

Basics of Algebra, Topology, and Differential Calculus

Jean Gallier
Department of Computer and Information Science
University of Pennsylvania
Philadelphia, PA 19104, USA
e-mail: jean@cis.upenn.edu

© Jean Gallier

August 28, 2015

Contents

1	Introduction	9
2	Vector Spaces, Bases, Linear Maps	11
2.1	Groups, Rings, and Fields	11
2.2	Vector Spaces	22
2.3	Linear Independence, Subspaces	29
2.4	Bases of a Vector Space	34
2.5	Linear Maps	41
2.6	Quotient Spaces	48
2.7	Summary	49
3	Matrices and Linear Maps	51
3.1	Matrices	51
3.2	Haar Basis Vectors and a Glimpse at Wavelets	67
3.3	The Effect of a Change of Bases on Matrices	83
3.4	Summary	86
4	Direct Sums, The Dual Space, Duality	87
4.1	Sums, Direct Sums, Direct Products	87
4.2	The Dual Space E^* and Linear Forms	100
4.3	Hyperplanes and Linear Forms	113
4.4	Transpose of a Linear Map and of a Matrix	114
4.5	The Four Fundamental Subspaces	123
4.6	Summary	125
5	Determinants	129
5.1	Permutations, Signature of a Permutation	129
5.2	Alternating Multilinear Maps	133
5.3	Definition of a Determinant	136
5.4	Inverse Matrices and Determinants	142
5.5	Systems of Linear Equations and Determinants	146
5.6	Determinant of a Linear Map	147
5.7	The Cayley–Hamilton Theorem	147

5.8	Permanents	152
5.9	Further Readings	155
6	Gaussian Elimination, LU, Cholesky, Echelon Form	157
6.1	Motivating Example: Curve Interpolation	157
6.2	Gaussian Elimination and LU -Factorization	161
6.3	Gaussian Elimination of Tridiagonal Matrices	183
6.4	SPD Matrices and the Cholesky Decomposition	186
6.5	Reduced Row Echelon Form	190
6.6	Transvections and Dilatations	205
6.7	Summary	211
7	Vector Norms and Matrix Norms	213
7.1	Normed Vector Spaces	213
7.2	Matrix Norms	219
7.3	Condition Numbers of Matrices	232
7.4	An Application of Norms: Inconsistent Linear Systems	240
7.5	Summary	242
8	Eigenvectors and Eigenvalues	245
8.1	Eigenvectors and Eigenvalues of a Linear Map	245
8.2	Reduction to Upper Triangular Form	252
8.3	Location of Eigenvalues	256
8.4	Summary	258
9	Iterative Methods for Solving Linear Systems	259
9.1	Convergence of Sequences of Vectors and Matrices	259
9.2	Convergence of Iterative Methods	262
9.3	Methods of Jacobi, Gauss-Seidel, and Relaxation	264
9.4	Convergence of the Methods	269
9.5	Summary	276
10	Euclidean Spaces	277
10.1	Inner Products, Euclidean Spaces	277
10.2	Orthogonality, Duality, Adjoint of a Linear Map	285
10.3	Linear Isometries (Orthogonal Transformations)	296
10.4	The Orthogonal Group, Orthogonal Matrices	299
10.5	QR -Decomposition for Invertible Matrices	301
10.6	Some Applications of Euclidean Geometry	305
10.7	Summary	306

