

# Linear algebra for Machine Learning

Paul MINCHELLA, Stéphane CHRÉTIEN

October 14, 2025

# Contents

<b>1</b>	<b>Motivations and Review</b>	<b>3</b>
1.1	An Introductory Example . . . . .	3
1.2	Gaussian Elimination . . . . .	3
1.3	Matrices: Definition and Properties . . . . .	4
<b>2</b>	<b>Vector Spaces</b>	<b>10</b>
2.1	From Number Sets to Algebraic Structures . . . . .	10
2.2	Fundamental Algebraic Structures . . . . .	10
2.3	Vector Spaces . . . . .	12
2.3.1	Examples of Vector Spaces . . . . .	13
2.3.2	Geometric Intuition . . . . .	14
2.3.3	Subspaces of a Vector Space . . . . .	14
2.4	Span, Linear Independence, and Basis . . . . .	15
2.4.1	Linear Independence and Generating Families . . . . .	15
2.4.2	Basis and Dimension . . . . .	16
2.4.3	Dimension of Subspaces . . . . .	16
2.4.4	Extension and Extraction of Bases . . . . .	17
<b>3</b>	<b>Linear Transformations</b>	<b>18</b>
3.1	Definition and Motivation . . . . .	18
3.2	Examples in $\mathbb{R}^3$ . . . . .	18
3.3	Isomorphisms and Endomorphisms . . . . .	19
3.4	Kernel and Image . . . . .	19
3.5	Rank–Nullity Theorem . . . . .	19
3.6	Matrices Associated with Linear Maps . . . . .	21
3.6.1	From Linear Maps to Matrices . . . . .	21
<b>4</b>	<b>Change of Basis</b>	<b>23</b>
4.1	Motivation . . . . .	23
4.2	General Setup . . . . .	23
4.3	Worked Example in Dimension 3 . . . . .	24
4.4	Geometric Interpretation . . . . .	25
4.5	Statement and semonstration . . . . .	26
<b>5</b>	<b>Diagonalization</b>	<b>28</b>
5.1	Motivation: Why Diagonalize? . . . . .	28
5.2	Eigenvalues and Eigenvectors . . . . .	28
5.3	Characteristic Polynomial and Eigenspace . . . . .	28
5.4	Change of Basis and Similarity . . . . .	30
5.5	Diagonalization Formalism . . . . .	31
5.6	Important Operators on Matrices . . . . .	32
5.7	Application: Equilibrium States . . . . .	33
5.8	Orthogonal Matrices and the Spectral Theorem . . . . .	35
5.8.1	Orthogonal Matrices . . . . .	35
5.8.2	The Spectral Theorem . . . . .	36

5.9	Statement and demonstration . . . . .	37
5.9.1	Diagonalization criterion . . . . .	37
5.9.2	Similarity as an Equivalence Relation . . . . .	38
5.9.3	The Quotient Set . . . . .	38
5.9.4	Jordan Matrices . . . . .	38
5.10	Gram–Schmidt Orthogonalization . . . . .	40
<b>6</b>	<b>Application of Diagonalization</b> . . . . .	<b>41</b>
6.1	Ranking via Diagonalization: The Google Matrix . . . . .	41
6.1.1	Motivation . . . . .	41
6.1.2	From Graph to Matrices . . . . .	41
6.1.3	Ranking Vector . . . . .	42
6.1.4	Perron–Frobenius Theorem . . . . .	43
6.2	Spectral properties of stochastic matrices and convergence rates . . . . .	43
6.2.1	Spectrum localization . . . . .	43
6.2.2	Why the second eigenvalue controls convergence . . . . .	44
6.2.3	Quantifying the “speed”: inverse gap and log factor . . . . .	44
6.2.4	Aperiodicity, periodicity, and complex eigenvalues . . . . .	45
6.2.5	Practical takeaway . . . . .	45
6.3	Data Analysis via Diagonalization: Covariance and PCA . . . . .	45
6.3.1	The Covariance Matrix . . . . .	45
6.3.2	Variance and Eigenanalysis . . . . .	45
6.3.3	Principal Component Analysis (PCA) . . . . .	45
6.3.4	Worked Example . . . . .	46
<b>7</b>	<b>Application to Statistics: Least Squares and SVD</b> . . . . .	<b>47</b>
7.1	Orthogonality and Distance Minimization . . . . .	47
7.1.1	Law of Cosines and the Dot Product . . . . .	47
7.2	Least Squares Approximation . . . . .	48
7.2.1	Motivation . . . . .	48
7.2.2	Warm-up Example: Line of Best Fit in $\mathbb{R}^2$ . . . . .	48
7.3	Ordinary Least Squares (OLS) . . . . .	48
7.4	Connection to Singular Value Decomposition (SVD) . . . . .	49
7.5	Subspaces and Orthogonality . . . . .	50
7.5.1	Definitions . . . . .	50
7.5.2	Fundamental Subspaces of a Matrix . . . . .	50
7.5.3	Worked Example . . . . .	50
7.6	Least Squares and Projection Matrices . . . . .	51
7.6.1	Quadratic Least Squares . . . . .	51
7.7	Singular Value Decomposition (SVD) . . . . .	51
7.7.1	From $M^\top M$ to the SVD: symmetry, (semi)definiteness, and eigenvalues . . . . .	51
7.7.2	SVD Statement . . . . .	52
<b>8</b>	<b>Matrix Conditioning and Numerical Stability</b> . . . . .	<b>58</b>
8.1	Matrix and Vector Norms . . . . .	58
8.1.1	Euclidean norm and induced matrix norm . . . . .	58
8.1.2	Spectral radius and its relation to norms . . . . .	59
8.2	Definition of the Matrix Condition Number . . . . .	59
8.2.1	Relative error amplification . . . . .	60
8.2.2	Spectral expression (in 2-norm) . . . . .	60
8.3	Ill-conditioning and the Explosion of the Inverse . . . . .	62
8.3.1	An intuitive illustration . . . . .	62
8.3.2	Computing the inverse . . . . .	62
8.3.3	Rigorous formalism . . . . .	62
8.4	Main Results and Theorems . . . . .	62
8.5	Conclusion . . . . .	63

# Chapter 1

## Motivations and Review

Linear algebra plays a central role in many areas of mathematics and its applications. It can truly be regarded as a computational engine. It appears in numerical analysis, for instance in the finite element method, as well as in algebraic geometry, through concepts such as the Hodge decomposition. It is also at the heart of statistics, where covariance matrices and the analysis of data shape are fundamental.

In practice, data science practitioners often come from very diverse backgrounds. It is therefore important to rely on a shared foundation of linear algebraic tools and methods. Moreover, certain essential concepts — such as eigenvalues and eigenvectors — are frequently misunderstood or forgotten, and require regular review.

The goal of this course is to develop real dexterity with the core instruments of linear algebra. We will revisit the resolution of linear systems, Gaussian elimination, and the manipulation of matrices. We will study vector spaces, linear transformations, and basis changes. We will analyze diagonalization and its applications, ranging from eigenvalue theory to practical examples such as webpage ranking or the study of covariance matrices. Finally, we will emphasize orthogonality, least-squares methods, and the singular value decomposition (SVD).

### 1.1 An Introductory Example

Consider the following example, which illustrates how a real-world phenomenon can be modeled with linear algebra.

**Example 1.1.1.** *In the town of Smallville, each year 30% of nonsmokers begin smoking, while 20% of smokers quit. Initially, the population consists of 8000 smokers and 2000 nonsmokers. We can ask several natural questions: what will the numbers be after 100 years? More generally, what is the population after  $n$  years? Does the system tend to a stable equilibrium?*

Such dynamics can be modeled by transition matrices, and their analysis relies on the theory of eigenvalues and eigenvectors.

### 1.2 Gaussian Elimination

One of the most fundamental tools in linear algebra for solving systems is Gaussian elimination. The general idea is to systematically reduce a linear system to a simpler one, with fewer equations and variables, through a process of elimination.

Consider for example a system of two equations:

$$\begin{aligned}a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= b_1, \\a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n &= b_2.\end{aligned}$$

The first step is to choose a pivot, that is, to make sure that  $a_{11}$  is nonzero, possibly by swapping two rows. Then we normalize the first row by dividing through by  $a_{11}$ . Finally, we eliminate the entries below the pivot by subtracting from each lower row a suitable multiple of the first row.

This procedure is then repeated on the remaining submatrix until we obtain an upper-triangular system. The solution is then completed by back-substitution. If the algorithm is carried to the end, we can even obtain the reduced row echelon form (RREF).

It is important to note that the operations used in this process preserve the solution set of the system  $Ax = b$ . Indeed, each operation corresponds to a left-multiplication by an elementary matrix  $E$ , which is invertible. Therefore we have the equivalence

$$Ax = b \iff (EA)x = Eb.$$

The basic operations are:

- row swap:  $L_i \leftrightarrow L_j$ ,
- scaling:  $L_i \leftarrow \lambda L_i$  with  $\lambda \neq 0$ ,
- row addition:  $L_i \leftarrow L_i + \lambda L_j$ .

*Exercise 1.2.1.* Solve the following system using Gaussian elimination:

$$\begin{cases} x + y + z = 3, \\ 2x + y = 7, \\ 3x + 2z = 5. \end{cases}$$

This exercise will allow us to practice the steps of pivoting, normalization, elimination, and back-substitution.

## 1.3 Matrices: Definition and Properties

Let's introduce one of the most important object in linear algebra and in mathematics in general: the matrix.

**Definition 1.3.1** (Matrix). A matrix is an  $m \times n$  array of elements, where  $m$  is the number of rows and  $n$  is the number of columns.

Formally,

$$A \in \mathcal{M}_{m \times n}(\mathbb{K}),$$

where  $\mathbb{K}$  is a field (for example  $\mathbb{R}$  or  $\mathbb{C}$ ). In the case  $\mathbb{K} = \mathbb{R}$ , one also writes  $A \in \mathbb{R}^{m \times n}$ .

**Property 1.3.1** (Basic properties of  $\mathcal{M}_{m \times n}(\mathbb{K})$ ). The set of  $m \times n$  matrices with coefficients in a field  $\mathbb{K}$  satisfies the following properties:

- It forms a vector space over  $\mathbb{K}$ , with entrywise addition and scalar multiplication.
- Its dimension is  $\dim \mathcal{M}_{m \times n}(\mathbb{K}) = mn$ .
- Matrix multiplication is defined if the inner dimensions match:

$$A \in \mathcal{M}_{m \times n}(\mathbb{K}), \quad B \in \mathcal{M}_{n \times p}(\mathbb{K}) \implies AB \in \mathcal{M}_{m \times p}(\mathbb{K}).$$

- Multiplication is associative but not commutative in general.

**Matrix Multiplication as Composition of Linear Systems** Matrix multiplication corresponds to the composition of two linear maps:

$$C = AB \iff \text{Apply } B \text{ first, then } A.$$

To illustrate that, let consider a two  $2 \times 2$  systems. Denoting

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

Applying  $B$  to  $\begin{pmatrix} x \\ y \end{pmatrix}$  gives

$$\begin{cases} z = 2x, \\ w = x + 3y. \end{cases}$$

Then applying  $A$  to  $\begin{pmatrix} z \\ w \end{pmatrix}$ :

$$\begin{cases} u = z + 2w = 4x + 6y, \\ v = w = x + 3y. \end{cases}$$

Thus, the combined transformation  $\begin{pmatrix} x \\ y \end{pmatrix} \mapsto \begin{pmatrix} u \\ v \end{pmatrix}$  is given by

$$C = AB = \begin{bmatrix} 4 & 6 \\ 1 & 3 \end{bmatrix}.$$

This illustrates that matrix multiplication is equivalent to chaining linear systems.

## Dot Product and Matrix Multiplication

**Definition 1.3.2** (Dot Product). For two vectors  $\mathbf{v} = [a_1, \dots, a_n]$  and  $\mathbf{w} = [b_1, \dots, b_n]$ , the dot product is defined as

$$\mathbf{v} \cdot \mathbf{w} = \sum_{i=1}^n a_i b_i.$$

The norm of a vector is given by

$$\|\mathbf{v}\| = \sqrt{\mathbf{v} \cdot \mathbf{v}}.$$

By the law of cosines,  $\mathbf{v}$  and  $\mathbf{w}$  are orthogonal if and only if  $\mathbf{v} \cdot \mathbf{w} = 0$ .

*Remark 1.* The dot product is also called the *scalar product*, since it associates to two vectors a real scalar. It is commonly denoted using angle brackets:

$$\langle \mathbf{v}, \mathbf{w} \rangle = \mathbf{v} \cdot \mathbf{w}, \quad \forall \mathbf{v}, \mathbf{w} \in \mathbb{R}^n.$$

The coordinate expression  $\sum_{i=1}^n a_i b_i$  is valid only when  $\mathbf{v}$  and  $\mathbf{w}$  are expressed in the canonical basis (or more generally, in an orthonormal basis). The inner product itself, however, is basis-independent.

**Property 1.3.2** (Matrix multiplication). Let  $A \in \mathcal{M}_{m \times n}(\mathbb{K})$  and  $B \in \mathcal{M}_{p \times q}(\mathbb{K})$ . Matrix multiplication  $AB$  is defined when  $n = p$ .

If  $(AB)_{ij}$  denotes the  $(i, j)$ -entry of the product, then

$$(AB)_{ij} = \text{row}_i(A) \cdot \text{col}_j(B).$$

*Interpretation:* Each matrix represents a linear map with respect to a chosen basis. Hence, multiplication of matrices (composition of linear maps) and the dot product (row  $\cdot$  column) only make sense in the same basis. This idea will be further formalized with the notions of linear maps and bases.

*Exercise 1.3.1.* Compute the missing entries in the following product:

$$\begin{bmatrix} 2 & 7 \\ 3 & 3 \\ 1 & 5 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \end{bmatrix} = \begin{bmatrix} 37 & 46 & 55 & 64 \\ * & * & * & * \\ * & * & * & * \end{bmatrix}.$$

**Transpose, Symmetry, and Inverses** We now introduce several fundamental notions related to matrices: the transpose, symmetry, diagonal and identity matrices, and the concept of an inverse, providing both algebraic structure and geometric intuition.

**Definition 1.3.3** (Transpose). Let  $A \in \mathcal{M}_{m \times n}(\mathbb{K})$ . The transpose of  $A$ , denoted  $A^\top$ , is the  $n \times m$  matrix defined by

$$(A^\top)_{ij} = A_{ji}.$$

The transpose simply flips a matrix along its main diagonal: rows become columns, and columns become rows. This operation is extremely natural: it exchanges the perspective of reading the coefficients horizontally or vertically.

**Property 1.3.3** (Properties of the transpose). For all  $A, B \in \mathcal{M}_{m \times n}(\mathbb{K})$  and  $\alpha, \beta \in \mathbb{K}$ :

- *Linearity:*

$$(\alpha A + \beta B)^\top = \alpha A^\top + \beta B^\top.$$

- *Involution:*

$$(A^\top)^\top = A.$$

- *Product rule (note the order reversal):*

$$(AB)^\top = B^\top A^\top.$$

These rules are straightforward but powerful. The reversal in the product rule is particularly important: it reflects the fact that transposition changes the order of composition of linear maps.

**Definition 1.3.4** (Symmetric matrix). A square matrix  $A \in \mathcal{M}_{n \times n}(\mathbb{K})$  is called symmetric if

$$A^\top = A.$$

Symmetric matrices are fundamental because they correspond, in geometry, to linear transformations that preserve certain inner products. They are also the matrices that can be diagonalized in an orthonormal basis when working over  $\mathbb{R}$ .

**Definition 1.3.5** (Diagonal matrix). A square matrix  $A \in \mathcal{M}_{n \times n}(\mathbb{K})$  is called diagonal if

$$A_{ij} \neq 0 \quad \Rightarrow \quad i = j.$$

For two diagonal matrices  $A, B \in \mathcal{M}_{n \times n}(\mathbb{K})$ , multiplication is commutative:

$$AB = BA,$$

and moreover

$$(AB)_{ii} = a_{ii} b_{ii}.$$

Diagonal matrices are the simplest possible ones: they act by rescaling each coordinate independently. Because of this simplicity, they commute with each other, which is very rare in general for matrices.

**Definition 1.3.6** (Identity matrix). The identity matrix of size  $n$  is

$$I_n = \text{diag}(1, \dots, 1).$$

It satisfies, for any  $A \in \mathcal{M}_{n \times n}(\mathbb{K})$

$$I_n A = A I_n = A.$$

**Definition 1.3.7** (Inverse matrix). *A square matrix  $A \in \mathcal{M}_{n \times n}(\mathbb{K})$  is invertible if there exists a matrix  $B \in \mathcal{M}_{n \times n}(\mathbb{K})$  such that*

$$BA = AB = I_n.$$

*In this case,  $B$  is unique and is called the inverse of  $A$ , denoted  $A^{-1}$ .*

The existence of an inverse means that the linear system associated with  $A$  can always be solved uniquely. In practical terms, applying  $A^{-1}$  precisely cancels the effect of  $A$ .

**Non-commutativity, Linear Systems, and Diagonalization** We now turn to some key structural properties of matrices. We start with an exercise that illustrates an important fact: matrix multiplication is not commutative in general. This contrasts with the familiar multiplication of real numbers, and it has deep consequences for linear algebra.

*Exercise 1.3.2.* We propose here two exercise.

1. Find  $2 \times 2$  matrices  $(A, B)$  such that  $AB \neq BA$ .
2. Show that for any matrices  $A, B$  of compatible dimensions:

$$(AB)^\top = B^\top A^\top.$$

Deduce that  $A^\top A$  is always symmetric.

The first item encourages you to experiment with small matrices to see concretely that commutativity fails in general. The second item recalls a fundamental property of the transpose: it reverses the order of multiplication. From this rule, we immediately see that

$$(A^\top A)^\top = A^\top (A^\top)^\top = A^\top A,$$

hence  $A^\top A$  is symmetric for every matrix  $A$ .

**Definition 1.3.8** (Matrix form of a linear system). *A system of  $n$  linear equations in  $m$  unknowns can be written compactly in matrix form:*

$$A\mathbf{x} = \mathbf{b}, \quad A \in \mathcal{M}_{n \times m}(\mathbb{K}), \quad \mathbf{x} = [x_1, \dots, x_m]^\top.$$

This representation highlights the power of matrices: rather than working with each equation individually, the entire system is captured by a single matrix-vector equation. Solving the system amounts to understanding the properties of the matrix  $A$ .

**Diagonalization: A Motivating Example** One of the most powerful tools in linear algebra is *diagonalization*. To motivate it, we consider a simple real-world model.



**Example 1.3.1** (Smoking in Smallville). Suppose we study the evolution of smokers and non-smokers in a small town. Let

$$(n_t, s_t) = (\text{number of nonsmokers, number of smokers}) \text{ at year } t.$$

The yearly transition is modeled by the matrix

$$\begin{bmatrix} n_{t+1} \\ s_{t+1} \end{bmatrix} = \begin{bmatrix} 0.7 & 0.2 \\ 0.3 & 0.8 \end{bmatrix} \begin{bmatrix} n_t \\ s_t \end{bmatrix}.$$

By iterating the process, after  $t$  years we obtain

$$\begin{bmatrix} n_t \\ s_t \end{bmatrix} = \left( \begin{bmatrix} 0.7 & 0.2 \\ 0.3 & 0.8 \end{bmatrix} \right)^t \begin{bmatrix} n_0 \\ s_0 \end{bmatrix}.$$

Thus the long-term behavior of the population is governed by powers of the matrix

$$A = \begin{bmatrix} 0.7 & 0.2 \\ 0.3 & 0.8 \end{bmatrix}.$$

The central problem is now clear: to study the system for large  $t$ , we need to compute  $A^t$ . But repeated multiplication is computationally expensive, and it does not reveal structural properties. Is there a better way?

**Motivation for diagonalization** If  $A$  is diagonalizable, we will show in this course that there exists an invertible matrix  $P$  and a diagonal matrix  $D$  such that

$$PAP^{-1} = D.$$

Then, for any integer  $m \geq 1$ ,

$$A^m = P^{-1}D^mP.$$

Since  $D$  is diagonal, computing  $D^m$  is straightforward: one simply raises each diagonal entry to the  $m$ -th power.

**Intuition.** Diagonalization provides the key to understanding the long-term dynamics of linear systems. Instead of iterating  $A$  directly, we transfer the problem into a basis where  $A$  acts by simple scaling along independent directions. In that basis, the action of  $A$  is completely transparent.

**Proposition 1.3.1.** Let

$$D = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n) \in M_n(\mathbb{R}).$$

Then for every integer  $m \geq 0$ ,

$$D^m = \text{diag}(\lambda_1^m, \lambda_2^m, \dots, \lambda_n^m).$$

**Proof.** We proceed by induction on  $m$ .

*Base case:* For  $m = 0$ , we have  $D^0 = I_n$ , while

$$\text{diag}(\lambda_1^0, \dots, \lambda_n^0) = \text{diag}(1, \dots, 1) = I_n.$$

So the formula holds.

*Inductive step:* Suppose the formula holds for some  $m \geq 0$ , i.e.,

$$D^m = \text{diag}(\lambda_1^m, \dots, \lambda_n^m).$$

Then

$$D^{m+1} = D^m D = \text{diag}(\lambda_1^m, \dots, \lambda_n^m) \cdot \text{diag}(\lambda_1, \dots, \lambda_n).$$

Since the product of two diagonal matrices is diagonal with entries multiplied componentwise, we obtain

$$D^{m+1} = \text{diag}(\lambda_1^{m+1}, \dots, \lambda_n^{m+1}).$$

Thus, by induction, the formula holds for all  $m \in \mathbb{N}$ .

□

# Chapter 2

## Vector Spaces

### 2.1 From Number Sets to Algebraic Structures

We begin by recalling the basic sets of numbers that underpin all of algebra. Each extension of number systems responds to a need: negative numbers for subtraction, fractions for division, irrational numbers for completeness, and complex numbers for solving polynomial equations.

