \frac{2}{35} \begin{bmatrix} 11 \\ 4 \end{bmatrix} \Rightarrow \boxed{\mathbf{b}_n = \frac{2}{35} \begin{bmatrix} 19 \\ 18 \\ 15 \end{bmatrix}}.$$

We now check if the above result is true. If  $\mathbf{b}_n$  is correct, then  $(\mathbf{b} - \mathbf{b}_n) \in R(\mathbf{A})^\perp$ . Let us verify this last condition: We first compute the vector  $(\mathbf{b} - \mathbf{b}_n)$ , which is given by

$$(\mathbf{b} - \mathbf{b}_n) = \frac{1}{35} \left( \begin{bmatrix} 35 \\ 35 \\ 35 \end{bmatrix} - \begin{bmatrix} 38 \\ 36 \\ 30 \end{bmatrix} \right) \Rightarrow (\mathbf{b} - \mathbf{b}_n) = \frac{1}{35} \begin{bmatrix} -3 \\ -1 \\ 5 \end{bmatrix}.$$

It is simple to see that

$$\begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} -3 \\ -1 \\ 5 \end{bmatrix} = -3 - 2 + 5 = 0, \quad \begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} -3 \\ -1 \\ 5 \end{bmatrix} = -6 + 1 + 5 = 0.$$

◁

7. Consider the inner product space  $(\mathbb{C}^3, \cdot)$  and the subspace  $W \subset \mathbb{C}^3$  given by

$$W = \text{Span} \left( \left\{ \begin{bmatrix} i \\ 1 \\ i \end{bmatrix} \right\} \right).$$

- (a) Find a basis for  $W^\perp$ , the orthogonal complement of  $W$ .  
 (b) Find an orthonormal basis for  $W^\perp$ .

**SOLUTION:**(a) The vector  $\mathbf{w} \in W^\perp$  iff holds

$$0 = \begin{bmatrix} i \\ 1 \\ i \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = [-i \quad 1 \quad -i] \cdot \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = -iw_1 + w_2 - iw_3 \Rightarrow w_1 = -iw_2 - w_3.$$

Therefore,

$$\mathbf{w} = \begin{bmatrix} -i \\ 1 \\ 0 \end{bmatrix} w_2 + \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} w_3.$$

A basis for  $W^\perp$  is the set  $\mathcal{U}$  given by

$$\mathcal{U} = \left\{ \mathbf{u}_1 = \begin{bmatrix} -i \\ 1 \\ 0 \end{bmatrix}, \mathbf{u}_2 = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} \right\}.$$

(b) The basis above is not orthonormal, since

$$\begin{bmatrix} -i \\ 1 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} = [i \quad 1 \quad 0] \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} = -i \neq 0.$$

We use the Gram-Schmidt method to find an orthonormal basis,

$$\mathbf{u}_{2\perp} = \mathbf{u}_2 - \frac{\mathbf{u}_1 \cdot \mathbf{u}_2}{\|\mathbf{u}_1\|^2} \mathbf{u}_1.$$

We already computed  $\mathbf{u}_1 \cdot \mathbf{u}_2 = -i$ , while we also need

$$\|\mathbf{u}_1\|^2 = \begin{bmatrix} -i \\ 1 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} -i \\ 1 \\ 0 \end{bmatrix} = [i \quad 1 \quad 0] \begin{bmatrix} -i \\ 1 \\ 0 \end{bmatrix} = 2.$$

Then,

$$\mathbf{u}_{2\perp} = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} - \frac{(-i)}{2} \begin{bmatrix} -i \\ 1 \\ 0 \end{bmatrix} = \frac{1}{2} \left( \begin{bmatrix} -2 \\ 0 \\ 2 \end{bmatrix} + \begin{bmatrix} 1 \\ i \\ 0 \end{bmatrix} \right) \Rightarrow \mathbf{u}_{2\perp} = \frac{1}{2} \begin{bmatrix} -1 \\ i \\ 2 \end{bmatrix}.$$

It is simple to see that an orthonormal basis of  $W^\perp$  is given by

$$\mathcal{V} = \left\{ \mathbf{v}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} -i \\ 1 \\ 0 \end{bmatrix}, \mathbf{v}_2 = \frac{1}{\sqrt{6}} \begin{bmatrix} -1 \\ i \\ 2 \end{bmatrix} \right\}.$$

◀

8. (a) Find the eigenvalues and eigenvectors of matrix  $\mathbf{A} = \begin{bmatrix} 5 & 4 \\ 4 & 5 \end{bmatrix}$ , and show that  $\mathbf{A}$  is diagonalizable.  
 (b) Knowing that matrix  $\mathbf{A}$  above is diagonalizable, explicitly find a square root of matrix  $\mathbf{A}$ , that is, find a matrix  $\mathbf{X}$  such that  $\mathbf{X}^2 = \mathbf{A}$ . How many square roots does matrix  $\mathbf{A}$  have?