11	<i>QR</i>-Decomposition for Arbitrary Matrices	309
11.1	Orthogonal Reflections	309
11.2	<i>QR</i> -Decomposition Using Householder Matrices	312
11.3	Summary	317
12	Hermitian Spaces	319
12.1	Hermitian Spaces, Pre-Hilbert Spaces	319
12.2	Orthogonality, Duality, Adjoint of a Linear Map	328
12.3	Linear Isometries (Also Called Unitary Transformations)	333
12.4	The Unitary Group, Unitary Matrices	335
12.5	Orthogonal Projections and Involutions	338
12.6	Dual Norms	341
12.7	Summary	344
13	Spectral Theorems	347
13.1	Introduction	347
13.2	Normal Linear Maps	347
13.3	Self-Adjoint and Other Special Linear Maps	356
13.4	Normal and Other Special Matrices	363
13.5	Conditioning of Eigenvalue Problems	366
13.6	Rayleigh Ratios and the Courant-Fischer Theorem	369
13.7	Summary	377
14	Bilinear Forms and Their Geometries	379
14.1	Bilinear Forms	379
14.2	Sesquilinear Forms	386
14.3	Orthogonality	390
14.4	Adjoint of a Linear Map	394
14.5	Isometries Associated with Sesquilinear Forms	397
14.6	Totally Isotropic Subspaces. Witt Decomposition	400
14.7	Witt's Theorem	413
14.8	Symplectic Groups	417
14.9	Orthogonal Groups	421
15	Introduction to The Finite Elements Method	429
15.1	A One-Dimensional Problem: Bending of a Beam	429
15.2	A Two-Dimensional Problem: An Elastic Membrane	440
15.3	Time-Dependent Boundary Problems	443
16	Singular Value Decomposition and Polar Form	451
16.1	Singular Value Decomposition for Square Matrices	451
16.2	Singular Value Decomposition for Rectangular Matrices	459
16.3	Ky Fan Norms and Schatten Norms	462

16.4 Summary	463
17 Applications of SVD and Pseudo-Inverses	465
17.1 Least Squares Problems and the Pseudo-Inverse	465
17.2 Data Compression and SVD	475
17.3 Principal Components Analysis (PCA)	476
17.4 Best Affine Approximation	483
17.5 Summary	486
18 Quadratic Optimization Problems	489
18.1 Quadratic Optimization: The Positive Definite Case	489
18.2 Quadratic Optimization: The General Case	497
18.3 Maximizing a Quadratic Function on the Unit Sphere	501
18.4 Summary	506
19 Basics of Affine Geometry	507
19.1 Affine Spaces	507
19.2 Examples of Affine Spaces	514
19.3 Chasles's Identity	516
19.4 Affine Combinations, Barycenters	517
19.5 Affine Subspaces	520
19.6 Affine Independence and Affine Frames	525
19.7 Affine Maps	530
19.8 Affine Groups	537
19.9 Affine Geometry: A Glimpse	539
19.10 Affine Hyperplanes	543
19.11 Intersection of Affine Spaces	545
19.12 Problems	547
20 Polynomials, Ideals and PID's	561
20.1 Multisets	561
20.2 Polynomials	562
20.3 Euclidean Division of Polynomials	568
20.4 Ideals, PID's, and Greatest Common Divisors	570
20.5 Factorization and Irreducible Factors in $K[X]$	578
20.6 Roots of Polynomials	582
20.7 Polynomial Interpolation (Lagrange, Newton, Hermite)	589
21 UFD's, Noetherian Rings, Hilbert's Basis Theorem	595
21.1 Unique Factorization Domains (Factorial Rings)	595
21.2 The Chinese Remainder Theorem	609
21.3 Noetherian Rings and Hilbert's Basis Theorem	614

21.4	Futher Readings	618
22	Annihilating Polynomials; Primary Decomposition	619
22.1	Annihilating Polynomials and the Minimal Polynomial	619
22.2	Minimal Polynomials of Diagonalizable Linear Maps	621
22.3	The Primary Decomposition Theorem	627
22.4	Nilpotent Linear Maps and Jordan Form	633
23	Tensor Algebras	639
23.1	Tensors Products	639
23.2	Bases of Tensor Products	647
23.3	Some Useful Isomorphisms for Tensor Products	649
23.4	Duality for Tensor Products	650
23.5	Tensor Algebras	653
23.6	Symmetric Tensor Powers	658
23.7	Bases of Symmetric Powers	662
23.8	Some Useful Isomorphisms for Symmetric Powers	664
23.9	Duality for Symmetric Powers	664
23.10	Symmetric Algebras	666
23.11	Exterior Tensor Powers	668
23.12	Bases of Exterior Powers	672
23.13	Some Useful Isomorphisms for Exterior Powers	674
23.14	Duality for Exterior Powers	675
23.15	Exterior Algebras	677
23.16	The Hodge *-Operator	680
23.17	Testing Decomposability; Left and Right Hooks	682
23.18	Vector-Valued Alternating Forms	689
23.19	The Pfaffian Polynomial	692
24	Introduction to Modules; Modules over a PID	697
24.1	Modules over a Commutative Ring	697
24.2	Finite Presentations of Modules	705
24.3	Tensor Products of Modules over a Commutative Ring	710
24.4	Extension of the Ring of Scalars	713
24.5	The Torsion Module Associated With An Endomorphism	716
24.6	Torsion Modules over a PID; Primary Decomposition	724
24.7	Finitely Generated Modules over a PID	729
25	Normal Forms; The Rational Canonical Form	745
25.1	The Rational Canonical Form	745
25.2	The Rational Canonical Form, Second Version	750
25.3	The Jordan Form Revisited	751
25.4	The Smith Normal Form	753