**Definition 2.1.1.** *Number sets*

- **Natural numbers:**

$$\mathbb{N} = \{0, 1, 2, 3, \dots\},$$

*the set of non-negative integers.*

- **Integers:**

$$\mathbb{Z} = \{\dots, -2, -1, 0, 1, 2, \dots\},$$

*all negative and positive whole numbers.*

- **Rational numbers:**

$$\mathbb{Q} = \left\{ \frac{p}{q} : p \in \mathbb{Z}, q \in \mathbb{Z}^*, q \neq 0 \right\},$$

*ratios of integers.*

- **Real numbers:**  $\mathbb{R}$  is the completion of  $\mathbb{Q}$ : a totally ordered complete field. Famous irrational examples include  $\pi, e, \sqrt{2}, \varphi = \frac{1+\sqrt{5}}{2}$ .

- **Complex numbers:**

$$\mathbb{C} = \{x + iy : x, y \in \mathbb{R}, i^2 = -1\},$$

*which extend  $\mathbb{R}$  to ensure that every polynomial equation has a solution.*

**Example 2.1.1** (Hierarchy of number sets). *These sets are naturally nested:*

$$\mathbb{N} \subset \mathbb{Z} \subset \mathbb{Q} \subset \mathbb{R} \subset \mathbb{C}.$$

This hierarchy illustrates how mathematics evolves by progressively enlarging its domain to resolve limitations of the previous set.

### 2.2 Fundamental Algebraic Structures

Beyond numbers themselves, algebra is concerned with the *operations* defined on sets and the structures they create.

**Definition 2.2.1** (Binary operation). Let  $E$  be a set. An internal binary operation on  $E$  is a map

$$\star : E \times E \rightarrow E, \quad (x, y) \mapsto x \star y.$$

For example,  $+$  and  $\times$  are binary operations on  $\mathbb{Z}$ .

Binary operations allow us to combine elements of a set. Some operations have particularly rich properties, which lead to well-known algebraic structures.

**Definition 2.2.2** (Group). A group is a pair  $(G, \star)$  where  $\star$  is a binary operation satisfying:

- **Associativity:**  $(x \star y) \star z = x \star (y \star z)$ .
- **Identity element:** there exists  $e \in G$  such that  $x \star e = e \star x = x$  for all  $x$ .
- **Inverse:** every  $x \in G$  has an inverse  $x^{-1}$  such that  $x \star x^{-1} = e$ .

If  $x \star y = y \star x$  for all  $x, y$ , the group is called abelian.

**Example 2.2.1** (Group structure).  $(\mathbb{Z}, +)$  is an abelian group: the identity is 0, and every integer has an additive inverse. By contrast,  $(\mathbb{Z}, \times)$  is not a group: not every integer admits a multiplicative inverse within  $\mathbb{Z}$  (for instance 2 has no integer reciprocal).

**Definition 2.2.3** (Ring). A ring  $(A, +, \times)$  is a set equipped with two operations such that:

- $(A, +)$  is an abelian group,
- multiplication  $\times$  is associative and has an identity element 1,
- multiplication distributes over addition.

**Definition 2.2.4** (Field). A field is a ring  $(K, +, \times)$  in which every nonzero element has a multiplicative inverse.

**Example 2.2.2** (Examples of fields and rings).  $\mathbb{Q}, \mathbb{R}, \mathbb{C}$  are fields, since every nonzero element admits a multiplicative inverse. On the other hand,  $\mathbb{Z}$  is a ring but not a field, as only  $\pm 1$  have inverses in  $\mathbb{Z}$ .

To test our understanding of these definitions, let us explore an unusual operation on the integers.

*Exercise 2.2.1* (A nonstandard binary operation on  $\mathbb{Z}$ ). Define, for  $a, b \in \mathbb{Z}$ , the operation

$$a \star b = a + b + 1.$$

1. Show that  $\star$  is an *internal* binary operation on  $\mathbb{Z}$ .
2. Check whether  $\star$  is associative and whether it is commutative.
3. Determine the identity element  $e$  for  $\star$ .
4. For a given  $a \in \mathbb{Z}$ , find the inverse of  $a$  with respect to  $\star$ .
5. Conclude: is  $(\mathbb{Z}, \star)$  a group? Is it abelian?

**Solution.** (1) *Internal law.* For  $a, b \in \mathbb{Z}$ , we have  $a \star b = a + b + 1 \in \mathbb{Z}$  since  $\mathbb{Z}$  is closed under  $+$  and contains 1. Hence  $\star$  is an internal binary operation on  $\mathbb{Z}$ .

(2) *Associativity.* For  $a, b, c \in \mathbb{Z}$ ,

$$(a \star b) \star c = (a + b + 1) \star c = (a + b + 1) + c + 1 = a + b + c + 2,$$

while

$$a \star (b \star c) = a \star (b + c + 1) = a + (b + c + 1) + 1 = a + b + c + 2.$$

Thus  $(a \star b) \star c = a \star (b \star c)$ , so  $\star$  is associative.

(3) *Commutativity.*

$$a \star b = a + b + 1 = b + a + 1 = b \star a,$$

so  $\star$  is commutative.

(4) *Identity element.* We seek  $e \in \mathbb{Z}$  such that for all  $a \in \mathbb{Z}$ ,

$$a \star e = a \quad \text{and} \quad e \star a = a.$$

From  $a \star e = a + e + 1 = a$  we get  $e = -1$ . Conversely,  $e \star a = (-1) + a + 1 = a$ . Hence the identity is  $e = -1$ .

(5) *Inverse of  $a$ .* Given  $a \in \mathbb{Z}$ ,  $b$  is an inverse if  $a \star b = e = -1$ :

$$a \star b = a + b + 1 = -1 \iff b = -a - 2.$$

By commutativity, also  $b \star a = -1$ . Therefore the  $\star$ -inverse of  $a$  is  $a^{-1\star} = -a - 2$ .

(6) *Conclusion.* The operation  $\star$  is associative and commutative, admits identity  $-1$ , and every  $a \in \mathbb{Z}$  has a  $\star$ -inverse  $-a - 2$ . Hence  $(\mathbb{Z}, \star)$  is an abelian group.

(Optional isomorphism insight). The map  $\varphi : \mathbb{Z} \rightarrow \mathbb{Z}$ ,  $\varphi(a) = a + 1$ , is a bijection and satisfies

$$\varphi(a \star b) = \varphi(a + b + 1) = (a + b + 1) + 1 = (a + 1) + (b + 1) = \varphi(a) + \varphi(b),$$

so  $\varphi$  is a group isomorphism  $(\mathbb{Z}, \star) \cong (\mathbb{Z}, +)$ . □

## 2.3 Vector Spaces

In order to understand and fully master the notion of matrices, it is necessary to introduce the abstract framework in which they naturally live: that of **vector spaces** and **linear maps**.

A vector space is a mathematical structure that generalizes and unifies very different kinds of objects, such as:

- vectors in the plane  $\mathbb{R}^2$  or in space  $\mathbb{R}^3$ , often represented as arrows with direction and magnitude;
- polynomials in one variable, which can be added together and multiplied by real numbers;
- continuous functions defined on an interval, which can also be added and multiplied by scalars.

All these examples share the same logic: we can *add* objects of the same nature, and we can *scale* them by numbers coming from a field (such as  $\mathbb{R}$  or  $\mathbb{C}$ ). This double structure — internal addition and external scalar multiplication — is precisely what defines a vector space.

**Definition 2.3.1** (Vector Space). Let  $\mathbb{K}$  be a field (for example  $\mathbb{R}$  or  $\mathbb{C}$ ). A vector space over  $\mathbb{K}$  is a set  $V$  equipped with two operations:

- an internal addition

$$(x, y) \mapsto x + y \quad \text{in } V,$$

- an external scalar multiplication

$$(\lambda, x) \mapsto \lambda x \quad \text{from } \mathbb{K} \times V \rightarrow V,$$

which satisfy the following axioms for all  $x, y, z \in V$  and  $\lambda, \mu \in \mathbb{K}$ :

1.  $(V, +)$  is an abelian group: addition is associative and commutative, there exists a neutral element 0, and each vector  $x$  has an additive inverse  $-x$ .

2. Compatibility of scalars:

$$(\lambda\mu)x = \lambda(\mu x).$$

3. Neutral element of scalars:

$$1_{\mathbb{K}}x = x.$$

4. Distributivity:

$$(\lambda + \mu)x = \lambda x + \mu x, \quad \lambda(x + y) = \lambda x + \lambda y.$$

Thus, a vector space provides an abstract framework in which one can manipulate objects “as if they were vectors,” even when they do not have an obvious geometric representation. This abstraction is powerful because it allows us to study very diverse mathematical situations with the same unified set of tools. In this setting, matrices will naturally appear as the *concrete representations* of linear maps.

### 2.3.1 Examples of Vector Spaces

Vector spaces come in many flavors, from finite-dimensional coordinate spaces to spaces of functions. They all share the same axioms: closure under addition and scalar multiplication, distributivity, etc. The following examples are canonical and will be used throughout the course.

**Example 2.3.1.** Some vector spaces:

- **Coordinate space.** For a field  $\mathbb{K}$  and  $n \in \mathbb{N}$ , the set  $\mathbb{K}^n$  of  $n$ -tuples with entries in  $\mathbb{K}$  is a vector space with entrywise operations:

$$\begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix} + \lambda \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} = \begin{pmatrix} a_1 + \lambda b_1 \\ \vdots \\ a_n + \lambda b_n \end{pmatrix} \in \mathbb{K}^n.$$

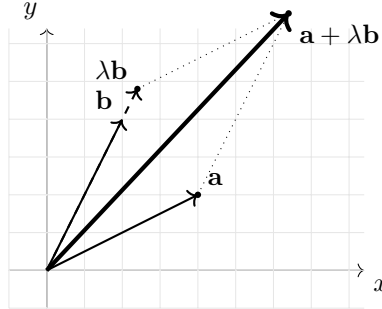
- **Polynomials.** The set  $\mathbb{R}[X]$  of polynomials with real coefficients is a vector space over  $\mathbb{R}$ ; addition and scalar multiplication are defined coefficientwise.
- **Continuous functions.** The set  $C^0([a, b], \mathbb{R})$  of all real-valued continuous functions on  $[a, b]$  is a vector space (pointwise operations).

These examples illustrate a key point: vectors need not be geometric arrows in the plane. A “vector” can be a list of numbers, a polynomial, or even a function—what matters are the axioms.

### 2.3.2 Geometric Intuition

In low dimensions (e.g.,  $\mathbb{R}^2$ ), linear combinations have a vivid geometric meaning: the vector  $\mathbf{a} + \lambda\mathbf{b}$  slides along the line through  $\mathbf{a}$  in the direction of  $\mathbf{b}$ , and the parallelogram rule encodes vector addition.

$$\begin{pmatrix} a_1 \\ a_2 \end{pmatrix} + \lambda \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} = \begin{pmatrix} a_1 + \lambda b_1 \\ a_2 + \lambda b_2 \end{pmatrix} \in \mathbb{R}^2.$$



### 2.3.3 Subspaces of a Vector Space

Subspaces are the “linear” subsets that themselves form vector spaces with the inherited operations. They are the basic building blocks for decomposing vector spaces and for understanding solution sets of linear systems.

**Definition 2.3.2** (Subspace). *Let  $V$  be a vector space over  $\mathbb{K}$ . A subset  $F \subset V$  is a subspace if*

- $0_V \in F$ ,
- for all  $x, y \in F$ , we have  $x + y \in F$  (closure under addition),
- for all  $\lambda \in \mathbb{K}$  and  $x \in F$ , we have  $\lambda x \in F$  (closure under scalar multiplication).

**Example 2.3.2.** *Here some basic examples of subspaces:*

- In  $\mathbb{R}^3$ , the set  $\{(x, y, 0) : x, y \in \mathbb{R}\}$  is a subspace (a coordinate plane).
- The sets  $\{0\}$  and  $V$  are always subspaces (the trivial subspaces).

**Proposition 2.3.1** (Intersection of subspaces). *If  $F$  and  $G$  are subspaces of  $V$ , then  $F \cap G$  is also a subspace of  $V$ .*

**Proof.** Since  $F$  and  $G$  are subspaces,  $0_V \in F$  and  $0_V \in G$ , hence  $0_V \in F \cap G$ . Let  $x, y \in F \cap G$ . Then  $x, y \in F$  and  $x, y \in G$ . Because  $F$  and  $G$  are subspaces,  $x + y \in F$  and  $x + y \in G$ , so  $x + y \in F \cap G$ . Similarly, for any  $\lambda \in \mathbb{K}$ , we have  $\lambda x \in F$  and  $\lambda x \in G$ , hence  $\lambda x \in F \cap G$ . Therefore  $F \cap G$  satisfies the three subspace axioms and is a subspace of  $V$ .  $\square$

**Proposition 2.3.2** (Sum of subspaces). *If  $F$  and  $G$  are subspaces of  $V$ , their sum is*

$$F + G = \{x + y : x \in F, y \in G\}.$$

*Then  $F + G$  is a subspace of  $V$ .*

**Proof.** First,  $0_V = 0_V + 0_V \in F + G$  since  $0_V \in F$  and  $0_V \in G$ . Let  $u, v \in F + G$ . Then there exist  $x_1, x_2 \in F$  and  $y_1, y_2 \in G$  such that  $u = x_1 + y_1$  and  $v = x_2 + y_2$ . Hence

$$u + v = (x_1 + y_1) + (x_2 + y_2) = (x_1 + x_2) + (y_1 + y_2).$$

Because  $F$  and  $G$  are subspaces,  $x_1 + x_2 \in F$  and  $y_1 + y_2 \in G$ , so  $u + v \in F + G$ . For  $\lambda \in \mathbb{K}$ , we have

$$\lambda u = \lambda(x_1 + y_1) = (\lambda x_1) + (\lambda y_1) \in F + G,$$

since  $\lambda x_1 \in F$  and  $\lambda y_1 \in G$ . Thus  $F + G$  contains  $0_V$  and is closed under addition and scalar multiplication, hence is a subspace of  $V$ .  $\square$

*Remark 2.* In general, the union  $F \cup G$  is *not* a subspace unless one subspace is contained in the other. Moreover,  $F + G$  is the smallest subspace of  $V$  containing  $F \cup G$  (i.e., any subspace that contains both  $F$  and  $G$  must also contain  $F + G$ ).

## Direct Sum of Subspaces

*Beyond the mere addition of subspaces, an important concept is when this addition is "without overlap": this leads to the notion of a **direct sum**.*

**Definition 2.3.3** (Direct Sum). *Let  $V$  be a vector space over  $\mathbb{K}$ , and let  $F, G \subset V$  be two subspaces. We say that  $V$  is the direct sum of  $F$  and  $G$ , and we write*

$$V = F \oplus G,$$

*if every vector  $v \in V$  can be written uniquely as*

$$v = f + g, \quad f \in F, \quad g \in G.$$

*Remark 3* (Equivalent characterization). The condition of uniqueness is equivalent to requiring that

$$F \cap G = \{0\}.$$

Thus, the direct sum means that the two subspaces intersect only at the zero vector, and together they generate the whole space.

**Example 2.3.3.** *In  $\mathbb{R}^2$ , the  $x$ -axis  $F = \{(x, 0)\}$  and the  $y$ -axis  $G = \{(0, y)\}$  satisfy*

$$\mathbb{R}^2 = F \oplus G.$$

*Any vector  $(a, b)$  has the unique decomposition  $(a, 0) + (0, b)$ .*

## 2.4 Span, Linear Independence, and Basis

**Definition 2.4.1** (Span of a Set). *Given a subset  $A \subset V$  of a vector space  $V$ , one can form all finite linear combinations of elements of  $A$ :*

$$\lambda_1 v_1 + \cdots + \lambda_k v_k, \quad k \in \mathbb{N}, \quad v_i \in A, \quad \lambda_i \in \mathbb{K}.$$

*The set of all such combinations is itself a subspace of  $V$ , denoted*

$$\text{Span}(A).$$

*It is called the subspace spanned by  $A$ .*

Intuitively, the span is the "smallest" subspace of  $V$  that contains the set  $A$ .

### 2.4.1 Linear Independence and Generating Families

When studying a vector space, two fundamental questions naturally arise: *which vectors are essential to describe the whole space, and which ones are redundant?*



**Definition 2.4.2** (Linear Independence). *A finite family of vectors  $(v_1, \dots, v_p)$  in  $V$  is said to be linearly independent if the only linear relation of the form*

$$\lambda_1 v_1 + \dots + \lambda_p v_p = 0$$

*is the trivial one where all coefficients vanish:  $\lambda_1 = \dots = \lambda_p = 0$ . If there exists a non-trivial relation, the family is said to be linearly dependent.*

Linear independence expresses the idea that none of the vectors in the family can be obtained by combining the others. Each vector therefore brings "new information" to the family. If a family is linearly dependent, then at least one of its members is redundant, since it can be expressed as a combination of the others.

**Definition 2.4.3** (Generating Families). *Conversely, the family  $(v_1, \dots, v_p)$  is called a generating family of  $V$  if*

$$\text{Span}(v_1, \dots, v_p) = V.$$

*This means that every vector of  $V$  can be written as a linear combination of the  $v_i$ .*

A generating family provides a complete "toolbox" for building the entire vector space. Once such a family is known, any vector of  $V$  can be produced from it. However, a generating family is not necessarily efficient: it may contain more vectors than strictly necessary, and some may even be redundant. The challenge, therefore, is to find generating families that are as small as possible, while still spanning the space.

## 2.4.2 Basis and Dimension

A **basis** of a vector space  $V$  is a family of vectors  $(v_1, \dots, v_n)$  that is both linearly independent and generating. Equivalently, a basis is a minimal generating family (no vector can be removed without losing the generating property), or a maximal independent family (no vector can be added without losing independence).

**Definition 2.4.4** (Dimension of a Space Vector). *The **dimension** of  $V$ , denoted  $\dim(V)$ , is defined as the number of vectors in any basis of  $V$ . This number is well-defined: all bases of a finite dimensional vector space have the same cardinality.*

## 2.4.3 Dimension of Subspaces

If  $F$  is a subspace of a finite-dimensional vector space  $V$ , then

$$\dim(F) \leq \dim(V),$$

with strict inequality if  $F \neq V$ . This reflects the fact that no subspace can have greater dimension than the space that contains it.

An important relation involving dimensions is the **Grassmann formula**.

**Theorem 2.4.1** (Grassmann formula). *If  $F$  and  $G$  are finite-dimensional subspaces of  $V$ , then*

$$\dim(F + G) = \dim(F) + \dim(G) - \dim(F \cap G).$$

This identity provides a way to compute the dimension of the sum of two subspaces using their intersection.

**Example 2.4.1.** Consider the set

$$\mathcal{F} = \{(1, 0), (0, 1), (2, 3)\} \subset \mathbb{R}^2.$$

- Is  $\mathcal{F}$  linearly independent?
- Does  $\mathcal{F}$  generate  $\mathbb{R}^2$ ?
- What is the dimension of  $\text{Span}(\mathcal{F})$ ?

**Example 2.4.2** (Canonical Basis). In  $\mathbb{R}^n$ , the family

$$e_1 = (1, 0, \dots, 0), \quad e_j = (0, \dots, 1, \dots, 0), \quad e_n = (0, \dots, 0, 1),$$

where the 1 appears at the  $j$ -th position, is called the canonical basis. It is linearly independent and generates  $\mathbb{R}^n$ . Therefore, it is a basis and  $\dim(\mathbb{R}^n) = n$ .

## 2.4.4 Extension and Extraction of Bases

Two fundamental theorems clarify the relationship between bases, independent families, and generating families.

**Theorem 2.4.2** (Incomplete Basis Theorem). Let  $V$  be a vector space of dimension  $n$ . If  $(v_1, \dots, v_p)$  is a linearly independent family with  $p < n$ , then there exist additional vectors  $v_{p+1}, \dots, v_n$  in  $V$  such that the enlarged family

$$(v_1, \dots, v_p, v_{p+1}, \dots, v_n)$$

is a basis of  $V$ .

In other words, any linearly independent family can be extended to a basis.

For instance, in  $\mathbb{R}^3$ , the family  $\{(1, 0, 0), (0, 1, 0)\}$  is linearly independent but not a basis. By adding  $(0, 0, 1)$ , we obtain  $\{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}$ , which is the canonical basis of  $\mathbb{R}^3$ .

**Theorem 2.4.3** (Extracted Basis Theorem). Let  $V$  be a vector space of dimension  $n$ . If  $(v_1, \dots, v_p)$  is a generating family with  $p \geq n$ , then there exists a subfamily of exactly  $n$  vectors that forms a basis of  $V$ .

In other words, any generating family contains a basis.

In  $\mathbb{R}^3$ , the family

$$\{(1, 0, 0), (0, 1, 0), (0, 0, 1), (1, 1, 1)\}$$

is generating. By removing the redundant vector  $(1, 1, 1)$ , we recover the canonical basis

$$\{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}.$$

## Chapter 3

# Linear Transformations

Linear transformations are one of the central objects of linear algebra. They are the natural maps between vector spaces, because they preserve the algebraic structure of addition and scalar multiplication. In this chapter, we introduce their definition, explore fundamental examples, and study key concepts such as kernel, image, and rank.

### 3.1 Definition and Motivation

A linear transformation can be thought of as a "structure-preserving" map: it takes vectors as inputs and produces other vectors in such a way that the two operations of the vector space are respected.

**Definition 3.1.1** (Linear Transformation). *Let  $V$  and  $W$  be vector spaces over a field  $\mathbb{K}$ . A map  $T : V \rightarrow W$  is called a linear transformation if, for all  $v_1, v_2 \in V$  and  $c \in \mathbb{K}$ ,*

$$T(cv_1 + v_2) = cT(v_1) + T(v_2).$$

*Equivalently,  $T$  preserves vector addition and scalar multiplication:*

$$T(v_1 + v_2) = T(v_1) + T(v_2), \quad T(cv) = cT(v).$$

**Remark (intuition).** A linear transformation is fully determined by its action on a basis of the domain. Once  $T$  is known on the basis vectors, its values on all other vectors follow by linearity.