**SOLUTION:**

(a) The eigenvalues are the roots of the characteristic polynomial,

$$p(\lambda) = \begin{vmatrix} 5-\lambda & 4 \\ 4 & 5-\lambda \end{vmatrix} = (\lambda-5)^2 - 16 = \lambda^2 - 10\lambda + 9.$$

Therefore, the eigenvalues are computed by the well-known formula

$$\lambda_{\pm} = \frac{1}{2}(10 \pm \sqrt{100 - 36}) = \frac{1}{2}(10 \pm \sqrt{64}) = 5 \pm 4 \Rightarrow \begin{bmatrix} \lambda_+ = 9, \\ \lambda_- = 1. \end{bmatrix}$$

The corresponding eigenvectors are the following: For  $\lambda_+ = 9$  we obtain,

$$A - 9I_2 = \begin{bmatrix} -4 & 4 \\ 4 & -4 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -1 \\ 0 & 0 \end{bmatrix} \Rightarrow v_1^+ = v_2^+ \Rightarrow \boxed{v^+ = \begin{bmatrix} 1 \\ 1 \end{bmatrix}}.$$

For  $\lambda_- = 1$  we obtain

$$A - I_2 = \begin{bmatrix} 4 & 4 \\ 4 & 4 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} \Rightarrow v_1^- = -v_2^- \Rightarrow \boxed{v^- = \begin{bmatrix} -1 \\ 1 \end{bmatrix}}.$$

Introducing the matrices

$$P = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}, \quad D = \begin{bmatrix} 9 & 0 \\ 0 & 1 \end{bmatrix}, \quad P^{-1} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix},$$

it is simple to see that  $A$  is diagonalizable, that is,

$$\boxed{A = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 9 & 0 \\ 0 & 1 \end{bmatrix} \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}}.$$

(b) Since matrix  $X$  satisfies that  $X^2 = A = PDP^{-1}$ , matrix  $X^2$  is diagonalizable, sharing the eigenvectors and eigenvalues with matrix  $A$ , that is,

$$X^2 = P \begin{bmatrix} 9 & 0 \\ 0 & 1 \end{bmatrix} P^{-1},$$

We also know that a matrix  $X$  and its squared  $X^2$  share their eigenvectors with the eigenvalues of the latter being the squares of the eigenvalues of the former, that is,

$$X = P \begin{bmatrix} X_{11} & 0 \\ 0 & X_{22} \end{bmatrix} P^{-1}, \quad \text{with } x_{11} = \pm 3, \quad X_{22} = \pm 1.$$

The equation above says that there are four square-roots of matrix  $A$ . These matrices are the following,

$$\begin{aligned} \boxed{X_1 = P \begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix} P^{-1}}, & \quad \boxed{X_2 = P \begin{bmatrix} -3 & 0 \\ 0 & 1 \end{bmatrix} P^{-1}}, \\ \boxed{X_3 = P \begin{bmatrix} 3 & 0 \\ 0 & -1 \end{bmatrix} P^{-1}}, & \quad \boxed{X_4 = P \begin{bmatrix} -3 & 0 \\ 0 & -1 \end{bmatrix} P^{-1}}. \end{aligned}$$

◀

9. Consider the matrix  $A = \begin{bmatrix} 8 & -18 \\ 3 & -7 \end{bmatrix}$ .

(a) Find the eigenvalues and eigenvectors of  $A$ .

(b) Compute explicitly the matrix-valued function  $e^{At}$  for  $t \in \mathbb{R}$ .

**SOLUTION:**

(a) The eigenvalues are the roots of the characteristic polynomial,

$$p(\lambda) = \begin{vmatrix} 8 - \lambda & -18 \\ 3 & -7 - \lambda \end{vmatrix} = (\lambda + 7)(\lambda - 8) + 54 = \lambda^2 - \lambda - 2.$$

Therefore, the eigenvalues are computed by the well-known formula

$$\lambda_{\pm} = \frac{1}{2}(1 \pm \sqrt{1+8}) = \frac{1}{2}(1 \pm 3) \Rightarrow \boxed{\begin{matrix} \lambda_+ = 2, \\ \lambda_- = -1. \end{matrix}}$$

The corresponding eigenvectors are the following: For  $\lambda_+ = 2$  we obtain,

$$A - 2I_2 = \begin{bmatrix} 6 & -18 \\ 3 & -9 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -3 \\ 0 & 0 \end{bmatrix} \Rightarrow v_1^+ = 3v_2^+ \Rightarrow \boxed{v^+ = \begin{bmatrix} 3 \\ 1 \end{bmatrix}}.$$