26 Topology	767
26.1 Metric Spaces and Normed Vector Spaces	767
26.2 Topological Spaces	771
26.3 Continuous Functions, Limits	776
26.4 Connected Sets	781
26.5 Compact Sets	787
26.6 Continuous Linear and Multilinear Maps	801
26.7 Normed Affine Spaces	806
26.8 Futher Readings	806
27 A Detour On Fractals	807
27.1 Iterated Function Systems and Fractals	807
28 Differential Calculus	815
28.1 Directional Derivatives, Total Derivatives	815
28.2 Jacobian Matrices	823
28.3 The Implicit and The Inverse Function Theorems	828
28.4 Tangent Spaces and Differentials	832
28.5 Second-Order and Higher-Order Derivatives	833
28.6 Taylor's formula, Faà di Bruno's formula	839
28.7 Vector Fields, Covariant Derivatives, Lie Brackets	843
28.8 Futher Readings	845
29 Extrema of Real-Valued Functions	847
29.1 Local Extrema and Lagrange Multipliers	847
29.2 Using Second Derivatives to Find Extrema	856
29.3 Using Convexity to Find Extrema	859
29.4 Summary	867
30 Newton's Method and its Generalizations	869
30.1 Newton's Method for Real Functions of a Real Argument	869
30.2 Generalizations of Newton's Method	870
30.3 Summary	876
31 Appendix: Zorn's Lemma; Some Applications	877
31.1 Statement of Zorn's Lemma	877
31.2 Proof of the Existence of a Basis in a Vector Space	878
31.3 Existence of Maximal Proper Ideals	879
Bibliography	879

Chapter 1

Introduction

Chapter 2

Vector Spaces, Bases, Linear Maps

2.1 Groups, Rings, and Fields

In the following three chapters, the basic algebraic structures (groups, rings, fields, vector spaces) are reviewed, with a major emphasis on vector spaces. Basic notions of linear algebra such as vector spaces, subspaces, linear combinations, linear independence, bases, quotient spaces, linear maps, matrices, change of bases, direct sums, linear forms, dual spaces, hyperplanes, transpose of a linear maps, are reviewed.

The set \mathbb{R} of real numbers has two operations $+: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ (addition) and $*: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ (multiplication) satisfying properties that make \mathbb{R} into an abelian group under $+$, and $\mathbb{R} - \{0\} = \mathbb{R}^*$ into an abelian group under $*$. Recall the definition of a group.

Definition 2.1. A *group* is a set G equipped with a binary operation $\cdot: G \times G \rightarrow G$ that associates an element $a \cdot b \in G$ to every pair of elements $a, b \in G$, and having the following properties: \cdot is associative, has an identity element $e \in G$, and every element in G is invertible (w.r.t. \cdot). More explicitly, this means that the following equations hold for all $a, b, c \in G$:

$$(G1) \quad a \cdot (b \cdot c) = (a \cdot b) \cdot c. \quad (\text{associativity});$$

$$(G2) \quad a \cdot e = e \cdot a = a. \quad (\text{identity});$$

$$(G3) \quad \text{For every } a \in G, \text{ there is some } a^{-1} \in G \text{ such that } a \cdot a^{-1} = a^{-1} \cdot a = e. \quad (\text{inverse}).$$

A group G is *abelian* (or *commutative*) if

$$a \cdot b = b \cdot a \quad \text{for all } a, b \in G.$$

A set M together with an operation $\cdot: M \times M \rightarrow M$ and an element e satisfying only conditions (G1) and (G2) is called a *monoid*. For example, the set $\mathbb{N} = \{0, 1, \dots, n, \dots\}$ of natural numbers is a (commutative) monoid under addition. However, it is not a group.