**Example 3.1.1.** *Consider  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  defined by*

$$T(x, y) = (2x + y, -x + 3y).$$

*To check whether  $T$  is linear, one must verify that both additivity and homogeneity are satisfied. Expanding the definitions shows that  $T$  is indeed a linear transformation, and in fact corresponds to the matrix*

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

### 3.2 Examples in $\mathbb{R}^3$

Linear transformations can often be represented concretely by matrices. Familiar geometric transformations in  $\mathbb{R}^3$  provide useful illustrations.

**Example 3.2.1** (Canonical Examples in  $\mathbb{R}^3$ ). • **Identity:**  $\text{Id}(x, y, z) = (x, y, z)$ . Matrix:  $I_3$ .

• **Scaling (Homothety):**  $H_\alpha(x, y, z) = \alpha(x, y, z)$  for  $\alpha \in \mathbb{R}$ . Matrix:  $\alpha I_3$ .

• **Rotation about the  $z$ -axis:** for angle  $\theta$ ,

$$R_z(\theta) = \begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

which rotates the  $xy$ -plane while leaving the  $z$ -axis fixed.

### 3.3 Isomorphisms and Endomorphisms

**Definition 3.3.1** (Types of Linear Maps). Let  $V, W$  be vector spaces.

- $\mathcal{L}(V, W)$  denotes the set of linear maps  $T : V \rightarrow W$ .
- An **endomorphism** is a map in  $\mathcal{L}(V, V)$ .
- An **isomorphism** is a bijective linear transformation  $T \in \mathcal{L}(V, W)$ .
- An **automorphism** is a bijective endomorphism, i.e., an isomorphism  $T : V \rightarrow V$ .

**Algebraic structure.** The set  $\mathcal{L}(V) = \mathcal{L}(V, V)$  carries two operations: addition of maps and composition. With these,

$$(\mathcal{L}(V), +, \circ)$$

forms a ring (not necessarily commutative), with the identity map  $I_d$  as multiplicative unit. Distributivity of composition over addition holds:

$$T \circ (S_1 + S_2) = T \circ S_1 + T \circ S_2, \quad (S_1 + S_2) \circ T = S_1 \circ T + S_2 \circ T.$$

### 3.4 Kernel and Image

**Definition 3.4.1** (Kernel and Image). For  $T \in \mathcal{L}(V, W)$ , we define:

$$\ker(T) := \{v \in V : T(v) = 0_W\}, \quad \text{Im}(T) := \{T(v) : v \in V\} \subseteq W.$$

**Proposition 3.4.1** (Key properties). These equivalences give characterizations:

- $T$  is injective  $\iff \ker(T) = \{0_V\}$ .
- If  $(v_i)_{i \in I}$  generates  $V$ , then  $\text{Im}(T) = \text{Span}(T(v_i) : i \in I)$ .

### 3.5 Rank–Nullity Theorem

One of the fundamental results of linear algebra connects dimension, kernel, and image.

**Theorem 3.5.1** (Rank–Nullity Theorem). Let  $T \in \mathcal{L}(V, W)$  with  $\dim(V) < \infty$ . Then

$$\dim(\ker T) + \dim(\text{Im } T) = \dim(V).$$

**Definition 3.5.1.** Nullity and rank are simply but important notions.

- The dimension of the kernel is called the nullity of  $T$ .
- The dimension of the image is called the rank of  $T$ , and is denoted  $\text{rank}(T)$ .

**Proposition 3.5.1** (Consequences of rank-nullity theorem.). In finite dimension, we have the equivalences:

- $T$  injective  $\iff \dim(\ker T) = 0 \iff \text{rank}(T) = \dim(V)$ .
- If  $\dim(W) < \infty$ , then  $T$  surjective  $\iff \text{rank}(T) = \dim(W)$ .
- If  $\dim(V) = \dim(W)$ , then

$$T \text{ injective} \iff T \text{ surjective} \iff T \text{ is an isomorphism.}$$

**Rank-1 Linear Maps** Among all linear transformations, those of rank 1 play a very special role. They are the simplest non-trivial linear maps, and their structure can be described explicitly.

**Proposition 3.5.2** (Structure of Rank-1 Maps). Let  $T \in \mathcal{L}(\mathbb{R}^n)$  be a linear map with  $\text{rank}(T) = 1$ . Then there exist two nonzero vectors  $u, v \in \mathbb{R}^n$  such that

$$T = v u^\top, \quad \text{i.e.,} \quad T(x) = (u^\top x) v, \quad \forall x \in \mathbb{R}^n.$$

**Interpretation.** A rank-1 map has a one-dimensional image: every vector  $x$  is first "tested" against  $u$  by the scalar product  $u^\top x$ , producing a scalar, and then this scalar is used to scale the vector  $v$ . In other words:

$$\ker(T) = \{x \in \mathbb{R}^n : u^\top x = 0\}, \quad \text{Im}(T) = \text{Span}(v).$$

Thus the kernel is a hyperplane of dimension  $n - 1$ , and the image is a line through the origin.

**Remark.** Any matrix of rank 1 can be written as the outer product  $vu^\top$ , which explains why such matrices are sometimes called *outer product matrices*. They also form the basic building blocks for low-rank approximations in applications such as numerical linear algebra and data science.

**A Worked Example in  $\mathbb{R}^2$**  We illustrate the notions of kernel, image, and rank with a simple but instructive example.

*Exercise 3.5.1.* Consider the linear map  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  defined by

$$T(x, y) = (x + y, 2x + 2y).$$

1. **Linearity.** For  $u = (x_1, y_1)$ ,  $v = (x_2, y_2)$  and  $\lambda \in \mathbb{R}$ :

$$T(\lambda u + v) = T(\lambda(x_1, y_1) + (x_2, y_2)) = T(\lambda x_1 + x_2, \lambda y_1 + y_2).$$

Expanding gives

$$(\lambda x_1 + x_2 + \lambda y_1 + y_2, 2(\lambda x_1 + x_2) + 2(\lambda y_1 + y_2)),$$

which is equal to

$$\lambda(x_1 + y_1, 2x_1 + 2y_1) + (x_2 + y_2, 2x_2 + 2y_2) = \lambda T(u) + T(v).$$

Hence  $T$  is linear.

2. **Kernel and image.** We look for  $(x, y)$  such that  $T(x, y) = (0, 0)$ :

$$x + y = 0, \quad 2x + 2y = 0.$$

Both equations are equivalent to  $y = -x$ . Thus

$$\ker(T) = \text{Span}\{(1, -1)\}.$$

Next, observe that

$$T(x, y) = (x + y, 2(x + y)) = (x + y)(1, 2).$$

Hence the image is the line generated by  $(1, 2)$ :

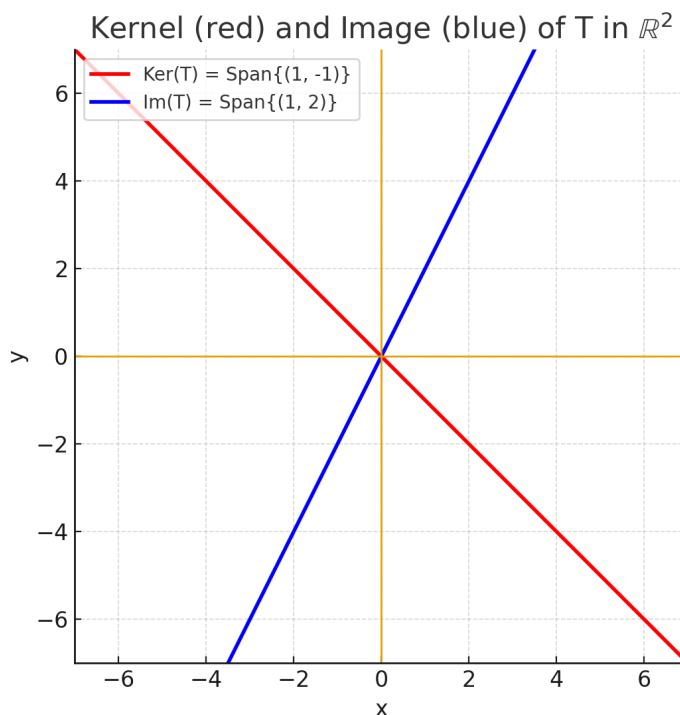
$$\text{Im}(T) = \text{Span}\{(1, 2)\}.$$

Since  $\dim \ker(T) = 1$  and  $\dim \text{Im}(T) = 1$ , we verify

$$\dim \ker(T) + \dim \text{Im}(T) = 1 + 1 = 2 = \dim(\mathbb{R}^2),$$

in accordance with the rank–nullity theorem.

3. **Geometric interpretation.** The kernel is the red line  $\text{Span}\{(1, -1)\}$ , while the image is the blue line  $\text{Span}\{(1, 2)\}$ . The transformation  $T$  "collapses" the entire plane onto the blue line: every vector is mapped onto this line, and those vectors lying on the red line are mapped to the origin.



This example shows that  $T$  is a rank–1 linear map: it reduces the two–dimensional space  $\mathbb{R}^2$  to a one–dimensional line, with the kernel providing the direction along which the collapse occurs.

## 3.6 Matrices Associated with Linear Maps

### 3.6.1 From Linear Maps to Matrices

A linear transformation between finite-dimensional vector spaces can always be represented by a matrix, once bases have been chosen in the domain and the codomain. This correspondence is fundamental, as it allows us to replace abstract linear maps by concrete arrays of numbers.

**Definition 3.6.1** (Matrix of a Linear Map). Let  $T : V \rightarrow W$  be a linear transformation, where  $V$  and  $W$  are vector spaces with bases  $\mathcal{B}_V = \{\mathbf{e}_1, \dots, \mathbf{e}_n\}$  and  $\mathcal{B}_W = \{\mathbf{f}_1, \dots, \mathbf{f}_m\}$ .

The matrix of  $T$  with respect to these bases is the  $m \times n$  matrix  $M = (M_{ij})$  defined by

$$T(\mathbf{e}_j) = \sum_{i=1}^m M_{ij} \mathbf{f}_i.$$

Equivalently, the  $j$ -th column of  $M$  is the coordinate vector of  $T(\mathbf{e}_j)$  in the basis  $\mathcal{B}_W$ .

Thus, matrices are simply the "coordinate avatars" of linear maps. This point of view will allow us to apply algebraic operations (such as matrix multiplication) to study compositions of maps.

**Example 3.6.1.** Let us compute some matrices explicitly, in the canonical bases of the respective spaces.

1.  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ ,  $T(x, y) = (2x + y, -x + 3y)$  has matrix

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

2.  $T : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ ,  $T(x, y, z) = (x + z, y, 2z)$  has matrix

$$\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 2 \end{bmatrix}.$$

3.  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ ,  $T(x, y) = (x, y, x + y)$  has matrix

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}.$$

**Semantic Interpretation** The matrix associated to a linear map is not just an array of coefficients: it encodes how the transformation acts geometrically.

- If  $T : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  is **bijective** (*i.e.*, invertible), its image is the whole space  $\mathbb{R}^3$ : any direction can be reached.
- If  $T$  is **not bijective**, then the image is a subspace of lower dimension: a plane (dimension 2), a line (dimension 1), or even just  $\{0\}$  (dimension 0). In this case,  $T$  "collapses" some directions.

These abstract ideas find immediate applications in various domains.

1. **Projection.** The projection of a 3D object onto a screen corresponds to a linear map  $\mathbb{R}^3 \rightarrow \mathbb{R}^2$ . This transformation is linear but not invertible: depth information is lost.
2. **Coordinate changes in meteorology.** Converting a wind vector  $(u, v, w)$  into rotated or polar coordinates is achieved via an invertible linear map: a rotation (*i.e.*, a change of basis).
3. **Unit conversions in mechanics or thermodynamics.** Transformations such as Joules  $\leftrightarrow$  kilocalories, or forces  $\leftrightarrow$  stresses, correspond to diagonal scaling matrices. These are invertible linear maps, since no information is lost.

Matrices are the concrete representation of linear maps in chosen bases. They not only allow explicit computations but also carry geometric meaning, revealing whether a transformation preserves dimensions, collapses directions, or stretches space.

# Chapter 4

## Change of Basis

### 4.1 Motivation

A linear transformation is an abstract object, but once a basis is chosen it can be represented concretely by a matrix. Different choices of bases yield different matrix representations of the same linear map. This naturally motivates the study of the *change of basis*: a bijective correspondence between the coordinates of the same vector expressed in two different bases.

### 4.2 General Setup

Let  $\mathcal{B}_1 = (b_1, b_2, b_3)$  and  $\mathcal{B}_2 = (c_1, c_2, c_3)$  be two bases of  $\mathbb{R}^3$ .

We seek the change-of-basis matrix  $P_{21}$  which converts coordinates from basis  $\mathcal{B}_1$  to basis  $\mathcal{B}_2$ :

$$[\mathbf{x}]_{\mathcal{B}_2} = P_{21} [\mathbf{x}]_{\mathcal{B}_1}.$$

#### Linear System Approach

Each old basis vector  $b_j$  must be expressed as a linear combination of the new basis vectors  $c_1, c_2, c_3$ . That is,

$$b_j = \alpha_{1j}c_1 + \alpha_{2j}c_2 + \alpha_{3j}c_3, \quad j = 1, 2, 3.$$

This gives three independent linear systems, one for each  $b_j$ .

#### Matrix Formulation

Let

$$C = [c_1 \ c_2 \ c_3], \quad B = [b_1 \ b_2 \ b_3],$$

and let  $A = (\alpha_{ij})_{1 \leq i, j \leq 3}$  be the coordinate matrix. The previous relations compactly read

$$CA = B.$$

Since  $C$  is invertible (since its columns form a basis), we obtain

$$A = C^{-1}B.$$

Hence the change-of-basis matrix is precisely

$$P_{21} := A = C^{-1}B.$$



## Augmented-Matrix Method

The computation of  $P_{21}$  can be carried out by Gaussian elimination. One sets up the augmented block matrix

$$[C \mid B],$$

and row-reduces it to

$$[I_d \mid C^{-1}B],$$

where  $I_d$  is the identity matrix. The right-hand block is then exactly the change-of-basis matrix  $P_{21}$ .

## 4.3 Worked Example in Dimension 3

**Example 4.3.1.** *Let*

$$B_1 = \{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}, \quad B_2 = \{(1, 1, 1), (1, -1, 1), (2, 1, 1)\}.$$

*We want the coordinates of  $\mathbf{b} = (2, 1, -1)$  in the basis  $B_2$ .*

### Method 1: Direct System

We look for  $\alpha, \beta, \gamma$  such that

$$\alpha(1, 1, 1) + \beta(1, -1, 1) + \gamma(2, 1, 1) = (2, 1, -1).$$

This yields the system

$$\begin{cases} \alpha + \beta + 2\gamma = 2, \\ \alpha - \beta + \gamma = 1, \\ \alpha + \beta + \gamma = -1. \end{cases}$$

Solving gives  $\alpha = -3$ ,  $\beta = -1$ ,  $\gamma = 3$ . Hence

$$[\mathbf{b}]_{B_2} = \begin{pmatrix} -3 \\ -1 \\ 3 \end{pmatrix}.$$

### Method 2: Matrix Inversion

Form the matrix of the new basis:

$$C = \begin{bmatrix} 1 & 1 & 2 \\ 1 & -1 & 1 \\ 1 & 1 & 1 \end{bmatrix}, \quad C^{-1} = \begin{bmatrix} -1 & \frac{1}{2} & \frac{3}{2} \\ 0 & -\frac{1}{2} & \frac{1}{2} \\ 1 & 0 & -1 \end{bmatrix}.$$

Since  $B_1$  is the canonical basis,  $P_{21} = C^{-1}$ , and thus

$$[\mathbf{b}]_{B_2} = P_{21} [\mathbf{b}]_{B_1} = \begin{bmatrix} -3 \\ -1 \\ 3 \end{bmatrix},$$

in agreement with the direct method.

**Solution.**

$$\begin{aligned}
[\mathcal{B}_2 \mid \mathcal{B}_1] &= \left[ \begin{array}{ccc|ccc} 1 & 1 & 2 & 1 & 0 & 0 \\ 1 & -1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 \end{array} \right] \sim \left[ \begin{array}{ccc|ccc} 1 & 1 & 2 & 1 & 0 & 0 \\ 0 & -2 & -1 & -1 & 1 & 0 \\ 0 & 0 & -1 & -1 & 0 & 1 \end{array} \right] \quad (L_2 \leftarrow L_2 - L_1, L_3 \leftarrow L_3 - L_1) \\
&\sim \left[ \begin{array}{ccc|ccc} 1 & 1 & 2 & 1 & 0 & 0 \\ 0 & -2 & -1 & -1 & 1 & 0 \\ 0 & 0 & 1 & 1 & 0 & -1 \end{array} \right] \quad (L_3 \leftarrow -L_3) \\
&\sim \left[ \begin{array}{ccc|ccc} 1 & 1 & 2 & 1 & 0 & 0 \\ 0 & -2 & 0 & 0 & 1 & -1 \\ 0 & 0 & 1 & 1 & 0 & -1 \end{array} \right] \quad (L_2 \leftarrow L_2 + L_3) \\
&\sim \left[ \begin{array}{ccc|ccc} 1 & 1 & 2 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & -\frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 1 & 1 & 0 & -1 \end{array} \right] \quad (L_2 \leftarrow -\frac{1}{2}L_2) \\
&\sim \left[ \begin{array}{ccc|ccc} 1 & 0 & 0 & -1 & \frac{1}{2} & \frac{3}{2} \\ 0 & 1 & 0 & 0 & -\frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 1 & 1 & 0 & -1 \end{array} \right] \quad (L_1 \leftarrow L_1 - L_2 - 2L_3)
\end{aligned}$$

□

## 4.4 Geometric Interpretation

The same vector has different coordinate representations depending on the chosen basis. The change-of-basis matrix allows us to pass from one representation to another without ambiguity.

**Example 4.4.1** (Exercise). *Let*

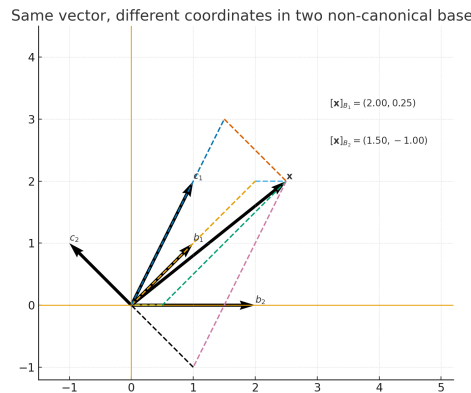
$$\mathcal{B}_1 = \left\{ \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 2 \\ 0 \end{pmatrix} \right\}, \quad \mathcal{B}_2 = \left\{ \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \begin{pmatrix} -1 \\ 1 \end{pmatrix} \right\},$$

*and suppose*

$$[\mathbf{x}]_{\mathcal{B}_1} = \begin{pmatrix} 2 \\ 0.25 \end{pmatrix}.$$

*Compute the change-of-basis matrix  $P_{21}$  and deduce  $[\mathbf{x}]_{\mathcal{B}_2}$ .*

The change of basis is a fundamental operation in linear algebra: it allows us to reinterpret the same vector in different coordinate systems, and to relate different matrix representations of the same linear transformation.



## 4.5 Statement and semonstration

**Proposition 4.5.1.** *Let  $\mathcal{B}_1 = (b_1, \dots, b_n)$  and  $\mathcal{B}_2 = (c_1, \dots, c_n)$  be two bases of the same vector space  $V$ . The change-of-basis matrix from  $\mathcal{B}_2$  to  $\mathcal{B}_1$ , denoted  $P$ , is defined by*

$$c_j = \sum_{i=1}^n p_{ij} b_i \quad (j = 1, \dots, n).$$

*Then, for every  $x \in V$ ,*

$$[x]_{\mathcal{B}_1} = P [x]_{\mathcal{B}_2}.$$

**Proof.** First, write each vector of the basis  $\mathcal{B}_2$  in terms of the basis  $\mathcal{B}_1$ : for each  $j$ , there exist scalars  $(p_{1j}, \dots, p_{nj})$  such that

$$c_j = \sum_{i=1}^n p_{ij} b_i.$$

By definition, the matrix  $P = (p_{ij})_{1 \leq i, j \leq n}$  has its  $j$ -th column equal to the coordinate vector of  $c_j$  with respect to  $\mathcal{B}_1$ , i.e.,  $[c_j]_{\mathcal{B}_1}$ .

Now let  $x \in V$ . Write  $x$  in the basis  $\mathcal{B}_2$ :

$$x = \sum_{j=1}^n \alpha_j c_j, \quad \alpha_j \in \mathbb{K},$$

so that  $[x]_{\mathcal{B}_2} = (\alpha_1, \dots, \alpha_n)^\top$ .

Substituting the expression of each  $c_j$  in terms of  $\mathcal{B}_1$ , we get

$$x = \sum_{j=1}^n \alpha_j \left( \sum_{i=1}^n p_{ij} b_i \right) = \sum_{i=1}^n \left( \sum_{j=1}^n p_{ij} \alpha_j \right) b_i.$$

By uniqueness of coordinates in a basis, the coefficient of  $b_i$  in the basis  $\mathcal{B}_1$  is

$$([x]_{\mathcal{B}_1})_i = \sum_{j=1}^n p_{ij} \alpha_j.$$

In matrix notation, this is exactly

$$[x]_{\mathcal{B}_1} = P [x]_{\mathcal{B}_2}.$$

Thus completing the proof.  $\square$

*Remark 4.* Since  $P$  maps one basis to another, it is invertible. Its inverse  $P^{-1}$  is the change-of-basis matrix from  $\mathcal{B}_1$  to  $\mathcal{B}_2$  and satisfies

$$[x]_{\mathcal{B}_2} = P^{-1} [x]_{\mathcal{B}_1}$$

for every  $x \in V$ .