For  $\lambda_- = -1$  we obtain

$$\mathbf{A} + \mathbf{I}_2 = \begin{bmatrix} 9 & -18 \\ 3 & -6 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -2 \\ 0 & 0 \end{bmatrix} \Rightarrow v_1^- = 2v_2^- \Rightarrow \boxed{\mathbf{v}^- = \begin{bmatrix} 2 \\ 1 \end{bmatrix}}.$$

(b) Since matrix  $\mathbf{A}$  is diagonalizable, we know that

$$e^{\mathbf{A}t} = \mathbf{P} e^{\mathbf{D}t} \mathbf{P}^{-1}, \quad \text{where } \mathbf{P} = \begin{bmatrix} 3 & 2 \\ 1 & 1 \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} 2 & 0 \\ 0 & -1 \end{bmatrix},$$

that is,

$$\boxed{e^{\mathbf{A}t} = \begin{bmatrix} 3 & 2 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} e^{2t} & 0 \\ 0 & e^{-t} \end{bmatrix} \begin{bmatrix} 1 & -2 \\ -1 & 3 \end{bmatrix}}.$$

◁

**10.** Find the function  $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^2$  solution of the initial value problem

$$\frac{d}{dt} \mathbf{x}(t) = \mathbf{A} \mathbf{x}(t), \quad \mathbf{x}(0) = \mathbf{x}_0,$$

where the matrix  $\mathbf{A} = \begin{bmatrix} -7 & 12 \\ -4 & 7 \end{bmatrix}$  and the vector  $\mathbf{x}_0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ .

**SOLUTION:** We first need to compute the eigenvalues and eigenvectors of matrix  $\mathbf{A}$ . The eigenvalues are the roots of the characteristic polynomial,

$$p(\lambda) = \begin{vmatrix} (-7 - \lambda) & 12 \\ -4 & (7 - \lambda) \end{vmatrix} = (\lambda + 7)(\lambda - 7) + 48 = \lambda^2 - 1.$$

Therefore, the eigenvalues are

$$\boxed{\lambda_+ = 1}, \quad \boxed{\lambda_- = -1}.$$

The corresponding eigenvectors are the following: For  $\lambda_+ = 1$  we obtain,

$$\mathbf{A} - \mathbf{I}_2 = \begin{bmatrix} -8 & 12 \\ -4 & 6 \end{bmatrix} \rightarrow \begin{bmatrix} 2 & -3 \\ 0 & 0 \end{bmatrix} \Rightarrow 2v_1^+ = 3v_2^+ \Rightarrow \boxed{\mathbf{v}^+ = \begin{bmatrix} 3 \\ 2 \end{bmatrix}}.$$

For  $\lambda_- = -1$  we obtain

$$\mathbf{A} + \mathbf{I}_2 = \begin{bmatrix} -6 & 12 \\ -4 & 8 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & -2 \\ 0 & 0 \end{bmatrix} \Rightarrow v_1^- = 2v_2^- \Rightarrow \boxed{\mathbf{v}^- = \begin{bmatrix} 2 \\ 1 \end{bmatrix}}.$$

Since the general solution to the differential equation above is

$$\mathbf{x}(t) = c_+ \mathbf{v}^+ e^{\lambda_+ t} + c_- \mathbf{v}^- e^{\lambda_- t},$$

we conclude that

$$\mathbf{x}(t) = c_+ \begin{bmatrix} 3 \\ 2 \end{bmatrix} e^t + c_- \begin{bmatrix} 2 \\ 1 \end{bmatrix} e^{-t}.$$

The initial condition implies that

$$\begin{bmatrix} 1 \\ 1 \end{bmatrix} = \mathbf{x}_0 = \mathbf{x}(0) = c_+ \begin{bmatrix} 3 \\ 2 \end{bmatrix} + c_- \begin{bmatrix} 2 \\ 1 \end{bmatrix}.$$

We need to solve the linear system

$$\begin{bmatrix} 3 & 2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} c_+ \\ c_- \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow \begin{bmatrix} c_+ \\ c_- \end{bmatrix} = \frac{1}{(-1)} \begin{bmatrix} 1 & -2 \\ -2 & 3 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow \begin{bmatrix} c_+ \\ c_- \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}.$$

So the solution to the initial value problem above is

$$\boxed{\mathbf{x}(t) = \begin{bmatrix} 3 \\ 2 \end{bmatrix} e^t - \begin{bmatrix} 2 \\ 1 \end{bmatrix} e^{-t}}.$$

◁

## REFERENCES

- [1] S. Hassani. *Mathematical physics*. Springer, New York, 2000. Corrected second printing.
- [2] D. Lay. *Linear algebra and its applications*. Addison Wesley, New York, 2005. Third updated edition.
- [3] C. Meyer. *Matrix analysis and applied linear algebra*. SIAM, Philadelphia, 2000.
- [4] G. Strang. *linear algebra and its applications*. Brooks/Cole, India, 2005. Fourth edition.