Some examples of groups are given below.

Example 2.1.

1. The set $\mathbb{Z} = \{\dots, -n, \dots, -1, 0, 1, \dots, n, \dots\}$ of integers is a group under addition, with identity element 0. However, $\mathbb{Z}^* = \mathbb{Z} - \{0\}$ is not a group under multiplication.
2. The set \mathbb{Q} of rational numbers (fractions p/q with $p, q \in \mathbb{Z}$ and $q \neq 0$) is a group under addition, with identity element 0. The set $\mathbb{Q}^* = \mathbb{Q} - \{0\}$ is also a group under multiplication, with identity element 1.
3. Given any nonempty set S , the set of bijections $f: S \rightarrow S$, also called *permutations of S* , is a group under function composition (i.e., the multiplication of f and g is the composition $g \circ f$), with identity element the identity function id_S . This group is not abelian as soon as S has more than two elements.
4. The set of $n \times n$ invertible matrices with real (or complex) coefficients is a group under matrix multiplication, with identity element the identity matrix I_n . This group is called the *general linear group* and is usually denoted by $\mathbf{GL}(n, \mathbb{R})$ (or $\mathbf{GL}(n, \mathbb{C})$).

It is customary to denote the operation of an abelian group G by $+$, in which case the inverse a^{-1} of an element $a \in G$ is denoted by $-a$.

The identity element of a group is *unique*. In fact, we can prove a more general fact:

Fact 1. If a binary operation $\cdot: M \times M \rightarrow M$ is associative and if $e' \in M$ is a left identity and $e'' \in M$ is a right identity, which means that

$$e' \cdot a = a \quad \text{for all } a \in M \tag{G2l}$$

and

$$a \cdot e'' = a \quad \text{for all } a \in M, \tag{G2r}$$

then $e' = e''$.

Proof. If we let $a = e''$ in equation (G2l), we get

$$e' \cdot e'' = e'',$$

and if we let $a = e'$ in equation (G2r), we get

$$e' \cdot e'' = e',$$

and thus

$$e' = e' \cdot e'' = e'',$$

as claimed. □

Fact 1 implies that the identity element of a monoid is unique, and since every group is a monoid, the identity element of a group is unique. Furthermore, every element in a group has a *unique inverse*. This is a consequence of a slightly more general fact:

Fact 2. In a monoid M with identity element e , if some element $a \in M$ has some left inverse $a' \in M$ and some right inverse $a'' \in M$, which means that

$$a' \cdot a = e \tag{G3l}$$

and

$$a \cdot a'' = e, \tag{G3r}$$

then $a' = a''$.

Proof. Using (G3l) and the fact that e is an identity element, we have

$$(a' \cdot a) \cdot a'' = e \cdot a'' = a''.$$

Similarly, Using (G3r) and the fact that e is an identity element, we have

$$a' \cdot (a \cdot a'') = a' \cdot e = a'.$$

However, since M is monoid, the operation \cdot is associative, so

$$a' = a' \cdot (a \cdot a'') = (a' \cdot a) \cdot a'' = a'',$$

as claimed. □

Remark: Axioms (G2) and (G3) can be weakened a bit by requiring only (G2r) (the existence of a right identity) and (G3r) (the existence of a right inverse for every element) (or (G2l) and (G3l)). It is a good exercise to prove that the group axioms (G2) and (G3) follow from (G2r) and (G3r).

If a group G has a finite number n of elements, we say that G is a group of *order* n . If G is infinite, we say that G has *infinite order*. The order of a group is usually denoted by $|G|$ (if G is finite).

Given a group G , for any two subsets $R, S \subseteq G$, we let

$$RS = \{r \cdot s \mid r \in R, s \in S\}.$$

In particular, for any $g \in G$, if $R = \{g\}$, we write

$$gS = \{g \cdot s \mid s \in S\},$$

and similarly, if $S = \{g\}$, we write

$$Rg = \{r \cdot g \mid r \in R\}.$$

From now on, we will drop the multiplication sign and write g_1g_2 for $g_1 \cdot g_2$.