**Proposition 4.5.2** (Change of basis). *Let  $T : V \rightarrow V$  be a linear endomorphism on a finite-dimensional vector space  $V$ . Let  $\mathcal{B}_1 = (b_1, \dots, b_n)$  and  $\mathcal{B}_2 = (c_1, \dots, c_n)$  be two bases of  $V$ . Denote by*

$$M_{11} = [T]_{\mathcal{B}_1 \rightarrow \mathcal{B}_1}, \quad M_{22} = [T]_{\mathcal{B}_2 \rightarrow \mathcal{B}_2}$$

*the matrices of  $T$  in these bases. Let  $P$  be the change-of-basis matrix from  $\mathcal{B}_2$ -coordinates to  $\mathcal{B}_1$ -coordinates, i.e.*

$$[x]_{\mathcal{B}_1} = P [x]_{\mathcal{B}_2} \quad \text{for all } x \in V.$$

*Then*

$$M_{22} = P^{-1} M_{11} P \quad (\text{equivalently } M_{11} = P M_{22} P^{-1}).$$

**Proof.** By definition of matrix of a linear map in a basis, for any  $x \in V$  we have

$$[T(x)]_{\mathcal{B}_1} = M_{11} [x]_{\mathcal{B}_1}, \quad [T(x)]_{\mathcal{B}_2} = M_{22} [x]_{\mathcal{B}_2}.$$

Using the relation  $[x]_{\mathcal{B}_1} = P[x]_{\mathcal{B}_2}$  and linearity,

$$[T(x)]_{\mathcal{B}_1} = M_{11} [x]_{\mathcal{B}_1} = M_{11} P [x]_{\mathcal{B}_2}.$$

On the other hand,

$$[T(x)]_{\mathcal{B}_1} = P [T(x)]_{\mathcal{B}_2} = P M_{22} [x]_{\mathcal{B}_2}.$$

Since this holds for all  $[x]_{\mathcal{B}_2}$ , we get  $M_{11}P = PM_{22}$ , hence  $M_{22} = P^{-1}M_{11}P$  (because  $P$  is invertible as a change-of-basis matrix).  $\square$

*Remark 5.* Two matrices representing the same linear map in different bases are *similar*. Similarity preserves the characteristic polynomial, eigenvalues, trace, determinant, rank, and (more generally) the Jordan canonical structure.

# Chapter 5

## Diagonalization

### 5.1 Motivation: Why Diagonalize?

In many applications, we want to understand the long-term behavior of a system described by repeated applications of a linear map, *i.e.*, by powers of a matrix  $A^t$ . For instance, in population models (like the "Smallville" example 5.7.1), we want to know:

- What happens after  $n$  steps?
- Does the system converge to an equilibrium?

The central idea: matrices represent linear maps *in a chosen basis*. By selecting a more convenient basis (made of eigenvectors), the matrix may become diagonal, making computations much easier.

### 5.2 Eigenvalues and Eigenvectors

**Definition 5.2.1** (Eigenvalue and Eigenvector). *Let  $T : V \rightarrow V$  be a linear operator on  $V = \mathbb{R}^n$ . A scalar  $\lambda \in \mathbb{R}$  is an eigenvalue of  $T$  if there exists a nonzero vector  $\mathbf{v} \in V$  such that*

$$T(\mathbf{v}) = \lambda \mathbf{v}.$$

*Such a nonzero vector  $\mathbf{v}$  is called an eigenvector associated with  $\lambda$ .*

*Remark 6.* If a basis of  $V$  can be formed entirely of eigenvectors of  $T$ , then the matrix of  $T$  in this basis is diagonal, with the eigenvalues on the diagonal. This is the essence of diagonalization.

### 5.3 Characteristic Polynomial and Eigenspace

**Definition 5.3.1** (Characteristic Polynomial). *For a square matrix  $A \in M_n(\mathbb{R})$ , the polynomial*

$$\chi_A(\lambda) = \det(A - \lambda I_n)$$

*is called the characteristic polynomial of  $A$ .*

*Remark 7.* A scalar  $\lambda$  is an eigenvalue of  $A$  if and only if it is a root of  $\chi_A(\lambda)$ . Indeed,

$$A\mathbf{v} = \lambda \mathbf{v} \iff (A - \lambda I)\mathbf{v} = 0,$$

and this has a nonzero solution precisely when  $A - \lambda I$  is singular, *i.e.*,  $\det(A - \lambda I) = 0$ .

**Definition 5.3.2** (Eigenspace). *Given an eigenvalue  $\lambda$ , the associated eigenspace is*

$$E_\lambda = \ker(A - \lambda I) = \{ \mathbf{v} \in \mathbb{R}^n \mid A\mathbf{v} = \lambda\mathbf{v} \}.$$

*It is a subspace of  $\mathbb{R}^n$ , spanned by all eigenvectors of  $A$  with eigenvalue  $\lambda$ .*

**Definition 5.3.3** (Algebraic vs. Geometric Multiplicity). *Let  $\lambda$  be an eigenvalue of  $A$ .*

- *The algebraic multiplicity of  $\lambda$  is its multiplicity as a root of  $\chi_A(\lambda)$ .*
- *The geometric multiplicity of  $\lambda$  is  $\dim(E_\lambda)$ .*

**Property 5.3.1** (Diagonalizability criterion). *We always have:*

$$1 \leq \dim(E_\lambda) \leq \text{algebraic multiplicity of } \lambda.$$

*Therefore, a matrix  $A$  is **diagonalizable** if and only if, for each eigenvalue  $\lambda$ , the **geometric and algebraic multiplicities for each eigenvalue coincide**, and the **sum of dimensions of eigenspaces is  $n$** .*

Here we propose a worked example:

**Example 5.3.1.** *Let*

$$A = \begin{bmatrix} 7 & 2 \\ 3 & 8 \end{bmatrix}.$$

*We compute its characteristic polynomial:*

$$\chi_A(\lambda) = \det \begin{bmatrix} 7 - \lambda & 2 \\ 3 & 8 - \lambda \end{bmatrix} = (7 - \lambda)(8 - \lambda) - 6.$$

*Simplifying,*

$$\chi_A(\lambda) = \lambda^2 - 15\lambda + 50.$$

*The roots are  $\lambda = 10$  and  $\lambda = 5$ , hence the eigenvalues are 5 and 10.*

*For  $\lambda = 10$ , we need to solve  $(A - 10I)\mathbf{v} = 0$ :*

$$\begin{bmatrix} -3 & 2 \\ 3 & -2 \end{bmatrix} \mathbf{v} = 0.$$

*This gives the relation  $-3x + 2y = 0$ , i.e.,  $y = \frac{3}{2}x$ . So  $E_{10} = \text{Span}\{(2, 3)\}$ .*

*For  $\lambda = 5$ , solve  $(A - 5I)\mathbf{v} = 0$ :*

$$\begin{bmatrix} 2 & 2 \\ 3 & 3 \end{bmatrix} \mathbf{v} = 0,$$

*giving  $x = -y$ . So  $E_5 = \text{Span}\{(1, -1)\}$ .*

We then propose some examples to study:

**Example 5.3.2.** Consider the matrices

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

Determine which are diagonalizable.

**Solution.** *Case A.*  $A$  is upper triangular, so its eigenvalues are on the diagonal: 1, 7, 2 (all distinct). Distinct eigenvalues always yield linearly independent eigenvectors, hence  $A$  is diagonalizable.

*Case B.* Eigenvalues: 7 with algebraic multiplicity 3. Compute  $E_7 = \ker(B - 7I)$ . We have

$$B - 7I = \begin{bmatrix} 0 & 2 & 3 \\ 0 & 0 & -1 \\ 0 & 0 & 0 \end{bmatrix}.$$

The rank is 2, so  $\dim(E_7) = 1$ . Since the geometric multiplicity 1 is strictly less than the algebraic multiplicity 3,  $B$  is *not* diagonalizable.

*Case C.* The rows are linearly dependent, so  $\text{rank}(C) = 1$ . Its characteristic polynomial is  $\chi_C(\lambda) = \lambda^2(\lambda + 1)$ . Hence eigenvalues:  $\lambda = 0$  (algebraic multiplicity 2) and  $\lambda = -1$  (algebraic multiplicity 1). The null space of  $C$  has dimension 2, so the geometric multiplicity of  $\lambda = 0$  is 2. Thus the sum of geometric multiplicities is  $2 + 1 = 3$ , which equals the size of the matrix. Therefore  $C$  is diagonalizable, even though  $C$  itself is not invertible.  $\square$

## 5.4 Change of Basis and Similarity

**Definition 5.4.1** (Matrices in Different Bases). Let  $T : V \rightarrow V$  be a linear endomorphism of a finite-dimensional vector space  $V$ . If  $\mathcal{B}_1$  and  $\mathcal{B}_2$  are two bases of  $V$ , then the matrix of  $T$  depends on the chosen basis. We denote by  $M_{\mathcal{B}_1\mathcal{B}_1}$  and  $M_{\mathcal{B}_2\mathcal{B}_2}$  the matrices of  $T$  expressed in these bases.

**Proposition 5.4.1** (Change of Basis Relation). Let  $P_{12}$  be the change-of-basis matrix from  $\mathcal{B}_2$  to  $\mathcal{B}_1$ , so that

$$[\mathbf{v}]_{\mathcal{B}_1} = P_{12} [\mathbf{v}]_{\mathcal{B}_2}.$$

Then the two matrices of  $T$  are related by

$$M_{\mathcal{B}_2\mathcal{B}_2} = P_{21} M_{\mathcal{B}_1\mathcal{B}_1} P_{12}, \quad \text{with } P_{21} = P_{12}^{-1}.$$

**Definition 5.4.2** (Similarity). Two matrices  $A, B \in M_n(\mathbb{R})$  are said to be similar if there exists an invertible matrix  $P$  such that

$$B = P^{-1}AP.$$

In this case,  $A$  and  $B$  represent the same linear transformation in two different bases. This can be denoted as

$$A \sim B$$

*Remark 8.* Similarity preserves many invariants: the characteristic polynomial, the eigenvalues, the determinant, the trace, the rank, and more. However, the explicit entries of the matrices may differ, depending on the chosen basis.

## 5.5 Diagonalization Formalism

**Definition 5.5.1** (Diagonalization). A matrix  $A \in M_n(\mathbb{R})$  is said to be diagonalizable if there exists a basis of  $\mathbb{R}^n$  consisting entirely of eigenvectors of  $A$ . In this eigenbasis, the matrix of  $A$  is diagonal:

$$D = \text{diag}(\lambda_1, \dots, \lambda_n),$$

where  $\lambda_i$  are the eigenvalues of  $A$ .

**Proposition 5.5.1** (Change of Basis for Diagonalization). Let  $P$  be the change-of-basis matrix from the eigenbasis to the reference basis. Then

$$A = PDP^{-1}.$$

**Theorem 5.5.1** (Powers of a Diagonalizable Matrix). If  $A = PDP^{-1}$ , then for all integers  $p \geq 0$ ,

$$A^p = PD^pP^{-1}.$$

Since  $D^p = \text{diag}(\lambda_1^p, \dots, \lambda_n^p)$ , powers of  $A$  are easily computed.

*Remark 9.* This is one of the main motivations for diagonalization: simplifying computations involving matrix powers, exponentials, and long-term behaviors of dynamical systems.

**Theorem 5.5.2** (Characterizations of Diagonalizability). For  $A \in M_n(\mathbb{R})$ , the following are equivalent:

- $A$  is diagonalizable  $\iff$  there exists a basis of  $\mathbb{R}^n$  made of eigenvectors of  $A$ .
- $A$  is similar to a diagonal matrix:

$$\exists P \text{ invertible, } P^{-1}AP = D.$$

- The sum of the dimensions of all eigenspaces of  $A$  equals  $n$ .

*Remark 10.* A practical criterion: if  $A$  has  $n$  distinct eigenvalues, then it is diagonalizable. If some eigenvalues are repeated, one must check whether their geometric multiplicities equal their algebraic multiplicities.

**Eigenspaces: Direct Sums and Diagonalizability** A fundamental property of eigenspaces is that, when they correspond to distinct eigenvalues, they behave very well with respect to linear independence and direct sums.

**Proposition 5.5.2.** Let  $A \in \mathbb{R}^{n \times n}$  and let  $\lambda_1, \dots, \lambda_k$  be distinct eigenvalues of  $A$ , with corresponding eigenspaces  $E_{\lambda_1}, \dots, E_{\lambda_k}$ . Then:

1. The family  $(E_{\lambda_1}, \dots, E_{\lambda_k})$  is linearly independent, i.e.

$$E_{\lambda_i} \cap \left( \sum_{j \neq i} E_{\lambda_j} \right) = \{0\}, \quad \text{for each } i.$$

2. Equivalently, the sum of these subspaces is a **direct sum**:

$$E_{\lambda_1} \oplus \dots \oplus E_{\lambda_k}.$$



**Sketch of proof.** Suppose  $v_1 + \cdots + v_k = 0$  with  $v_i \in E_{\lambda_i}$ . Applying  $A$  to this relation gives

$$\lambda_1 v_1 + \cdots + \lambda_k v_k = 0.$$

By combining these relations and using the distinctness of the eigenvalues, one shows successively that each  $v_i = 0$ . Hence the only linear relation is the trivial one, which proves independence.  $\square$

**Proposition 5.5.3** (A corollary). *If the sum of the dimensions of the eigenspaces equals  $n$ , i.e.*

$$\dim(E_{\lambda_1}) + \cdots + \dim(E_{\lambda_k}) = n,$$

*then*

$$\mathbb{R}^n = E_{\lambda_1} \oplus \cdots \oplus E_{\lambda_k}.$$

*In this case,  $A$  admits a basis of eigenvectors and is therefore **diagonalizable**.*

*Remark 11.* Diagonalizability occurs precisely when, for every eigenvalue, the geometric multiplicity (dimension of its eigenspace) equals its algebraic multiplicity (multiplicity as a root of the characteristic polynomial). In such a case, the direct sum of all eigenspaces gives the whole space.

## 5.6 Important Operators on Matrices

### The Trace of a Matrix

**Definition 5.6.1** (Trace). *For a square matrix  $A = (a_{ij}) \in M_n(\mathbb{R})$ , the trace is defined as*

$$\text{tr}(A) = \sum_{i=1}^n a_{ii}.$$

*Remark 12* (Intuition and Properties). The trace has several fundamental properties:

1. **Invariant under change of basis:** for any invertible  $P$ ,  $\text{tr}(P^{-1}AP) = \text{tr}(A)$ .
2. **Spectral link:** the trace equals the sum of eigenvalues (counted with multiplicity),  $\text{tr}(A) = \sum_{i=1}^n \lambda_i$ .
3. **Cyclic property:** for any  $A, B \in M_n(\mathbb{R})$ ,  $\text{tr}(AB) = \text{tr}(BA)$ .
4. **Associated norm:** the Frobenius norm satisfies  $\|A\|_F^2 = \text{tr}(A^\top A)$ .
5. **Geometric interpretation:** if  $A$  is the Jacobian of a linear vector field, then  $\text{tr}(A)$  is its divergence.

*Remark 13* (Historical Note). The term *trace* comes from the German word *Spur* (“footprint, track”), introduced in the 19<sup>th</sup> century by von Staudt and later used by Frobenius. The idea is that the trace is the *footprint* left by a linear transformation, independent of the chosen basis.

**Proposition 5.6.1** (Trace of a Diagonalizable Matrix). *If  $A$  is diagonalizable, i.e.,  $A = PDP^{-1}$  with  $D = \text{diag}(\lambda_1, \dots, \lambda_n)$ , then*

$$\text{tr}(A) = \text{tr}(D) = \lambda_1 + \cdots + \lambda_n.$$

## The Determinant of a Matrix

**Definition 5.6.2** (Determinant). *The determinant of a square matrix  $A = (a_{ij}) \in M_n(\mathbb{R})$  is a scalar denoted  $\det(A)$ .*

- For  $n = 1$ ,  $\det([a_{11}]) = a_{11}$ .
- For  $n \geq 2$ , it is defined recursively by the Laplace expansion:

$$\det(A) = \sum_{j=1}^n (-1)^{1+j} a_{1j} \det(M_{1j}),$$

where  $M_{1j}$  is the  $(n-1) \times (n-1)$  minor obtained by deleting row 1 and column  $j$ .

*Remark 14* (Geometric Meaning). The determinant measures the volume scaling factor of the linear transformation associated with  $A$ , with the sign encoding whether orientation is preserved ( $\det(A) > 0$ ) or reversed ( $\det(A) < 0$ ).

*Remark 15* (Properties of the Determinant). Let  $A, B \in M_n(\mathbb{R})$  and  $\lambda \in \mathbb{R}$ . The determinant satisfies:

1. **Multiplicativity:**  $\det(AB) = \det(A) \det(B)$ .
2. **Transpose invariance:**  $\det(A^\top) = \det(A)$ .
3. **Invertibility criterion:**  $A$  is invertible  $\iff \det(A) \neq 0$ , and in this case  $\det(A^{-1}) = \frac{1}{\det(A)}$ .
4. **Effect of scaling:** For scalar multiplication,  $\det(\lambda A) = \lambda^n \det(A)$ .
5. **Row/column operations:**
  - Swapping two rows (or columns) multiplies  $\det(A)$  by  $-1$ .
  - Multiplying one row (or column) by  $\alpha$  multiplies  $\det(A)$  by  $\alpha$ .
  - Adding a multiple of one row (or column) to another leaves  $\det(A)$  unchanged.
6. **Triangular matrices:** If  $A$  is triangular (upper or lower), then  $\det(A) = a_{11}a_{22} \cdots a_{nn}$ , i.e., the product of diagonal entries.

**Proposition 5.6.2** (Determinant and Eigenvalues). *If  $A$  is diagonalizable as  $A = PDP^{-1}$  with  $D = \text{diag}(\lambda_1, \dots, \lambda_n)$ , then*

$$\det(A) = \det(D) = \lambda_1 \cdot \lambda_2 \cdots \lambda_n.$$

**Spectrum and the Transpose** Combining these results, we obtain the following important observation:

**Proposition 5.6.3** (Spectrum Invariance). *Since*

$$\det(A - \lambda I) = \det((A - \lambda I)^\top) = \det(A^\top - \lambda I),$$

*we deduce that*

$$\text{Sp}(A) = \text{Sp}(A^\top).$$

## 5.7 Application: Equilibrium States

In this section, we illustrate how diagonalization provides a powerful tool for studying the long-term behavior of linear dynamical systems. A classical example comes from modeling population transitions with a stochastic matrix.

**The Smallville Example** We consider again the population of Smallville 5.7.1, split into two categories: nonsmokers and smokers. Each year, a fraction of people change category. The transition is described by the matrix

$$A = \begin{bmatrix} 0.7 & 0.2 \\ 0.3 & 0.8 \end{bmatrix}.$$

If  $n_t$  and  $s_t$  denote the number of nonsmokers and smokers at year  $t$ , then

$$\begin{bmatrix} n_{t+1} \\ s_{t+1} \end{bmatrix} = A \begin{bmatrix} n_t \\ s_t \end{bmatrix}.$$

Starting from

$$\begin{bmatrix} n_0 \\ s_0 \end{bmatrix} = \begin{bmatrix} 2000 \\ 8000 \end{bmatrix},$$

we ask: what happens as  $t \rightarrow \infty$ ?

**Proposition 5.7.1.** *We have*

$$\lim_{t \rightarrow \infty} A^t = \begin{bmatrix} 2/5 & 2/5 \\ 3/5 & 3/5 \end{bmatrix}.$$

**Solution.** The eigenvalues of  $A$  are obtained from the characteristic polynomial

$$\det(A - \lambda I) = (\lambda - 1)(\lambda - 0.5).$$

Thus,  $\lambda_1 = 1$  and  $\lambda_2 = 0.5$ . The eigenspaces are

$$E_{\lambda=1} = \text{Span} \begin{bmatrix} 2 \\ 3 \end{bmatrix}, \quad E_{\lambda=0.5} = \text{Span} \begin{bmatrix} 1 \\ -1 \end{bmatrix}.$$

Since  $|\lambda_2| < 1$ , the component along  $E_{\lambda=0.5}$  vanishes as  $t \rightarrow \infty$ . Hence,

$$\lim_{t \rightarrow \infty} A^t = \frac{1}{5} \begin{bmatrix} 2 & 2 \\ 3 & 3 \end{bmatrix}.$$

Applying this matrix to the initial state gives

$$\begin{bmatrix} 2/5 & 2/5 \\ 3/5 & 3/5 \end{bmatrix} \begin{bmatrix} 2000 \\ 8000 \end{bmatrix} = \begin{bmatrix} 4000 \\ 6000 \end{bmatrix}.$$

Therefore, the equilibrium population is 4000 nonsmokers and 6000 smokers. □

This illustrates the central role of the dominant eigenvalue  $\lambda = 1$  in stochastic matrices: it governs the equilibrium distribution.

**A Worked Example in Dimension Three** We now turn to a more algebraic problem, focusing on diagonalization of a  $3 \times 3$  matrix.

**Example 5.7.1.** *Consider the linear map  $f : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  represented in the canonical basis by*

$$M = \begin{pmatrix} 3 & 1 & 3 \\ 1 & 3 & 3 \\ 3 & 3 & 1 \end{pmatrix}.$$

1. *Compute the eigenvalues of  $M$ .*
2. *For each eigenvalue, determine a basis of the corresponding eigenspace.*
3. *Decide whether  $M$  is diagonalizable, and if so, exhibit an explicit diagonalization.*

**Solution. Step 1. Eigenvalues.** We compute the characteristic polynomial:

$$\det(M - \lambda I) = -(\lambda - 7)(\lambda - 1)(\lambda + 4).$$

Hence, the eigenvalues are  $\lambda_1 = 7$ ,  $\lambda_2 = 1$ ,  $\lambda_3 = -4$ .

**Step 2. Eigenspaces.**

- For  $\lambda = 7$ : solving  $(M - 7I)x = 0$  yields  $E_7 = \text{Span}\{(1, 1, 1)\}$ .
- For  $\lambda = 1$ : solving  $(M - I)x = 0$  yields  $E_1 = \text{Span}\{(-1, 1, 0)\}$ .
- For  $\lambda = -4$ : solving  $(M + 4I)x = 0$  yields  $E_{-4} = \text{Span}\{(-1, 0, 1)\}$ .

**Step 3. Diagonalization.** The three eigenvectors are linearly independent. Therefore  $M$  is diagonalizable, with

$$P = \begin{pmatrix} 1 & -1 & -1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix}, \quad D = \text{diag}(7, 1, -4).$$

We obtain

$$M = PDP^{-1}.$$

□

This example confirms that once distinct eigenvalues are found, the associated eigenvectors directly provide a diagonalization, greatly simplifying the computation of powers of  $M$ .

## 5.8 Orthogonal Matrices and the Spectral Theorem

*Remark 16.* If  $v$  is an eigenvector of a matrix  $A$  with eigenvalue  $\lambda$ , then its normalized version  $u = \frac{v}{\|v\|}$  is also an eigenvector associated with the same eigenvalue. Since  $Av = \lambda v$ , we compute

$$Au = A\left(\frac{v}{\|v\|}\right) = \frac{1}{\|v\|}Av = \frac{1}{\|v\|}(\lambda v) = \lambda\left(\frac{v}{\|v\|}\right) = \lambda u.$$

Thus  $u$  is an eigenvector with eigenvalue  $\lambda$ , and it has unit norm by construction.

### 5.8.1 Orthogonal Matrices

**Definition 5.8.1.** A square matrix  $Q \in \mathbb{R}^{n \times n}$  is called **orthogonal** if

$$Q^\top Q = QQ^\top = I_n.$$

Equivalently,  $Q^{-1} = Q^\top$ .

*Remark 17 (Geometric Interpretation).* Orthogonal matrices represent isometries of the Euclidean space. In other words, they preserve inner products, lengths, and angles. Geometrically, any orthogonal transformation is a composition of rotations and reflections.

**Proposition 5.8.1** (Isometry Property). For any vector  $\mathbf{x} \in \mathbb{R}^n$  and any orthogonal matrix  $Q$ ,

$$\|Q\mathbf{x}\|^2 = (Q\mathbf{x})^\top(Q\mathbf{x}) = \mathbf{x}^\top Q^\top Q\mathbf{x} = \mathbf{x}^\top \mathbf{x} = \|\mathbf{x}\|^2.$$

Thus, orthogonal transformations preserve Euclidean norms.

*Remark 18 (Orthonormal Basis).* If  $Q = [\mathbf{q}_1 \ \dots \ \mathbf{q}_n]$  is orthogonal, then its columns form an orthonormal basis of  $\mathbb{R}^n$ , since

$$\mathbf{q}_i \cdot \mathbf{q}_j = \delta_{ij}.$$

The same property holds for the rows of  $Q$ .

**Proof.** By definition,  $Q$  is orthogonal if  $Q^\top Q = I_n$ , leading the  $(i, j)$ -entry of  $Q^\top Q$  to be equal to

$$(Q^\top Q)_{ij} = \mathbf{q}_i^\top \mathbf{q}_j = \mathbf{q}_i \cdot \mathbf{q}_j.$$

But  $Q^\top Q = I_n$  means that

$$(Q^\top Q)_{ij} = \delta_{ij} = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j. \end{cases}$$

Thus the columns of  $Q$  form an orthonormal family of  $n$  vectors in  $\mathbb{R}^n$ , which is therefore an orthonormal basis. The same reasoning applies to the rows, since  $QQ^\top = I_n$ .  $\square$

**Example 5.8.1.** The rotation matrix in  $\mathbb{R}^2$ ,

$$R(\theta) = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix},$$

satisfies  $R(\theta)^\top R(\theta) = I_2$ . Hence, it is orthogonal.

*Remark 19 (Orthogonal Group).* The set of all orthogonal matrices forms a group under matrix multiplication, called the **orthogonal group**  $O(n)$ . The subgroup of orthogonal matrices with determinant 1 is denoted  $SO(n)$  and corresponds to pure rotations.

## 5.8.2 The Spectral Theorem

**Theorem 5.8.1** (Spectral Theorem). *Every real symmetric matrix is diagonalizable by an orthogonal change of basis. That is, if  $A \in \mathbb{R}^{n \times n}$  with  $A^\top = A$ , then there exists an orthogonal matrix  $Q$  such that*

$$Q^\top A Q = D,$$

*where  $D$  is diagonal. In particular,  $A$  admits an orthonormal basis of eigenvectors.*

**Idea of the proof.** The key steps are:

- Real symmetric matrices always have real eigenvalues.
- Eigenvectors corresponding to distinct eigenvalues are orthogonal.
- By choosing orthonormal eigenvectors in each eigenspace, one can construct an orthogonal basis of eigenvectors for  $\mathbb{R}^n$ .

$\square$

**Example 5.8.2.** Consider

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

Clearly  $A^\top = A$ , so  $A$  is symmetric.

- The characteristic polynomial is  $\det(A - \lambda I) = (\lambda - 3)(\lambda - 1)$ , giving eigenvalues  $\lambda_1 = 3$  and  $\lambda_2 = 1$ .
- For  $\lambda = 3$ , an eigenvector is  $(1, 1)$ ; for  $\lambda = 1$ , an eigenvector is  $(1, -1)$ .
- After normalization, we obtain an orthonormal eigenbasis:

$$u_1 = \frac{1}{\sqrt{2}}(1, 1), \quad u_2 = \frac{1}{\sqrt{2}}(1, -1).$$

- Taking  $Q = [u_1 \ u_2]$ , we check that  $Q^\top A Q = \text{diag}(3, 1)$ .

Thus  $A$  is diagonalizable by an orthogonal matrix.

*Remark 20.* The Spectral Theorem is a cornerstone of linear algebra. It explains why real symmetric matrices are so well-behaved: their eigenvalues are real, their eigenspaces are mutually orthogonal, and they admit an orthonormal basis of eigenvectors. As a consequence, symmetric matrices are always diagonalizable, unlike general square matrices.

For instance, the matrix  $M$  from Example 5.7.1 satisfies  $M^\top = M$ , and indeed it was diagonalizable with real eigenvalues.

## 5.9 Statement and demonstration

### 5.9.1 Diagonalization criterion

**Theorem 5.9.1** (Diagonalization criterion). Let  $A \in M_n(\mathbb{K})$  and let  $\lambda_1, \dots, \lambda_s$  be its distinct eigenvalues. For each  $i$ , denote by

$\alpha_i$  the algebraic multiplicity of  $\lambda_i$  in  $\chi_A$ ,  $g_i = \dim E_{\lambda_i}$  the geometric multiplicity.

Then:

1. For every  $i$ ,  $1 \leq g_i \leq \alpha_i$ .
2.  $A$  is diagonalizable  $\iff \sum_{i=1}^s g_i = n$  (equivalently,  $g_i = \alpha_i$  for all  $i$ ).

In that case, there exists an invertible  $P$  and a diagonal  $D = \text{diag}(\lambda_1, \dots, \lambda_s)$  such that

$$A = P D P^{-1}.$$

**Proof.** (1) *Inequality*  $g_i \leq \alpha_i$ . Fix  $\lambda \in \{\lambda_1, \dots, \lambda_s\}$  and set  $E_\lambda = \ker(A - \lambda I)$ . Choose a basis  $(v_1, \dots, v_g)$  of  $E_\lambda$  and extend it to a basis of  $V$ , say  $(v_1, \dots, v_g, w_{g+1}, \dots, w_n)$ . In this basis the matrix of  $A - \lambda I$  has the block form

$$[A - \lambda I] = \begin{bmatrix} 0_{g \times g} & * \\ 0 & B \end{bmatrix}.$$

Hence  $\chi_{A-\lambda I}(t) = t^g \chi_B(t)$ , and thus  $\chi_A(t) = \det(A - \lambda I - tI) = (t + \lambda)^g p(t)$  for some polynomial  $p$ . Therefore the factor  $(t - \lambda)$  appears at least  $g$  times in  $\chi_A$ , i.e.,  $\alpha_\lambda \geq g$ .

(2) *Sufficiency.* Assume  $\sum_i g_i = n$ . Since eigenspaces for distinct eigenvalues are linearly independent (a standard result), the direct sum  $E_{\lambda_1} \oplus \dots \oplus E_{\lambda_s}$  has dimension  $\sum_i g_i = n$ , hence equals  $V$ . Choosing a basis of each eigenspace and concatenating them gives a basis of  $V$  made entirely of eigenvectors. In that basis the matrix of  $A$  is diagonal with the corresponding eigenvalues on the diagonal. Thus  $A$  is diagonalizable.

(3) *Necessity.* Assume  $A$  is diagonalizable:  $A = PDP^{-1}$  with  $D = \text{diag}(\underbrace{\lambda_1, \dots, \lambda_1}_{\alpha_1}, \dots, \underbrace{\lambda_s, \dots, \lambda_s}_{\alpha_s})$ .

For a diagonal matrix, the eigenspace for  $\lambda_i$  has dimension exactly the number of times  $\lambda_i$  appears on the diagonal, namely  $\alpha_i$ . Similarity preserves eigenspace dimensions (indeed  $E_{\lambda_i}(A) = P E_{\lambda_i}(D)$ ), hence  $g_i = \alpha_i$  for all  $i$ , and in particular  $\sum_i g_i = \sum_i \alpha_i = n$ .  $\square$

**Property 5.9.1** (Diagonalization formula). *If  $A$  is diagonalizable, there exists an invertible matrix  $P$  whose columns form a basis of eigenvectors, and a diagonal  $D$  collecting the corresponding eigenvalues, such that*

$$A = PDP^{-1}, \quad A^m = PD^mP^{-1} \quad (m \in \mathbb{N}),$$

with  $D^m = \text{diag}(\lambda_1^m, \dots, \lambda_n^m)$ .

## 5.9.2 Similarity as an Equivalence Relation

**Proposition 5.9.1.** *Define a relation  $\sim$  on  $\mathbb{C}^{n \times n}$  by*

$$A \sim B \iff \exists P \in \text{GL}_n(\mathbb{C}), A = PBP^{-1}.$$

*Then  $\sim$  is an equivalence relation on  $\mathbb{C}^{n \times n}$ .*

**Proof.** We check the three axioms.

*Reflexivity.* For any  $A \in \mathbb{C}^{n \times n}$ , we have  $A = IAI^{-1}$  with  $I = I_n \in \text{GL}_n(\mathbb{C})$ . Thus  $A \sim A$ .

*Symmetry.* Suppose  $A \sim B$ , i.e.,  $A = PBP^{-1}$  for some  $P \in \text{GL}_n(\mathbb{C})$ . Then  $B = P^{-1}AP$ , and  $P^{-1} \in \text{GL}_n(\mathbb{C})$ . Hence  $B \sim A$ .

*Transitivity.* Suppose  $A \sim B$  and  $B \sim C$ . Then  $A = PBP^{-1}$  and  $B = QCQ^{-1}$  for some  $P, Q \in \text{GL}_n(\mathbb{C})$ . Substituting, we get

$$A = P(QCQ^{-1})P^{-1} = (PQ)C(PQ)^{-1},$$

with  $PQ \in \text{GL}_n(\mathbb{C})$ . Hence  $A \sim C$ .

Therefore,  $\sim$  is reflexive, symmetric, and transitive, so it is an equivalence relation.  $\square$

## 5.9.3 The Quotient Set

The equivalence relation  $\sim$  partitions  $\mathbb{C}^{n \times n}$  into *similarity classes*. The quotient  $\mathbb{C}^{n \times n} / \sim$  is the set of these classes, where each class corresponds to all matrices representing the same linear operator (up to choice of basis). Describing this quotient amounts to classifying matrices up to similarity. Over  $\mathbb{C}$ , this classification is achieved by the Jordan normal form.

## 5.9.4 Jordan Matrices

**Definition 5.9.1** (Jordan block and Jordan matrix). *For  $\lambda \in \mathbb{C}$  and an integer  $k \geq 1$ , the Jordan block of size  $k$  with eigenvalue  $\lambda$  is the  $k \times k$  matrix*

$$J_k(\lambda) = \begin{pmatrix} \lambda & 1 & 0 & \cdots & 0 \\ 0 & \lambda & 1 & \cdots & 0 \\ 0 & 0 & \lambda & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 1 \\ 0 & 0 & \cdots & 0 & \lambda \end{pmatrix}.$$

*A Jordan matrix is a block diagonal matrix whose diagonal blocks are Jordan blocks.*

**Example 5.9.1.** A Jordan block of size 4 with eigenvalue  $\lambda = 3$  is

$$J_4(3) = \begin{pmatrix} 3 & 1 & 0 & 0 \\ 0 & 3 & 1 & 0 \\ 0 & 0 & 3 & 1 \\ 0 & 0 & 0 & 3 \end{pmatrix}.$$

A Jordan matrix of size 5 could for instance be

$$\begin{pmatrix} J_2(2) & 0 \\ 0 & J_3(5) \end{pmatrix}.$$

**Theorem 5.9.2** (Jordan Normal Form). Let  $A \in \mathbb{C}^{n \times n}$ . There exists an invertible matrix  $P \in \text{GL}_n(\mathbb{C})$  such that

$$P^{-1}AP = J,$$

where  $J$  is a block-diagonal Jordan matrix, i.e., a direct sum of Jordan blocks  $J_k(\lambda)$  with  $\lambda \in \mathbb{C}$  and  $k \geq 1$ . Moreover,  $J$  is unique up to the order of its Jordan blocks. Consequently, the similarity classes in  $\mathbb{C}^{n \times n}$  (i.e., the elements of  $\mathbb{C}^{n \times n} / \sim$ ) are in bijection with Jordan matrices of size  $n$ .

**Sketch of proof.** *Existence.*

1. **Primary decomposition.** Over  $\mathbb{C}$  the characteristic (hence the minimal) polynomial of  $A$  splits as a product of linear factors. Let  $\lambda_1, \dots, \lambda_s$  be the distinct eigenvalues and set  $N_i := (A - \lambda_i I)$ . One has the *primary decomposition*

$$\mathbb{C}^n = \ker N_1^{m_1} \oplus \dots \oplus \ker N_s^{m_s},$$

for suitable integers  $m_i \geq 1$ . Each summand  $V_{\lambda_i} := \ker N_i^{m_i}$  is  $A$ -stable and is called the *generalized eigenspace* of  $\lambda_i$ .

2. **Reduction on each primary component.** On  $V_{\lambda_i}$ , the operator  $A$  decomposes as  $A = \lambda_i I + N_i$ , where  $N_i$  is nilpotent ( $N_i^{m_i} = 0$ ). Thus it suffices to put the nilpotent map  $N_i$  into a canonical form.
3. **Jordan chains for nilpotent operators.** For a nilpotent operator  $N$  on a finite-dimensional space, one constructs *Jordan chains* (also called *Jordan strings*): sequences  $v, Nv, N^2v, \dots$  with  $N^\ell v \neq 0$  and  $N^{\ell+1}v = 0$ . Choosing a maximal family of such chains with disjoint spans yields a basis in which  $N$  is a direct sum of Jordan blocks  $J_k(0)$ .
4. **Conclusion.** Applying the previous step to each  $N_i$  on  $V_{\lambda_i}$  and then adding  $\lambda_i I$  to each block transforms  $N_i$ -blocks  $J_k(0)$  into  $J_k(\lambda_i)$ . Taking the union of these bases over all  $i$  gives a basis of  $\mathbb{C}^n$  in which  $A$  is block-diagonal with Jordan blocks  $J_k(\lambda_i)$ .

*Uniqueness (up to block order).* The sizes and counts of the Jordan blocks for a fixed eigenvalue  $\lambda$  are determined by the sequence of dimensions

$$d_k(\lambda) := \dim \ker ((A - \lambda I)^k) \quad (k \geq 1).$$

Indeed,  $d_k(\lambda) - d_{k-1}(\lambda)$  equals the number of Jordan blocks of size at least  $k$  for the eigenvalue  $\lambda$ . Hence the Jordan structure is uniquely determined by  $A$  (the order of blocks along the diagonal being irrelevant). This proves uniqueness up to permutation of blocks.  $\square$



## 5.10 Gram–Schmidt Orthogonalization

**Proposition 5.10.1** (Gram–Schmidt procedure). *Let  $(u_1, \dots, u_m)$  be a family of vectors in  $\mathbb{R}^n$  with  $m \leq n$ . The Gram–Schmidt procedure constructs from it an orthogonal (and then orthonormal) family  $(v_1, \dots, v_m)$  as follows:*

$$\begin{aligned} v_1 &= u_1, \\ v_2 &= u_2 - \frac{\langle u_2, v_1 \rangle}{\langle v_1, v_1 \rangle} v_1, \\ v_3 &= u_3 - \frac{\langle u_3, v_1 \rangle}{\langle v_1, v_1 \rangle} v_1 - \frac{\langle u_3, v_2 \rangle}{\langle v_2, v_2 \rangle} v_2, \\ &\vdots \\ v_k &= u_k - \sum_{j=1}^{k-1} \frac{\langle u_k, v_j \rangle}{\langle v_j, v_j \rangle} v_j, \quad (k = 2, \dots, m). \end{aligned}$$

The vectors  $(v_1, \dots, v_m)$  are mutually orthogonal, and whenever  $v_k \neq 0$  we may normalize

$$e_k = \frac{v_k}{\|v_k\|}$$

to obtain an orthonormal family  $(e_1, \dots, e_m)$ .

**Sketch of proof.** We prove by induction on  $k$  that: (i)  $v_1, \dots, v_k$  are mutually orthogonal, and (ii)  $\text{Span}\{v_1, \dots, v_k\} = \text{Span}\{u_1, \dots, u_k\}$ .

*Base case  $k = 1$ .* Trivial:  $v_1 = u_1$ , hence orthogonality is vacuous and the spans coincide.

*Inductive step.* Assume (i)–(ii) hold up to  $k - 1$ . By definition,

$$v_k = u_k - \sum_{j=1}^{k-1} \frac{\langle u_k, v_j \rangle}{\langle v_j, v_j \rangle} v_j.$$

For any  $1 \leq i \leq k - 1$ , take the inner product with  $v_i$ :

$$\langle v_k, v_i \rangle = \langle u_k, v_i \rangle - \sum_{j=1}^{k-1} \frac{\langle u_k, v_j \rangle}{\langle v_j, v_j \rangle} \langle v_j, v_i \rangle = \langle u_k, v_i \rangle - \frac{\langle u_k, v_i \rangle}{\langle v_i, v_i \rangle} \langle v_i, v_i \rangle = 0,$$

since  $\langle v_j, v_i \rangle = 0$  for  $j \neq i$  by the induction hypothesis. Hence  $v_k$  is orthogonal to each  $v_i$  ( $i < k$ ), proving (i) at rank  $k$ .

For (ii), the construction writes  $v_k$  as  $u_k$  minus a linear combination of  $v_1, \dots, v_{k-1}$ , so  $v_k \in \text{Span}\{u_1, \dots, u_k\}$ . Conversely, the defining formula solves for

$$u_k = v_k + \sum_{j=1}^{k-1} \frac{\langle u_k, v_j \rangle}{\langle v_j, v_j \rangle} v_j \in \text{Span}\{v_1, \dots, v_k\},$$

and by the induction hypothesis  $\text{Span}\{u_1, \dots, u_{k-1}\} = \text{Span}\{v_1, \dots, v_{k-1}\}$ , hence the spans coincide at step  $k$ .

Finally, whenever  $v_k \neq 0$ , normalizing  $e_k = v_k / \|v_k\|$  preserves orthogonality and yields an orthonormal family  $(e_1, \dots, e_m)$  spanning the same subspace.  $\square$

*Remark 21.* At each step,  $v_k$  is obtained from  $u_k$  by subtracting its projections onto the previously constructed vectors  $v_1, \dots, v_{k-1}$ . Geometrically, this forces  $v_k$  to lie orthogonally to the span of the earlier vectors. If the initial family  $(u_1, \dots, u_m)$  is linearly independent, then  $(v_1, \dots, v_m)$  is a basis, and  $(e_1, \dots, e_m)$  is an orthonormal basis of the same subspace.

## Chapter 6

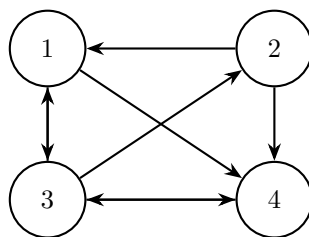
# Application of Diagonalization

Diagonalization is not only a theoretical tool but also a key technique in applied mathematics, data science, and computer science. In this chapter, we present two important applications: ranking in networks (PageRank) and data analysis via Principal Component Analysis (PCA). Both rely on eigenvalues and eigenvectors to extract fundamental information from matrices.

### 6.1 Ranking via Diagonalization: The Google Matrix

#### 6.1.1 Motivation

The World Wide Web can be modeled as a directed, weighted graph: vertices represent websites, and edges represent hyperlinks. If a website  $j$  has  $\ell_j$  outgoing links, each outgoing edge carries weight  $1/\ell_j$ . This construction yields a column-stochastic transition matrix  $T$  that describes the probability of moving from one page to another by following links. However, to allow for random jumps (so that the model does not get trapped in cycles), we add a "teleportation" mechanism, represented by a uniform matrix  $R$ .



#### 6.1.2 From Graph to Matrices

Here's a quick reminder on stochastic matrix:

**Definition 6.1.1** (Stochastic Matrix). A matrix  $M \in \mathbb{R}^{n \times n}$  is called **stochastic** if:

- all its entries are nonnegative:  $M_{ij} \geq 0$ ,
- the entries in each column (resp. row) sum to 1.

Convention: In this course, we adopt the **column-stochastic** convention:

$$\sum_{i=1}^n M_{ij} = 1 \quad \text{for every column } j.$$

Before going into the details, we need to define an important notion:

**Definition 6.1.2** (Ergodicity). A stochastic matrix (or a Markov chain) is called *ergodic* if it is irreducible and aperiodic.

- Irreducible means that every state can be reached from every other state (the Markov chain is fully connected).
- Aperiodic means that the system does not get trapped in cycles of fixed length: the return time to any state is not constrained to be a multiple of an integer greater than 1.