For any $g \in G$, define L_g , the *left translation by g* , by $L_g(a) = ga$, for all $a \in G$, and R_g , the *right translation by g* , by $R_g(a) = ag$, for all $a \in G$. Observe that L_g and R_g are bijections. We show this for L_g , the proof for R_g being similar.

If $L_g(a) = L_g(b)$, then $ga = gb$, and multiplying on the left by g^{-1} , we get $a = b$, so L_g is injective. For any $b \in G$, we have $L_g(g^{-1}b) = gg^{-1}b = b$, so L_g is surjective. Therefore, L_g is bijective.

Definition 2.2. Given a group G , a subset H of G is a *subgroup of G* iff

- (1) The identity element e of G also belongs to H ($e \in H$);
- (2) For all $h_1, h_2 \in H$, we have $h_1h_2 \in H$;
- (3) For all $h \in H$, we have $h^{-1} \in H$.

The proof of the following proposition is left as an exercise.

Proposition 2.1. *Given a group G , a subset $H \subseteq G$ is a subgroup of G iff H is nonempty and whenever $h_1, h_2 \in H$, then $h_1h_2^{-1} \in H$.*

If the group G is finite, then the following criterion can be used.

Proposition 2.2. *Given a finite group G , a subset $H \subseteq G$ is a subgroup of G iff*

- (1) $e \in H$;
- (2) H is closed under multiplication.

Proof. We just have to prove that condition (3) of Definition 2.2 holds. For any $a \in H$, since the left translation L_a is bijective, its restriction to H is injective, and since H is finite, it is also bijective. Since $e \in H$, there is a unique $b \in H$ such that $L_a(b) = ab = e$. However, if a^{-1} is the inverse of a in G , we also have $L_a(a^{-1}) = aa^{-1} = e$, and by injectivity of L_a , we have $a^{-1} = b \in H$. \square

Definition 2.3. If H is a subgroup of G and $g \in G$ is any element, the sets of the form gH are called *left cosets of H in G* and the sets of the form Hg are called *right cosets of H in G* .

The left cosets (resp. right cosets) of H induce an equivalence relation \sim defined as follows: For all $g_1, g_2 \in G$,

$$g_1 \sim g_2 \quad \text{iff} \quad g_1H = g_2H$$

(resp. $g_1 \sim g_2$ iff $Hg_1 = Hg_2$). Obviously, \sim is an equivalence relation.

Now, we claim that $g_1H = g_2H$ iff $g_2^{-1}g_1H = H$ iff $g_2^{-1}g_1 \in H$.

Proof. If we apply the bijection $L_{g_2^{-1}}$ to both g_1H and g_2H we get $L_{g_2^{-1}}(g_1H) = g_2^{-1}g_1H$ and $L_{g_2^{-1}}(g_2H) = H$, so $g_1H = g_2H$ iff $g_2^{-1}g_1H = H$. If $g_2^{-1}g_1H = H$, since $1 \in H$, we get $g_2^{-1}g_1 \in H$. Conversely, if $g_2^{-1}g_1 \in H$, since H is a group, the left translation $L_{g_2^{-1}g_1}$ is a bijection of H , so $g_2^{-1}g_1H = H$. Thus, $g_2^{-1}g_1H = H$ iff $g_2^{-1}g_1 \in H$. \square

It follows that the equivalence class of an element $g \in G$ is the coset gH (resp. Hg). Since L_g is a bijection between H and gH , the cosets gH all have the same cardinality. The map $L_{g^{-1}} \circ R_g$ is a bijection between the left coset gH and the right coset Hg , so they also have the same cardinality. Since the distinct cosets gH form a partition of G , we obtain the following fact:

Proposition 2.3. (*Lagrange*) *For any finite group G and any subgroup H of G , the order h of H divides the order n of G .*

The ratio n/h is denoted by $(G : H)$ and is called the *index of H in G* . The index $(G : H)$ is the number of left (and right) cosets of H in G . Proposition 2.3 can be stated as

$$|G| = (G : H)|H|.$$

The set of left cosets of H in G (which, in general, is **not** a group) is denoted G/H . The “points” of G/H are obtained by “collapsing” all the elements in a coset into a single element.