In this case, the chain admits a unique stationary distribution  $v_\infty$ , and the iteration

$$v(t+1) = Mv(t)$$

converges to  $v_\infty$  for any initial distribution  $v(0)$ .

The transition matrix of our worked example is defined as

$$T_{ij} = \begin{cases} \frac{1}{\ell_j} & \text{if there is an edge } j \rightarrow i, \\ 0 & \text{otherwise.} \end{cases}$$

leading to the following  $T$  for our application:

$$T = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{3} & 0 \\ 0 & 0 & \frac{1}{3} & 0 \\ \frac{1}{2} & 0 & 0 & 1 \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{3} & 0 \end{bmatrix}.$$

To guarantee ergodicity, we add

$$R = \frac{1}{n} \mathbf{1} \mathbf{1}^\top,$$

where  $\mathbf{1}$  is the all-ones column vector of size  $n$ . Thus  $R$  is the matrix where every entry is  $\frac{1}{n}$ , modeling the random jump to any website. The final Google matrix is

$$G = (1-p)T + pR, \quad \text{with a typical choice } p \approx 0.15.$$

The matrix  $G$  is stochastic, irreducible, and aperiodic, which ensures the existence of a unique stationary distribution.

### 6.1.3 Ranking Vector

Starting from the uniform distribution

$$v(0) = \left[ \frac{1}{n}, \quad \dots, \quad \frac{1}{n} \right]^\top,$$

we iterate

$$v(t+1) = Gv(t).$$

As  $t \rightarrow \infty$ , this converges to a unique fixed point  $v_\infty$ , satisfying

$$Gv_\infty = v_\infty.$$

The vector  $v_\infty$  is called the **PageRank vector**: its  $i$ -th entry gives the long-term probability of visiting page  $i$ . Pages with higher probability are ranked higher in search results.

### 6.1.4 Perron–Frobenius Theorem

**Theorem 6.1.1** (Perron–Frobenius). *If  $M$  is a column-stochastic matrix with all entries strictly positive, then:*

- 1 is an eigenvalue of  $M$ ,
- the associated eigenvector  $v_\infty$  has strictly positive entries,
- $v_\infty$  can be normalized so that its entries sum to 1,
- the sequence  $M^t v(0)$  converges to  $v_\infty$  for any initial vector  $v(0)$ .

*Remark 22.* The PageRank vector is precisely the eigenvector of  $G$  associated with eigenvalue 1. Direct computation is infeasible for the real web, where  $n$  is in the billions. Instead, the vector is approximated iteratively by successive multiplication by  $G$ , until convergence is reached.

## 6.2 Spectral properties of stochastic matrices and convergence rates

In this section we clarify where the spectrum of a real stochastic matrix lies, and why (under natural ergodicity assumptions) the *second* eigenvalue controls the rate at which the Markov iteration converges to stationarity.

### 6.2.1 Spectrum localization

**Definition 6.2.1** (Column-stochastic matrix). *A matrix  $M \in \mathbb{R}^{n \times n}$  is column-stochastic if*

$$\forall i, j, M_{ij} \geq 0$$

*and*

$$\sum_{i=1}^n M_{ij} = 1 \text{ for each column } j.$$

**Proposition 6.2.1** (Unit disk containment). *If  $M$  is column-stochastic, then 1 is an eigenvalue of  $M$ , and every eigenvalue  $\lambda$  of  $M$  satisfies  $|\lambda| \leq 1$ . Equivalently, the spectrum is contained in the closed unit disk:*

$$\text{Sp}(M) \subseteq \{z \in \mathbb{C} : |z| \leq 1\}.$$

**Proof.** Since the columns of  $M$  sum to 1, we have  $\mathbf{1}^\top M = \mathbf{1}^\top$ , so 1 is a *left* eigenvalue with eigenvector  $\mathbf{1}$ . For the bound  $|\lambda| \leq 1$ , apply Gershgorin’s circle theorem to  $M^\top$ , which is *row*-stochastic: each row of  $M^\top$  has nonnegative entries summing to 1. Thus every Gershgorin disk of  $M^\top$  is centered at a number in  $[0, 1]$  with radius at most 1; in particular all disks lie in the closed unit disk, hence so does the spectrum of  $M^\top$  and therefore of  $M$ .  $\square$

*Remark 23* (About real vs. complex eigenvalues). In general, a real stochastic matrix *need not* have only real eigenvalues; complex conjugate pairs inside the unit disk may occur (e.g. certain periodic chains). Real-line containment such as  $\text{Sp}(M) \subset [0, 1]$  *fails* in general. Two important cases do give real spectra:

- If  $M$  is *reversible* w.r.t. a strictly positive stationary distribution  $\pi$  (i.e.  $\pi_i M_{ij} = \pi_j M_{ji}$ ), then  $S := D_\pi^{1/2} M D_\pi^{-1/2}$  is *symmetric*, hence diagonalizable with real eigenvalues in  $[-1, 1]$ ;  $M$  is similar to  $S$ , so  $M$  has the same real spectrum.
- If  $M$  is *symmetric* (hence doubly stochastic), then it is self-adjoint in the Euclidean inner product and has real spectrum contained in  $[-1, 1]$ .

**Proposition 6.2.2** (Perron–Frobenius refinement for primitive chains). *If  $M$  is primitive (irreducible and aperiodic; e.g.  $M > 0$  entrywise), then*

$$1 = \rho(M) \text{ is a simple eigenvalue,} \quad |\lambda| < 1 \quad \text{for all other eigenvalues } \lambda \in \text{Sp}(M) \setminus \{1\}.$$

**Sketch.** This is the classical Perron–Frobenius theorem for nonnegative irreducible matrices with aperiodicity: primitivity implies the Perron root is 1, simple, with strictly positive right and left eigenvectors; all other eigenvalues have strictly smaller modulus.  $\square$

## 6.2.2 Why the second eigenvalue controls convergence

Let  $v(t) = M^t v(0)$  denote the state at time  $t$  in the column-stochastic convention. When  $M$  is primitive,  $v(t) \rightarrow v_\infty$  (the unique stationary state with  $M v_\infty = v_\infty$ ,  $\mathbf{1}^\top v_\infty = 1$ ) for any initial probability vector  $v(0)$ .

**Theorem 6.2.1** (Geometric convergence at the spectral-gap rate). *Assume  $M$  is primitive column-stochastic. Let the eigenvalues of  $M$  be  $1 = \lambda_1, \lambda_2, \dots, \lambda_n$  ordered so that  $|\lambda_2| \geq \dots \geq |\lambda_n|$ . Then there exists a constant  $C(v(0)) < \infty$  such that*

$$\|M^t v(0) - v_\infty\| \leq C(v(0)) |\lambda_2|^t \quad \text{for all } t \in \mathbb{N},$$

for any norm  $\|\cdot\|$  compatible with matrix multiplication. In particular, the mixing time to reach accuracy  $\varepsilon$  is

$$t_{\text{mix}}(\varepsilon) \lesssim \frac{\log(1/\varepsilon)}{-\log |\lambda_2|} \sim \frac{1}{1 - |\lambda_2|} \log \frac{1}{\varepsilon} \quad \text{when } |\lambda_2| \uparrow 1.$$

**Accessible proof sketches.** (*Diagonalizable / reversible case*). If  $M$  is reversible, then  $S = D_\pi^{1/2} M D_\pi^{-1/2}$  is symmetric with eigenvalues  $1 = \lambda_1 > \lambda_2 \geq \dots \geq \lambda_n > -1$  and an orthonormal eigenbasis. Decompose the initial error in the  $L^2(\pi)$ -orthonormal eigenbasis:

$$D_\pi^{1/2}(v(t) - v_\infty) = S^t D_\pi^{1/2}(v(0) - v_\infty) = \sum_{k=2}^n \lambda_k^t \alpha_k u_k,$$

so  $\|v(t) - v_\infty\|_{L^2(\pi)} \leq (\max_{k \geq 2} |\lambda_k|)^t \|v(0) - v_\infty\|_{L^2(\pi)}$ .

(*General primitive case*). By Jordan or Schur decomposition,

$$M = v_\infty \mathbf{1}^\top + R, \quad \mathbf{1}^\top R = 0, \quad R v_\infty = 0, \quad \rho(R) = |\lambda_2| < 1.$$

Then  $M^t = v_\infty \mathbf{1}^\top + R^t$ , hence  $M^t v(0) - v_\infty = R^t(v(0) - v_\infty)$  and  $\|M^t v(0) - v_\infty\| \leq C \|R^t\| \|v(0) - v_\infty\|$ . For any submultiplicative norm,  $\|R^t\| \leq C' |\lambda_2|^t$  (polynomial factors from Jordan blocks can appear but are dominated asymptotically by  $|\lambda_2|^t$ ). This yields the stated bound.  $\square$

*Remark 24* (Intuition: the “slowest mode”). The stochastic iteration splits into a stationary component  $v_\infty$  and transient “modes” corresponding to the other eigenvalues. Each mode decays like  $|\lambda_k|^t$ ; the slowest is determined by the largest modulus among  $k \geq 2$ , namely  $|\lambda_2|$ . A smaller  $|\lambda_2|$  (i.e. a larger spectral gap  $1 - |\lambda_2|$ ) means faster forgetting of initial conditions and quicker convergence to equilibrium.

## 6.2.3 Quantifying the “speed”: inverse gap and log factor

Two equivalent parametrizations of the speed are commonly used:

- **Geometric ratio:** the error decays roughly like  $|\lambda_2|^t$ .
- **Spectral gap:**  $1 - |\lambda_2|$ ; for small gaps,  $-\log |\lambda_2| \sim 1 - |\lambda_2|$  and

$$t_{\text{mix}}(\varepsilon) \approx \frac{1}{1 - |\lambda_2|} \log \frac{1}{\varepsilon}.$$

Thus, the “inverse gap”  $1/(1 - |\lambda_2|)$  sets the natural time scale of convergence: the closer  $|\lambda_2|$  is to 1, the slower the mixing.

### 6.2.4 Aperiodicity, periodicity, and complex eigenvalues

If the chain is irreducible but *periodic*, the spectrum still lies in the unit disk, but there can be eigenvalues on the unit circle other than 1 (e.g.  $-1$ ), leading to persistent oscillations. Aperiodicity rules out such extra unit-modulus eigenvalues, ensuring that  $|\lambda_2| < 1$  and that the geometric decay dominates without oscillatory persistence.

### 6.2.5 Practical takeaway

- *Always true (column-stochastic)*:  $\text{Sp}(M) \subseteq \{z : |z| \leq 1\}$  and  $1 \in \text{Sp}(M)$ .
- *Primitive*: 1 simple, all other eigenvalues strictly inside the unit disk  $\Rightarrow$  geometric convergence governed by  $|\lambda_2|$ .
- *Reversible (or symmetric)*: spectrum real in  $[-1, 1]$ , orthogonal eigendecomposition, sharp  $L^2(\pi)$  bounds.

These facts explain both *why* the iteration converges and *how fast* it does so, in terms of the second eigenvalue and its spectral gap from 1.

## 6.3 Data Analysis via Diagonalization: Covariance and PCA

### 6.3.1 The Covariance Matrix

Suppose we have  $k$  observations of  $m$  variables:

$$X = \{p_1, \dots, p_k\}, \quad p_i \in \mathbb{R}^m.$$

For each coordinate  $j$ , let  $\mu_j$  denote the mean of the  $j$ -th component across all observations. We define the centered data matrix

$$N_{ij} = p_{ij} - \mu_j.$$

The covariance matrix of the dataset is then

$$\text{cov}(X) = N^\top N.$$

### 6.3.2 Variance and Eigenanalysis

The eigenvalues of the covariance matrix capture the variance of the data along specific directions, and the eigenvectors give those directions.

**Theorem 6.3.1** (Variance along Eigenvectors). *Let  $\lambda_1 \leq \dots \leq \lambda_m$  be the eigenvalues of  $\text{cov}(X)$ . Then the variance of the data in the direction of the eigenvector associated with  $\lambda_i$  is proportional to  $\lambda_i$ .*

*Remark 25.* Some important remarks:

- The largest eigenvalue indicates the direction of maximum variance.
- Smaller eigenvalues correspond to directions with less spread.
- In two dimensions, the eigenvectors define the axes of the ellipse approximating the data cloud, and the eigenvalues determine the lengths of those axes.

### 6.3.3 Principal Component Analysis (PCA)

Theorem 6.3.1 is the foundation of **Principal Component Analysis (PCA)**. PCA reduces dimensionality by projecting data onto the first few principal directions, corresponding to the largest eigenvalues. This allows us to capture most of the variability in the data with fewer dimensions, a technique widely used in statistics, signal processing, and machine learning.

### 6.3.4 Worked Example

**Example 6.3.1.** *Consider the dataset*

$$X = \{(1, 1), (2, 2), (2, 3), (3, 2), (3, 3), (4, 4)\}.$$

1. *The coordinate-wise mean is  $\mu = (2.5, 2.5)$ .*
2. *The centered data matrix  $N$  has rows given by  $(p_i - \mu)$ .*
3. *The covariance matrix is*
$$\text{cov}(X) = N^\top N.$$
4. *Computing its eigenvalues and eigenvectors reveals the directions of maximal and minimal variance in the dataset.*

*This illustrates concretely how diagonalization of the covariance matrix leads to the principal axes of variation.*

## Chapter 7

# Application to Statistics: Least Squares and SVD

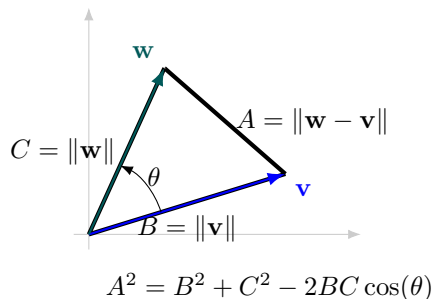
In this chapter we illustrate how diagonalization and orthogonality underpin fundamental techniques in statistics. We begin with the geometric intuition of orthogonality and projections, then study the least squares method, before linking it to the Singular Value Decomposition (SVD).

### 7.1 Orthogonality and Distance Minimization

Orthogonality is central to approximation problems in data science and statistics. Whenever we fit a model to data, we want to minimize the distance between the model's predictions and the actual data points.

Because squared distances expand into quadratic expressions, minimization naturally leads to linear algebra. Vector calculus reveals that minimizing a squared error is equivalent to computing an *orthogonal projection* of a vector onto a subspace.

#### 7.1.1 Law of Cosines and the Dot Product



The geometric relation between distance and orthogonality is already encoded in the law of cosines. For a triangle with sides  $A, B, C$  and opposite angle  $c$ ,

$$A^2 = B^2 + C^2 - 2BC \cos(c).$$

Applying this to the triangle formed by two vectors  $\mathbf{v}, \mathbf{w}$  at the origin, we obtain the fundamental formula:

$$\mathbf{v} \cdot \mathbf{w} = \|\mathbf{v}\| \|\mathbf{w}\| \cos(c),$$

where  $c$  is the angle between  $\mathbf{v}$  and  $\mathbf{w}$ . Thus, the dot product measures both orthogonality and similarity between vectors.



## 7.2 Least Squares Approximation

### 7.2.1 Motivation

In statistical modeling, we seek to approximate data without overfitting. The *least squares principle* states: among all possible models, choose the one that minimizes the sum of squared deviations between data and model predictions.

Algebraically, least squares is equivalent to projecting the data vector onto the column space of the design matrix.

### 7.2.2 Warm-up Example: Line of Best Fit in $\mathbb{R}^2$

Consider the dataset

$$X = \{(1, 6), (2, 5), (3, 7), (4, 10)\}.$$

We want to fit the line  $y = ax + b$  that minimizes

$$\text{error}^2 = \sum_{(x_i, y_i) \in X} (y_i - (ax_i + b))^2.$$

This leads to the linear system

$$\begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix} \begin{bmatrix} b \\ a \end{bmatrix} = \begin{bmatrix} 6 \\ 5 \\ 7 \\ 10 \end{bmatrix}.$$

Geometric interpretation: the right-hand side vector is projected onto the column space of the design matrix.

## 7.3 Ordinary Least Squares (OLS)

**Setup** We consider the linear regression model

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \quad i = 1, \dots, N,$$

where  $\varepsilon_i$  are independent random errors with mean zero. The goal is to estimate  $(\beta_0, \beta_1)$  by minimizing the sum of squared errors.

**Derivation** We minimize

$$f(\beta_0, \beta_1) = \sum_{i=1}^N (Y_i - (\beta_0 + \beta_1 X_i))^2.$$

Setting the derivatives to zero gives the normal equations:

$$\begin{cases} \sum (Y_i - \beta_0 - \beta_1 X_i) = 0, \\ \sum X_i (Y_i - \beta_0 - \beta_1 X_i) = 0. \end{cases}$$

Introducing the sample means

$$\bar{X} = \frac{1}{N} \sum X_i, \quad \bar{Y} = \frac{1}{N} \sum Y_i,$$

the equations simplify to

$$\begin{cases} \bar{Y} = \beta_0 + \beta_1 \bar{X}, \\ \frac{1}{N} \sum X_i Y_i = \beta_0 \bar{X} + \beta_1 \frac{1}{N} \sum X_i^2. \end{cases}$$

**Closed-Form Solution** We define the variance and covariance:

$$\text{Var}(X) = \frac{1}{N} \sum (X_i - \bar{X})^2, \quad \text{Cov}(X, Y) = \frac{1}{N} \sum (X_i - \bar{X})(Y_i - \bar{Y}).$$

Then the optimal coefficients are

$$\hat{\beta}_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}, \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}.$$

**Interpretation** The slope  $\hat{\beta}_1$  measures the average change in  $Y$  per unit change in  $X$ , normalized by the variability of  $X$ . The intercept  $\hat{\beta}_0$  adjusts the line so that it passes through the point  $(\bar{X}, \bar{Y})$ .

## 7.4 Connection to Singular Value Decomposition (SVD)

The least squares solution can also be expressed using the SVD. If the design matrix  $A$  admits an SVD

$$A = U\Sigma V^\top,$$

then the least squares solution to  $A\beta \approx Y$  is

$$\hat{\beta} = A^+ Y,$$

where  $A^+ = V\Sigma^+ U^\top$  is the Moore–Penrose pseudoinverse.

**Definition 7.4.1** (Moore–Penrose Pseudoinverse). *Let  $M \in \mathbb{R}^{m \times n}$ . The Moore–Penrose pseudoinverse of  $M$ , denoted  $M^+ \in \mathbb{R}^{n \times m}$ , is the unique matrix satisfying the following four conditions:*

1.  $MM^+M = M$ ,
2.  $M^+MM^+ = M^+$ ,
3.  $(MM^+)^\top = MM^+$ ,
4.  $(M^+M)^\top = M^+M$ .

*Remark 26.* The pseudoinverse extends the concept of the matrix inverse to non-square or singular matrices.

- If  $M$  is square and invertible, then  $M^+ = M^{-1}$ .
- If  $M^\top M$  is invertible (which happens exactly when  $M$  has full column rank), then

$$M^+ = (M^\top M)^{-1} M^\top.$$

Symmetrically, if  $MM^\top$  is invertible (full row rank), then

$$M^+ = M^\top (MM^\top)^{-1}.$$

- In the least-squares problem  $\min_x \|Mx - b\|_2$ , the solution of minimal Euclidean norm is given by  $x^* = M^+ b$ .
- Using the singular value decomposition (SVD)  $M = U\Sigma V^\top$ , the pseudoinverse can always be computed as

$$M^+ = V\Sigma^+ U^\top,$$

where  $\Sigma^+$  is obtained by inverting the nonzero singular values of  $\Sigma$  and transposing the matrix.

Thus, the pseudoinverse provides a systematic way to “invert” linear systems even when  $M$  is not square or not of full rank.

This formulation is particularly useful when  $A$  is not full rank, or when we need numerical stability in large-scale problems. It shows how least squares, orthogonal projections, and diagonalization come together under the umbrella of the SVD.

## 7.5 Subspaces and Orthogonality

We can then deepen the link between orthogonality, projections, least squares, and the Singular Value Decomposition (SVD).

### 7.5.1 Definitions

Two subspaces  $W, W' \subset V$  are said to be **orthogonal** if

$$\mathbf{w} \cdot \mathbf{w}' = 0 \quad \forall \mathbf{w} \in W, \mathbf{w}' \in W'.$$

A set  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is **orthonormal** if

$$\mathbf{v}_i \cdot \mathbf{v}_j = \delta_{ij}.$$

A matrix  $A$  is **orthogonal** if its columns form an orthonormal set.

### 7.5.2 Fundamental Subspaces of a Matrix

Let  $A \in \mathbb{R}^{m \times n}$ . Four subspaces naturally arise:

- The **column space**:

$$C(A) = \text{Im}(A) = \{A\mathbf{x} \mid \mathbf{x} \in \mathbb{R}^n\} \subset \mathbb{R}^m,$$

of dimension equal to the rank of  $A$ .

- The **null space**:

$$N(A) = \ker(A) = \{\mathbf{x} \in \mathbb{R}^n \mid A\mathbf{x} = 0\} \subset \mathbb{R}^n.$$

- The **row space**:

$$R(A) = \text{Im}(A^\top) = \text{span of the row vectors of } A \subset \mathbb{R}^n.$$

- The **orthogonal complement**: for a subspace  $W \subset V$ ,

$$W^\perp = \{\mathbf{v} \in V \mid \mathbf{v} \cdot \mathbf{w} = 0, \forall \mathbf{w} \in W\}.$$

The fundamental relation between these spaces is given by the **rank–nullity theorem**:

$$n = \dim N(A) + \dim R(A).$$

### 7.5.3 Worked Example

Take

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

Computation shows:

$$C(A) = \text{Span}\{(1, 2, 1)^\top, (2, 4, 1)^\top\}, \quad \dim = 2,$$

$$C(A)^\perp = \ker(A^\top) = \text{Span}\{(-2, 1, 0)^\top\},$$

$$N(A) = \ker(A) = \text{Span}\{(1, -2, 1)^\top\}, \quad \dim = 1,$$

$$R(A) = \text{Span}\{(1, 2, 3)^\top, (1, 1, 1)^\top\}, \quad \dim = 2.$$

Hence the rank–nullity theorem holds:

$$3 = \dim N(A) + \dim R(A) = 1 + 2.$$

*Exercise 7.5.1.* Show that  $N(A) = R(A)^\perp$  and  $N(A^\top) = C(A)^\perp$ . Prove that any  $\mathbf{v} \in V$  decomposes uniquely as

$$\mathbf{v} = \mathbf{w} + \mathbf{w}^\perp, \quad \mathbf{w} \in W, \mathbf{w}^\perp \in W^\perp,$$

and that  $\mathbf{w}$  is the closest vector in  $W$  to  $\mathbf{v}$ .

## 7.6 Least Squares and Projection Matrices

Consider again the design matrix from the line fitting problem:

$$A = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 6 \\ 5 \\ 7 \\ 10 \end{bmatrix}.$$

The normal equations are

$$A^\top A \mathbf{x} = A^\top \mathbf{b}.$$

Here

$$A^\top A = \begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix}, \quad (A^\top A)^{-1} = \begin{bmatrix} 3/2 & -1/2 \\ -1/2 & 1/5 \end{bmatrix}.$$

The projection matrix onto the column space of  $A$  is

$$P = A(A^\top A)^{-1}A^\top = \frac{1}{10} \begin{bmatrix} 7 & 4 & 1 & -2 \\ 4 & 3 & 2 & 1 \\ 1 & 2 & 3 & 4 \\ -2 & 1 & 4 & 7 \end{bmatrix}.$$

Hence

$$P \cdot \mathbf{b} = \begin{bmatrix} 4.9 \\ 6.3 \\ 7.7 \\ 9.1 \end{bmatrix},$$

which lies in  $\text{Col}(A)$  and corresponds to the fitted values.

### 7.6.1 Quadratic Least Squares

For a quadratic model  $y = ax^2 + bx + c$ , the design matrix is

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \\ 1 & 4 & 16 \end{bmatrix},$$

and the normal equations  $A^\top A \mathbf{y} = A^\top \mathbf{b}$  again yield the projection of  $\mathbf{b}$  onto  $\text{Col}(A)$ .

## 7.7 Singular Value Decomposition (SVD)

### 7.7.1 From $M^\top M$ to the SVD: symmetry, (semi)definiteness, and eigenvalues

**Proposition 7.7.1** (Basic properties of  $S = M^\top M$ ). *Let  $M \in \mathbb{R}^{m \times n}$  and set  $S := M^\top M \in \mathbb{R}^{n \times n}$ . Then*

1.  *$S$  is symmetric:  $(M^\top M)^\top = M^\top M$ .*
2.  *$S$  is positive semidefinite:  $x^\top S x \geq 0$  for all  $x \in \mathbb{R}^n$ .*
3.  *$\ker S = \ker M$ . In particular,  $S$  is positive definite iff  $M$  has full column rank (i.e.  $\ker M = \{0\}$ ).*

**Proof.** (1) Symmetry is immediate:  $(M^\top M)^\top = M^\top (M^\top)^\top = M^\top M$ .

(2) For any  $x \in \mathbb{R}^n$ ,

$$x^\top Sx = x^\top M^\top Mx = (Mx)^\top (Mx) = \sum_{i=1}^m ((Mx)_i)^2 \geq 0,$$

a sum of squares.

(3) If  $x \in \ker M$ , then  $Mx = 0$  and  $Sx = M^\top (Mx) = 0$ , so  $x \in \ker S$ . Conversely, if  $x \in \ker S$ , then

$$0 = x^\top Sx = (Mx)^\top (Mx) = \sum_{i=1}^m ((Mx)_i)^2,$$

hence  $Mx = 0$ . Therefore  $\ker S = \ker M$ . It follows that  $S$  is positive *definite* (i.e.  $x^\top Sx > 0$  for all  $x \neq 0$ ) exactly when  $\ker M = \{0\}$ , equivalently when the columns of  $M$  are linearly independent.  $\square$

**Proposition 7.7.2** (Eigenvalues of  $S$ ). *Let  $S = M^\top M$  as above. Then every eigenvalue of  $S$  is nonnegative. Moreover, if  $S$  is positive definite (equivalently,  $M$  has full column rank), then every eigenvalue is strictly positive.*

**Proof.** Let  $(\lambda, v)$  be an eigenpair with  $v \neq 0$ , so  $Sv = \lambda v$ . Taking the scalar product with  $v$ ,

$$\lambda \|v\|^2 = v^\top Sv = v^\top M^\top Mv = (Mv)^\top (Mv) = \sum_{i=1}^m ((Mv)_i)^2 \geq 0.$$

Since  $\|v\|^2 > 0$ , we get  $\lambda \geq 0$ . If  $S$  is positive definite, then  $v^\top Sv > 0$  for every nonzero  $v$ , hence  $\lambda > 0$  for every eigenpair.  $\square$

*Remark 27* (Link to the SVD). By the spectral theorem,  $S = M^\top M$  admits an orthonormal eigenbasis with eigenvalues  $\lambda_i \geq 0$ . Writing  $\sigma_i := \sqrt{\lambda_i}$  (the nonnegative square roots) yields the singular values of  $M$ . This is the starting point of the SVD construction: from  $M^\top M v_i = \sigma_i^2 v_i$ , define  $u_i := \frac{1}{\sigma_i} Mv_i$  (when  $\sigma_i > 0$ ), which are orthonormal, and assemble  $U, \Sigma, V$  so that  $M = U\Sigma V^\top$ .

## 7.7.2 SVD Statement

**Theorem 7.7.1** (Singular Value Decomposition (SVD)). *Let  $M \in \mathbb{R}^{m \times n}$ . There exist orthogonal matrices  $U \in \mathbb{R}^{m \times m}$  and  $V \in \mathbb{R}^{n \times n}$ , and a diagonal (“rectangular diagonal”) matrix  $\Sigma \in \mathbb{R}^{m \times n}$  with nonnegative entries on the diagonal,*

$$\Sigma = \begin{bmatrix} \text{diag}(\sigma_1, \dots, \sigma_r) & 0 \\ 0 & 0 \end{bmatrix}, \quad \sigma_1 \geq \dots \geq \sigma_r > 0,$$

*such that*

$$M = U\Sigma V^\top.$$

*The numbers  $\sigma_i$  are the singular values of  $M$ . They are uniquely determined by  $M$  (up to ordering). The columns  $u_i$  of  $U$  (resp.  $v_i$  of  $V$ ) are called left (resp. right) singular vectors; they are unique up to signs and orthogonal rotations inside eigenspaces corresponding to equal singular values.*

*Remark 28.* Key Ideas of the Proof

- $M^\top M$  is symmetric and diagonalizable with orthonormal eigenvectors.
- If  $M^\top M v_i = \lambda_i v_i$ , define  $\sigma_i = \sqrt{\lambda_i}$ .
- Set  $u_i = \frac{1}{\sigma_i} Mv_i$ , which are orthonormal in  $\mathbb{R}^m$ .
- Collecting  $u_i$  into  $U$  and  $v_i$  into  $V$  yields  $M = U\Sigma V^\top$ .

**Proof.** *Step 1: Spectral theorem for  $M^\top M$ .* Consider the symmetric positive semidefinite matrix  $B := M^\top M \in \mathbb{R}^{n \times n}$ . By the spectral theorem, there exists an orthonormal basis of eigenvectors  $\{v_1, \dots, v_n\}$  of  $\mathbb{R}^n$  and real eigenvalues  $\lambda_1, \dots, \lambda_n \geq 0$  such that

$$Bv_i = \lambda_i v_i, \quad v_i^\top v_j = \delta_{ij}.$$

Reorder so that  $\lambda_1 \geq \dots \geq \lambda_r > 0$  and  $\lambda_{r+1} = \dots = \lambda_n = 0$ , with  $r = \text{rank}(M)$  (see below).

Define  $\sigma_i := \sqrt{\lambda_i}$  for  $1 \leq i \leq r$ . These are nonzero and will be the (positive) singular values.

*Step 2: Define the left singular vectors.* For  $i = 1, \dots, r$ , set

$$u_i := \frac{1}{\sigma_i} M v_i \in \mathbb{R}^m.$$

Then  $u_i \neq 0$  (since  $\sigma_i > 0$ ) and the family  $\{u_i\}_{i=1}^r$  is orthonormal: for  $i, j \leq r$ ,

$$u_i^\top u_j = \frac{1}{\sigma_i \sigma_j} v_i^\top M^\top M v_j = \frac{1}{\sigma_i \sigma_j} v_i^\top (\lambda_j v_j) = \frac{\lambda_j}{\sigma_i \sigma_j} v_i^\top v_j = \frac{\sigma_j^2}{\sigma_i \sigma_j} \delta_{ij} = \delta_{ij}.$$

*Step 3: Complete to orthonormal bases.* Extend  $\{u_1, \dots, u_r\}$  to an orthonormal basis  $\{u_1, \dots, u_r, u_{r+1}, \dots, u_m\}$  of  $\mathbb{R}^m$  (e.g. by Gram–Schmidt). Similarly, keep  $\{v_1, \dots, v_n\}$  as the orthonormal basis from Step 1.

Define  $U := [u_1 \ \dots \ u_m] \in \mathbb{R}^{m \times m}$  and  $V := [v_1 \ \dots \ v_n] \in \mathbb{R}^{n \times n}$ ; then  $U$  and  $V$  are orthogonal.

*Step 4: The diagonal form  $U^\top M V$ .* Compute the action of  $M$  on the basis vectors  $v_i$ :

$$M v_i = \begin{cases} \sigma_i u_i, & 1 \leq i \leq r, \\ 0, & r < i \leq n, \end{cases}$$

because for  $i > r$ ,  $\lambda_i = 0$  implies  $M^\top M v_i = 0$  and hence  $M v_i = 0$ , since  $\|M v_i\|^2 = v_i^\top (M^\top M) v_i = 0$ . Therefore, in the bases given by  $U$  and  $V$ ,

$$U^\top M V = \begin{bmatrix} \text{diag}(\sigma_1, \dots, \sigma_r) & 0 \\ 0 & 0 \end{bmatrix} =: \Sigma.$$

Indeed, if one multiplies  $M v_i$  at left by  $u_j^\top$ , one gets  $\sigma_i \delta_{ij}$ . Equivalently,  $M = U \Sigma V^\top$ .

*Step 5: Rank identity and positivity.* The number  $r$  of positive singular values equals  $\text{rank}(A)$ . Indeed,  $\text{Im}(M) = \text{Span}\{M v_1, \dots, M v_n\} = \text{Span}\{u_1, \dots, u_r\}$ , so  $\dim \text{Im}(M) = r$ .

This proves existence.

*Uniqueness of singular values.* The multiset  $\{\sigma_i^2\}$  is the spectrum of  $M^\top M$ , hence uniquely determined by  $M$  (up to ordering). Taking nonnegative square roots yields unique  $\sigma_i$ 's (up to ordering).

*Uniqueness of singular vectors up to symmetries.* If  $\sigma_i$  is simple (strictly larger than neighboring singular values), then  $v_i$  and  $u_i$  are unique up to a simultaneous sign flip ( $v_i \mapsto -v_i$ ,  $u_i \mapsto -u_i$ ). If a singular value has multiplicity  $k > 1$ , any orthogonal change of basis within the corresponding  $k$ -dimensional right (resp. left) singular subspace yields another valid set of  $v_i$ 's (resp.  $u_i$ 's), producing the same  $\Sigma$ .

This completes the proof. □

*Remark 29 (Complex case).* For  $M \in \mathbb{C}^{m \times n}$ , replace transposes with conjugate transposes. The proof is identical, applying the spectral theorem to the Hermitian positive semidefinite matrix  $M^* M$ .

*Remark 30 (Characterizations).* The singular values admit the variational characterization (Courant–Fischer/Ky Fan)

$$\sigma_k = \max_{\substack{S \subset \mathbb{R}^n \\ \dim S = k}} \min_{\substack{x \in S \\ \|x\|=1}} \|M x\| = \min_{\substack{T \subset \mathbb{R}^n \\ \dim T = n-k+1}} \max_{\substack{x \perp T \\ \|x\|=1}} \|M x\|,$$

and satisfy  $\|M\|_2 = \sigma_1$ ,  $\|M\|_F^2 = \sum_i \sigma_i^2$ .

**Interpretation** SVD expresses any matrix as

(orthogonal change of basis)  $\times$  (scaling)  $\times$  (orthogonal change of basis).

It generalizes diagonalization to rectangular matrices.

**Proposition 7.7.3** (Relations between Singular Vectors). *Let  $M \in \mathbb{R}^{m \times n}$  and consider its Singular Value Decomposition*

$$M = U \Sigma V^\top,$$

*where  $U \in \mathbb{R}^{m \times m}$  and  $V \in \mathbb{R}^{n \times n}$  are **orthogonal**, and  $\Sigma \in \mathbb{R}^{m \times n}$  is block-diagonal with singular values  $\sigma_1 \geq \dots \geq \sigma_r > 0$  on the diagonal ( $r = \text{rank}(M)$ ).*

*Then, for each  $i = 1, \dots, r$ :*

$$M v_i = \sigma_i u_i \quad \text{and} \quad M^\top u_i = \sigma_i v_i,$$

*where  $u_i$  and  $v_i$  denote the  $i$ -th columns of  $U$  and  $V$ , respectively.*

**Proof.** By construction (check SVD proof). □

*Remark 31.* Practical Formulas:

- If  $v_i$  is known, the corresponding  $u_i$  is obtained as  $u_i = \frac{1}{\sigma_i} M v_i$ .
- If  $u_i$  is known, the corresponding  $v_i$  is obtained as  $v_i = \frac{1}{\sigma_i} M^\top u_i$ .

These formulas are valid whenever  $\sigma_i > 0$ .

*Remark 32* (Gram-matrix viewpoint). If  $M = [c_1 \dots c_n]$  with columns  $c_j \in \mathbb{R}^m$ , then

$$M^\top M = (\langle c_i, c_j \rangle)_{1 \leq i, j \leq n}$$

is the Gram matrix of  $\{c_j\}$ . For any coefficients  $\alpha \in \mathbb{R}^n$ ,

$$\alpha^\top (M^\top M) \alpha = \left\langle \sum_j \alpha_j c_j, \sum_k \alpha_k c_k \right\rangle = \langle v, v \rangle \geq 0,$$

with  $v = \sum_j \alpha_j c_j$ , showing again that  $M^\top M$  is positive semidefinite and positive definite exactly when the  $c_j$  are independent.

**Theorem 7.7.2** (Eckart–Young–Mirsky). *Let  $M \in \mathbb{R}^{m \times n}$  have singular values  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  ( $r = \text{rank}(M)$ ) and SVD  $M = U \Sigma V^\top$ . For any  $1 \leq k < r$ , the truncated SVD*

$$M_k := \sum_{i=1}^k \sigma_i u_i v_i^\top$$

*is a best rank- $k$  approximation of  $M$  simultaneously for the spectral norm and the Frobenius norm:*

$$\|M - M_k\|_2 = \min_{\text{rank}(X) \leq k} \|M - X\|_2 = \sigma_{k+1}, \quad \|M - M_k\|_F = \min_{\text{rank}(X) \leq k} \|M - X\|_F = \left( \sum_{i=k+1}^r \sigma_i^2 \right)^{1/2}.$$

*Moreover, every minimizer has the form  $U \begin{bmatrix} \Sigma_k & 0 \\ 0 & 0 \end{bmatrix} V^\top$  with  $\Sigma_k = \text{diag}(\sigma_1, \dots, \sigma_k)$  (up to rotations within equal singular values).*

**Proof sketch.** Both  $\|\cdot\|_2$  and  $\|\cdot\|_F$  are unitarily invariant. Thus  $\|M - X\| = \|U^\top(M - X)V\| = \|\Sigma - Y\|$  for  $Y := U^\top X V$  with  $\text{rank}(Y) \leq k$ . Hence the problem reduces to approximating the diagonal nonnegative matrix  $\Sigma$  by a rank- $k$  matrix  $Y$ .

*Spectral norm.* Since  $\|\cdot\|_2$  equals the largest singular value, and  $\Sigma - Y$  has singular values bounded below by the tail of those of  $\Sigma$ , the minimum achievable  $\|\Sigma - Y\|_2$  is  $\sigma_{k+1}$ , attained by taking  $Y = \text{diag}(\sigma_1, \dots, \sigma_k, 0, \dots)$ .

*Frobenius norm.* Because  $\|\cdot\|_F^2$  is the sum of squared singular values (and equals the sum of squared entries for diagonal matrices), the best choice is again to keep the first  $k$  diagonal entries and zero the rest, yielding error  $\sum_{i>k} \sigma_i^2$ . Formal justifications use von Neumann's trace inequality and variational characterizations (Ky Fan/Courant–Fischer) to certify optimality.  $\square$

*Remark 33* (Interpretation (explained variance and principal directions)). In Frobenius norm, the *explained variance* by the rank- $k$  approximation is

$$\frac{\|M_k\|_F^2}{\|M\|_F^2} = \frac{\sum_{i=1}^k \sigma_i^2}{\sum_{i=1}^r \sigma_i^2}.$$

Thus  $\sigma_1^2$  (and  $u_1, v_1$ ) captures the largest possible share of total energy/variance among all rank-one approximations. The spectral norm error equals  $\sigma_{k+1}$ , so the first singular value (not the spectral radius of  $M$  in general) governs the maximal stretch of  $M$  and the quality of the best rank-one fit. This is precisely the linear-algebraic backbone of PCA: principal components are the directions of  $u_i/v_i$ , ordered by decreasing captured variance  $\sigma_i^2$ .

### Exercises on a matrix $M$

*Exercise 7.7.1.* Compute the SVD of

$$M = \begin{bmatrix} 1 & 2 \\ 0 & 1 \\ 2 & 0 \end{bmatrix}.$$

**Solution.** First compute

$$M^\top M = \begin{bmatrix} 1 & 0 & 2 \\ 2 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 0 & 1 \\ 2 & 0 \end{bmatrix} = \begin{bmatrix} 5 & 2 \\ 2 & 5 \end{bmatrix}.$$

Eigenvalues of  $M^\top M$  solve  $\det \begin{bmatrix} 5-\lambda & 2 \\ 2 & 5-\lambda \end{bmatrix} = 0$ , i.e.  $(5-\lambda)^2 - 4 = 0$ , hence  $\lambda_1 = 7$  and  $\lambda_2 = 3$ .

Thus the *singular values* are

$$\sigma_1 = \sqrt{7}, \quad \sigma_2 = \sqrt{3}.$$

Associated (orthonormal) eigenvectors of  $M^\top M$  are

$$v_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad v_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix},$$

so we may take

$$V = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix}.$$

Next, define left singular vectors  $u_i = \frac{1}{\sigma_i} M v_i$ :

$$\begin{aligned} M v_1 &= \frac{1}{\sqrt{2}} \begin{bmatrix} 3 \\ 1 \\ 2 \end{bmatrix} \Rightarrow u_1 = \frac{1}{\sqrt{7}} M v_1 = \frac{1}{\sqrt{14}} \begin{bmatrix} 3 \\ 1 \\ 2 \end{bmatrix}, \\ M v_2 &= \frac{1}{\sqrt{2}} \begin{bmatrix} -1 \\ -1 \\ 2 \end{bmatrix} \Rightarrow u_2 = \frac{1}{\sqrt{3}} M v_2 = \frac{1}{\sqrt{6}} \begin{bmatrix} -1 \\ -1 \\ 2 \end{bmatrix}. \end{aligned}$$



Complete to an orthonormal basis of  $\mathbb{R}^3$  by choosing any  $u_3$  orthogonal to  $u_1, u_2$ , e.g. solve  $u_3 \cdot (3, 1, 2) = 0$  and  $u_3 \cdot (-1, -1, 2) = 0$ , giving  $u_3 \propto (-2, 4, 1)$ ; normalize:

$$u_3 = \frac{1}{\sqrt{21}} \begin{bmatrix} -2 \\ 4 \\ 1 \end{bmatrix}.$$

Thus

$$U = \begin{bmatrix} \frac{3}{\sqrt{14}} & -\frac{1}{\sqrt{6}} & -\frac{2}{\sqrt{21}} \\ \frac{1}{\sqrt{14}} & -\frac{1}{\sqrt{6}} & \frac{4}{\sqrt{21}} \\ \frac{2}{\sqrt{14}} & \frac{2}{\sqrt{6}} & \frac{1}{\sqrt{21}} \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \sqrt{7} & 0 \\ 0 & \sqrt{3} \\ 0 & 0 \end{bmatrix}.$$

One checks  $M = U\Sigma V^\top$ ,  $U^\top U = I_3$ ,  $V^\top V = I_2$ . □

*Exercise 7.7.2.* Verify that  $M = U\Sigma V^\top$  with the  $U, \Sigma, V$  found above. What is the best rank-one approximation of  $M$ ?

**Solution.** *Verification.* Compute  $U\Sigma$  (which has columns  $\sigma_1 u_1, \sigma_2 u_2$ ) and multiply by  $V^\top$ ; each entry matches  $M$ .

*Best rank-one approximation.* Since we denote by  $\sigma_1$  the largest singular value, the Eckart-Young-Mirsky theorem implies that the best rank-one approximation of  $M$  (in both the Frobenius and operator norms) is

$$M_1 = \sigma_1 u_1 v_1^\top.$$

Using  $u_1 = \frac{1}{\sqrt{14}}(3, 1, 2)^\top$  and  $v_1 = \frac{1}{\sqrt{2}}(1, 1)^\top$ ,

$$u_1 v_1^\top = \frac{1}{\sqrt{28}} \begin{bmatrix} 3 & 3 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} \Rightarrow M_1 = \sqrt{7} \cdot u_1 v_1^\top = \frac{\sqrt{7}}{\sqrt{28}} \begin{bmatrix} 3 & 3 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 3 & 3 \\ 1 & 1 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} 1.5 & 1.5 \\ 0.5 & 0.5 \\ 1 & 1 \end{bmatrix}.$$

□

*Exercise 7.7.3.* Write the rank- $r$  decomposition

$$M = \sum_{i=1}^r \sigma_i u_i v_i^\top,$$

and interpret it as a sum of rank-one matrices. Comment on its use in image compression.

**Solution.** For our  $3 \times 2$  matrix,  $r = \text{rank}(M) = 2$ , so

$$M = \sigma_1 u_1 v_1^\top + \sigma_2 u_2 v_2^\top.$$

Concretely,

$$\sigma_1 u_1 v_1^\top = \frac{1}{2} \begin{bmatrix} 3 & 3 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}, \quad \sigma_2 u_2 v_2^\top = \frac{1}{2} \begin{bmatrix} -1 & 1 \\ -1 & 1 \\ 2 & -2 \end{bmatrix},$$

and one checks their sum equals  $M$ . Each term  $\sigma_i u_i v_i^\top$  is a rank-one matrix (outer product), capturing “energy”  $\sigma_i^2$ . Truncating to the first  $k < r$  terms gives the best rank- $k$  approximation. In image compression, keeping only the largest singular values preserves most visual information while drastically reducing storage. □

*Exercise 7.7.4.* Show that the least-squares problem  $\min_x \|Mx - b\|_2$  can be solved via SVD as

$$x = V \Sigma^+ U^\top b,$$

where  $\Sigma^+$  is the Moore–Penrose pseudoinverse of  $\Sigma$ .

**Solution.** Let  $M = U \Sigma V^\top$  be an SVD with rank  $r$ , and set  $y = V^\top x$ . Since  $U$  is orthogonal,

$$\|Mx - b\|_2 = \|U^\top (Mx - b)\|_2 = \|\Sigma y - U^\top b\|_2.$$

Write  $U^\top b = (\beta_1, \dots, \beta_m)^\top$ . With  $\Sigma = \begin{bmatrix} \text{diag}(\sigma_1, \dots, \sigma_r) & 0 \\ 0 & 0 \end{bmatrix}$ , we minimize

$$\sum_{i=1}^r (\sigma_i y_i - \beta_i)^2 + \sum_{i=r+1}^m (\beta_i)^2.$$

The second sum is independent of  $y$ . The first sum is minimized coordinatewise by  $y_i = \beta_i / \sigma_i$  for  $i \leq r$ , while the minimum-norm choice is  $y_i = 0$  for  $i > r$ . Thus

$$y = \Sigma^+ U^\top b, \quad \Sigma^+ = \begin{bmatrix} \text{diag}(\sigma_1^{-1}, \dots, \sigma_r^{-1}) & 0 \\ 0 & 0 \end{bmatrix},$$

and hence

$$x = Vy = V \Sigma^+ U^\top b.$$

This is exactly the Moore–Penrose pseudoinverse solution. □

*Remark 34* (Why the SVD is Preferred for Computing the Pseudoinverse). In theory, if  $M$  has full column rank, one can write

$$M^+ = (M^\top M)^{-1} M^\top.$$

However, in practice this formula is numerically unstable and limited. The Singular Value Decomposition (SVD) provides a more robust and general approach:

- **Numerical stability:** Computing  $(M^\top M)^{-1}$  squares the condition number of  $M$ , which amplifies rounding errors. In contrast, the SVD works directly with singular values  $\sigma_i$  and avoids this instability.
- **Rank-deficient cases:** If  $M$  does not have full rank - which occurs, for example, when the number of variables exceeds the number of observations, a situation very common in practice -, then the formula  $(M^\top M)^{-1} M^\top$  does not apply, while the SVD naturally yields the Moore–Penrose pseudoinverse by inverting only the nonzero singular values.
- **Regularization:** The SVD makes it possible to ignore very small singular values (below a tolerance), thereby stabilizing the solution. This is known as truncated SVD and is widely used in statistics and data analysis.

For these reasons, modern numerical software (such as MATLAB, NumPy, or R) implements the pseudoinverse via the SVD rather than the normal equations.

## Chapter 8

# Matrix Conditioning and Numerical Stability

The concept of **matrix conditioning** is central to numerical analysis and applied linear algebra. It measures the *stability* of a linear problem and the *sensitivity* of its solution to small perturbations in the data. When solving a linear system

$$Ax = b,$$

a tiny perturbation in either  $A$  or  $b$  may cause a large variation in the computed solution  $x$ . The conditioning number quantifies this potential amplification of errors.

### 8.1 Matrix and Vector Norms

#### 8.1.1 Euclidean norm and induced matrix norm

For a vector  $x \in \mathbb{R}^n$ , the **Euclidean norm** is defined by:

$$\|x\|_2 = \sqrt{x^\top x} = \sqrt{\sum_{i=1}^n x_i^2}.$$

For a matrix  $A \in \mathbb{R}^{m \times n}$ , the **Frobenius norm** is given by:

$$\|A\|_F = \sqrt{\text{tr}(A^\top A)} = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2}.$$

This norm behaves like an Euclidean norm for matrices and is often used in practice.

**Definition 8.1.1.** Another important norm is the *induced matrix norm*:

$$\|A\|_2 = \sup_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2}.$$

**Property 8.1.1.** It equals the largest *singular value* of  $A$ , denoted  $\sigma_{\max}(A)$ .

**Proof.** By definition of the induced 2-norm,

$$\|A\|_2 = \sup_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \sup_{\|x\|_2=1} \|Ax\|_2 = \sup_{\|x\|_2=1} \sqrt{(Ax)^\top (Ax)} = \sup_{\|x\|_2=1} \sqrt{x^\top (A^\top A)x}.$$

The expression  $\frac{x^\top (A^\top A)x}{x^\top x}$  is the Rayleigh quotient of the symmetric positive semidefinite matrix  $A^\top A$ . Its maximum over  $x \neq 0$  equals the largest eigenvalue  $\lambda_{\max}(A^\top A)$ . Therefore

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^\top A)}.$$

By the singular value decomposition, the singular values of  $A$  are the square roots of the eigenvalues of  $A^\top A$ , hence  $\sqrt{\lambda_{\max}(A^\top A)} = \sigma_{\max}(A)$ .  $\square$

*Remark 35* (Geometric intuition). The induced 2-norm is the maximal radial stretching of the unit sphere by the linear map  $x \mapsto Ax$ . The directions of maximal stretching are given by the right singular vectors, and the stretch factor is the largest singular value  $\sigma_{\max}(A)$ .

### 8.1.2 Spectral radius and its relation to norms

**Definition 8.1.2.** The *spectral radius* of a square matrix  $A$  is:

$$\rho(A) = \max_i |\lambda_i(A)|,$$

where  $\lambda_i(A)$  are the eigenvalues of  $A$ .

The spectral radius captures the intrinsic "growth factor" carried by the eigenmodes of  $A$ . While a norm measures worst-case amplification over *all* directions,  $\rho(A)$  measures the amplification *along invariant directions*. Any submultiplicative norm bounds that growth, hence  $\rho(A) \leq \|A\|$ .

**Property 8.1.2.** For any submultiplicative matrix norm,

$$\rho(A) \leq \|A\|.$$

If  $A$  is normal (i.e.  $A^\top A = AA^\top$ ), equality holds for the spectral norm:

$$\|A\|_2 = \rho(A).$$

**Proof.** (i)  $\rho(A) \leq \|A\|$  for any submultiplicative norm. Let  $\lambda$  be any eigenvalue of  $A$  with eigenvector  $v \neq 0$ . Then  $A^k v = \lambda^k v$  for all  $k \in \mathbb{N}$ . Hence

$$\|A^k\| \geq \frac{\|A^k v\|}{\|v\|} = \frac{\|\lambda^k v\|}{\|v\|} = |\lambda|^k.$$

Taking  $k$ -th roots and letting  $k \rightarrow \infty$  gives  $|\lambda| \leq \|A\|$ . Maximizing over all eigenvalues yields  $\rho(A) \leq \|A\|$ .

(ii)  $\|A\|_2 = \rho(A)$  when  $A$  is normal. If  $A$  is normal, there exists an orthogonal matrix  $Q$  and a (possibly complex) diagonal  $\Lambda$  such that  $A = Q\Lambda Q^\top$  and the diagonal entries of  $\Lambda$  are the eigenvalues  $\lambda_i(A)$ . Normality implies the singular values of  $A$  are  $|\lambda_i(A)|$ . Therefore

$$\|A\|_2 = \sigma_{\max}(A) = \max_i |\lambda_i(A)| = \rho(A).$$

$\square$

*Remark 36* (Takeaway). Every submultiplicative norm upper-bounds the spectral radius. For normal matrices, the spectral norm is *exactly* the spectral radius, so operator growth and eigenmode growth coincide.

## 8.2 Definition of the Matrix Condition Number

Consider the linear system  $Ax = b$ , where  $A$  is invertible. A small perturbation  $\delta b$  in  $b$  produces a variation  $\delta x$  satisfying:

$$A(x + \delta x) = b + \delta b.$$

Hence,

$$A\delta x = \delta b \quad \Rightarrow \quad \delta x = A^{-1}\delta b.$$

### 8.2.1 Relative error amplification

We can relate the errors in  $x$  and  $b$  by:

$$\frac{\|\delta x\|}{\|x\|} = \frac{\|A^{-1}\delta b\|}{\|A^{-1}b\|} \leq \frac{\|A^{-1}\| \|A\| \|\delta b\|}{\|b\|}.$$

**Relative error bound (explicit steps).** Starting from  $A(x + \delta x) = b + \delta b$ , we have  $A\delta x = \delta b$ , hence

$$\|\delta x\| = \|A^{-1}\delta b\| \leq \|A^{-1}\| \|\delta b\| \quad (\text{definition of induced norm}).$$

On the other hand,  $b = Ax$  yields  $\|b\| = \|Ax\| \leq \|A\| \|x\|$ , hence  $\|x\| \geq \|b\|/\|A\|$ . Combining,

$$\frac{\|\delta x\|}{\|x\|} \leq \frac{\|A^{-1}\| \|\delta b\|}{\|b\|/\|A\|} = \|A^{-1}\| \|A\| \frac{\|\delta b\|}{\|b\|}.$$

**Definition 8.2.1.** The factor

$$\kappa(A) = \|A\| \|A^{-1}\|$$

is called the **condition number** of  $A$ .

**Property 8.2.1.** The condition number always satisfies  $\kappa(A) \geq 1$ .

**Proof.** By definition,

$$\kappa(A) = \|A\| \|A^{-1}\|.$$

Using the submultiplicative property of the norm,

$$\|AA^{-1}\| \leq \|A\| \|A^{-1}\|.$$

But  $AA^{-1} = I_d$ , and  $\|I_d\| = 1$  for any operator norm induced by a vector norm. Hence

$$1 = \|I_d\| \leq \|A\| \|A^{-1}\| = \kappa(A),$$

which proves that  $\kappa(A) \geq 1$ . □

*Remark 37* (Intuitive meaning). A condition number of 1 corresponds to an *isometric* linear transformation - one that preserves lengths exactly. Any deviation from this ideal case (stretching or compression in some directions) increases  $\kappa(A)$  beyond 1.

- If  $\kappa(A) \approx 1$ , the problem is *well-conditioned*.
- If  $\kappa(A)$  is large, the problem is *ill-conditioned*, meaning it is highly sensitive to perturbations.

### 8.2.2 Spectral expression (in 2-norm)

**Proposition 8.2.1** (Inverses of  $A^\top A$  and  $AA^\top$  via  $A^{-1}$ ). Let  $A \in \mathbb{R}^{n \times n}$  be invertible. Then

$$(A^\top A)^{-1} = A^{-1} (A^{-1})^\top \quad \text{and} \quad (AA^\top)^{-1} = (A^{-1})^\top A^{-1}.$$

In particular,  $A^\top A$  and  $AA^\top$  are symmetric positive definite (SPD), and the eigenvalues of their inverses are the reciprocals of the eigenvalues of  $A^\top A$  and  $AA^\top$ , respectively.

**Proof.** Use the identity  $(BC)^{-1} = C^{-1}B^{-1}$ , valid for any invertible matrices  $B, C$ .

*First identity.* Since  $A^\top A$  is invertible,

$$(A^\top A)^{-1} = A^{-1}(A^\top)^{-1} = A^{-1}(A^{-1})^\top,$$

because  $(A^\top)^{-1} = (A^{-1})^\top$ .

*Second identity.* Similarly,

$$(AA^\top)^{-1} = (A^\top)^{-1}A^{-1} = (A^{-1})^\top A^{-1}.$$

*SPD and reciprocal eigenvalues.* The matrices  $A^\top A$  and  $AA^\top$  are **Symmetric Positive Definite** (SPD): for any  $x \neq 0$ ,

$$x^\top (A^\top A) x = \|Ax\|_2^2 > 0$$

and

$$x^\top (AA^\top) x = \|A^\top x\|_2^2 > 0.$$

If  $M$  is SPD with eigenpair  $(\lambda, u)$ , then  $Mu = \lambda u$  implies  $M^{-1}u = \lambda^{-1}u$ . Applying this to  $M = A^\top A$  (resp.  $M = AA^\top$ ) shows that the eigenvalues of  $(A^\top A)^{-1}$  (resp.  $(AA^\top)^{-1}$ ) are precisely the reciprocals of those of  $A^\top A$  (resp.  $AA^\top$ ).  $\square$

*Remark 38.* Note that  $A^{-1}(A^{-1})^\top = (A^\top A)^{-1}$  and  $(A^{-1})^\top A^{-1} = (AA^\top)^{-1}$  are generally *different* matrices (they coincide iff  $A$  is normal and, in particular, symmetric). Both are SPD and share the same nonzero eigenvalues (since  $A^\top A$  and  $AA^\top$  do).

**Property 8.2.2.** *In the Euclidean norm,*

$$\|A\|_2 = \sigma_{\max}(A), \quad \|A^{-1}\|_2 = \frac{1}{\sigma_{\min}(A)}.$$

*Thus,*

$$\kappa_2(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}.$$

*For symmetric positive definite matrices, this simplifies to:*

$$\kappa_2(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}.$$

**Proof.** The first identity  $\|A\|_2 = \sigma_{\max}(A)$  follows from the Rayleigh–Ritz characterization proved above (8.1.1):

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^\top A)} = \sigma_{\max}(A).$$

If  $A$  is invertible, then  $(A^{-1})^\top A^{-1} = (A^\top A)^{-1}$  and the eigenvalues of  $(A^{-1})^\top A^{-1}$  are the reciprocals of those of  $A^\top A$ . Therefore

$$\|A^{-1}\|_2^2 = \lambda_{\max}((A^{-1})^\top A^{-1}) = \frac{1}{\lambda_{\min}(A^\top A)} \implies \|A^{-1}\|_2 = \frac{1}{\sqrt{\lambda_{\min}(A^\top A)}} = \frac{1}{\sigma_{\min}(A)}.$$

Multiplying the two expressions yields

$$\kappa_2(A) = \|A\|_2 \|A^{-1}\|_2 = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}.$$

If  $A$  is symmetric positive definite, then  $A$  is orthogonally diagonalizable with positive eigenvalues and  $A^\top A = A^2$ . The singular values of  $A$  are the absolute values of its eigenvalues, hence  $\sigma_{\max}(A) = \lambda_{\max}(A)$  and  $\sigma_{\min}(A) = \lambda_{\min}(A)$ , which gives

$$\kappa_2(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}.$$

$\square$

## 8.3 Ill-conditioning and the Explosion of the Inverse

### 8.3.1 An intuitive illustration

Consider a nearly singular matrix:

$$A_\varepsilon = \begin{pmatrix} 1 & 1 \\ 1 & 1 + \varepsilon \end{pmatrix},$$

whose determinant is  $\det(A_\varepsilon) = \varepsilon$ . As  $\varepsilon \rightarrow 0$ , the columns of  $A_\varepsilon$  become nearly linearly dependent, and  $A_\varepsilon$  becomes ill-conditioned.

### 8.3.2 Computing the inverse

A direct computation gives:

$$A_\varepsilon^{-1} = \frac{1}{\varepsilon} \begin{pmatrix} 1 + \varepsilon & -1 \\ -1 & 1 \end{pmatrix}.$$

Hence,

$$\|A_\varepsilon^{-1}\|_2 \underset{\varepsilon \rightarrow 0}{\sim} \frac{C}{|\varepsilon|}.$$

A tiny perturbation in the matrix causes a massive amplification in the inverse's norm. This is the hallmark of an *ill-conditioned matrix*.

### 8.3.3 Rigorous formalism

Let  $A$  be invertible and  $E$  a perturbation such that  $\|E\| < 1/\|A^{-1}\|$ . Then  $A + E$  is invertible and

$$(A + E)^{-1} = (I + A^{-1}E)^{-1}A^{-1} = \sum_{k=0}^{\infty} (-A^{-1}E)^k A^{-1}.$$

This leads to the bound

$$\|(A + E)^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\|\|E\|}.$$

As soon as  $\|A^{-1}\|\|E\|$  approaches 1, the denominator tends to zero, and thus the norm of the inverse **blows up**. This provides a rigorous explanation of why ill-conditioned matrices are numerically unstable.

## 8.4 Main Results and Theorems

**Theorem 8.4.1** (Condition number and sensitivity). *Let  $A$  be invertible and  $b$  a vector. For any perturbation  $\delta b$ , we have:*

$$\frac{\|\delta x\|}{\|x\|} \leq \kappa(A) \frac{\|\delta b\|}{\|b\|}.$$

**Sketch of proof.** From  $A(x + \delta x) = b + \delta b$  we get  $A\delta x = \delta b$ , hence

$$\|\delta x\| = \|A^{-1}\delta b\| \leq \|A^{-1}\| \|\delta b\|.$$

Also  $b = Ax$  implies  $\|b\| = \|Ax\| \leq \|A\| \|x\|$ , so  $\|x\| \geq \|b\|/\|A\|$ . Combine:

$$\frac{\|\delta x\|}{\|x\|} \leq \|A^{-1}\| \|\delta b\| \cdot \frac{\|A\|}{\|b\|} = \|A\| \|A^{-1}\| \frac{\|\delta b\|}{\|b\|} = \kappa(A) \frac{\|\delta b\|}{\|b\|}.$$

□

*Remark 39.* Solving  $Ax = b$  first stretches  $x$  by at most  $\|A\|$  to produce  $b$ , while inverting stretches back by at most  $\|A^{-1}\|$ . The product  $\kappa(A)$  is the worst-case amplification factor of relative errors from data to solution.

**Theorem 8.4.2** (Near singularity criterion). *A matrix  $A$  is close to singularity if and only if its condition number is large:*

$$\kappa(A) \gg 1 \iff \sigma_{\min}(A) \approx 0.$$

**Sketch of proof.** In the 2-norm,  $\kappa_2(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}$ .

( $\Rightarrow$ ) If  $\kappa_2(A)$  is large while  $\sigma_{\max}(A)$  is bounded (e.g. after scaling  $\|A\|_2$  to 1), then necessarily  $\sigma_{\min}(A)$  is small, hence  $A$  is close to losing invertibility.

( $\Leftarrow$ ) If  $\sigma_{\min}(A) \rightarrow 0$  and  $\sigma_{\max}(A) > 0$ , then  $\kappa_2(A) = \sigma_{\max}(A)/\sigma_{\min}(A) \rightarrow \infty$ . Thus small singular values force large condition numbers.  $\square$

*Remark 40.* The smallest singular value measures the least stretch of  $A$ . When that stretch collapses toward zero,  $A$  nearly flattens some direction, making inversion extremely sensitive; the ratio with the largest stretch explodes.

**Theorem 8.4.3** (Perturbation of the inverse). *If  $\|A^{-1}\|\|E\| < 1$ , then:*

$$(A + E)^{-1} - A^{-1} = -A^{-1}E(A + E)^{-1},$$

and consequently,

$$\frac{\|(A + E)^{-1} - A^{-1}\|}{\|A^{-1}\|} \leq \frac{\kappa(A)\|E\|/\|A\|}{1 - \kappa(A)\|E\|/\|A\|}.$$

**Sketch of proof.** Identity:

$$(A + E)^{-1} - A^{-1} = (A + E)^{-1}[(A + E) - A]A^{-1} = (A + E)^{-1}EA^{-1} = -A^{-1}E(A + E)^{-1}.$$

Norm bound: by submultiplicativity,

$$\|(A + E)^{-1} - A^{-1}\| \leq \|A^{-1}\| \|E\| \|(A + E)^{-1}\|.$$

If  $\|A^{-1}\|\|E\| < 1$ , Neumann series gives

$$(A + E)^{-1} = (I + A^{-1}E)^{-1}A^{-1}, \quad \|(A + E)^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\|\|E\|}.$$

Hence

$$\|(A + E)^{-1} - A^{-1}\| \leq \frac{\|A^{-1}\|^2 \|E\|}{1 - \|A^{-1}\|\|E\|}.$$

Divide by  $\|A^{-1}\|$  and rewrite  $\|A^{-1}\| = \kappa(A)/\|A\|$  to obtain

$$\frac{\|(A + E)^{-1} - A^{-1}\|}{\|A^{-1}\|} \leq \frac{\kappa(A)\|E\|/\|A\|}{1 - \kappa(A)\|E\|/\|A\|}.$$

$\square$

*Remark 41.* The inverse map is Lipschitz on a ball that stays away from singularity; the denominator  $1 - \|A^{-1}\|\|E\|$  is the safety margin to singularity. As you push  $A + E$  toward a singular matrix, that margin shrinks and the change in the inverse explodes.

## 8.5 Conclusion

The condition number of a matrix measures how **robust or fragile** a linear system is with respect to data perturbations. Well-conditioned matrices ensure that small input errors lead to small output errors, while ill-conditioned matrices can amplify even the tiniest numerical inaccuracies, leading to severe instability.

Understanding matrix conditioning is essential in modern numerical analysis, optimization, and machine learning, where large-scale linear problems must be solved reliably despite rounding and data uncertainties.