It is tempting to define a multiplication operation on left cosets (or right cosets) by setting

$$(g_1H)(g_2H) = (g_1g_2)H,$$

but this operation is not well defined in general, unless the subgroup H possesses a special property. This property is typical of the kernels of group homomorphisms, so we are led to

Definition 2.4. Given any two groups G and G' , a function $\varphi: G \rightarrow G'$ is a *homomorphism* iff

$$\varphi(g_1g_2) = \varphi(g_1)\varphi(g_2), \quad \text{for all } g_1, g_2 \in G.$$

Taking $g_1 = g_2 = e$ (in G), we see that

$$\varphi(e) = e',$$

and taking $g_1 = g$ and $g_2 = g^{-1}$, we see that

$$\varphi(g^{-1}) = \varphi(g)^{-1}.$$

If $\varphi: G \rightarrow G'$ and $\psi: G' \rightarrow G''$ are group homomorphisms, then $\psi \circ \varphi: G \rightarrow G''$ is also a homomorphism. If $\varphi: G \rightarrow G'$ is a homomorphism of groups, and $H \subseteq G$, $H' \subseteq G'$ are two subgroups, then it is easily checked that

$$\text{Im } H = \varphi(H) = \{\varphi(g) \mid g \in H\}$$

is a subgroup of G' called the *image of H by φ* , and

$$\varphi^{-1}(H') = \{g \in G \mid \varphi(g) \in H'\}$$

is a subgroup of G . In particular, when $H' = \{e'\}$, we obtain the *kernel* $\text{Ker } \varphi$ of φ . Thus,

$$\text{Ker } \varphi = \{g \in G \mid \varphi(g) = e'\}.$$

It is immediately verified that $\varphi: G \rightarrow G'$ is injective iff $\text{Ker } \varphi = \{e\}$. (We also write $\text{Ker } \varphi = (0)$.) We say that φ is an *isomorphism* if there is a homomorphism $\psi: G' \rightarrow G$, so that

$$\psi \circ \varphi = \text{id}_G \quad \text{and} \quad \varphi \circ \psi = \text{id}_{G'}.$$

In this case, ψ is unique and it is denoted φ^{-1} . When φ is an isomorphism we say the groups G and G' are *isomorphic*. It is easy to see that a bijective homomorphism is an isomorphism. When $G' = G$, a group isomorphism is called an *automorphism*. The left translations L_g and the right translations R_g are automorphisms of G .

We claim that $H = \text{Ker } \varphi$ satisfies the following property:

$$gH = Hg, \quad \text{for all } g \in G. \quad (*)$$

First, note that $(*)$ is equivalent to

$$gHg^{-1} = H, \quad \text{for all } g \in G,$$

and the above is equivalent to

$$gHg^{-1} \subseteq H, \quad \text{for all } g \in G. \quad (**)$$

This is because $gHg^{-1} \subseteq H$ implies $H \subseteq g^{-1}Hg$, and this for all $g \in G$. But,

$$\varphi(ghg^{-1}) = \varphi(g)\varphi(h)\varphi(g^{-1}) = \varphi(g)e'\varphi(g)^{-1} = \varphi(g)\varphi(g)^{-1} = e',$$

for all $h \in H = \text{Ker } \varphi$ and all $g \in G$. Thus, by definition of $H = \text{Ker } \varphi$, we have $gHg^{-1} \subseteq H$.

Definition 2.5. For any group G , a subgroup N of G is a *normal subgroup* of G iff

$$gNg^{-1} = N, \quad \text{for all } g \in G.$$

This is denoted by $N \triangleleft G$.

Observe that if G is abelian, then *every* subgroup of G is normal.

If N is a normal subgroup of G , the equivalence relation induced by left cosets is the same as the equivalence induced by right cosets. Furthermore, this equivalence relation \sim is a *congruence*, which means that: For all $g_1, g_2, g'_1, g'_2 \in G$,

- (1) If $g_1N = g'_1N$ and $g_2N = g'_2N$, then $g_1g_2N = g'_1g'_2N$, and
- (2) If $g_1N = g_2N$, then $g_1^{-1}N = g_2^{-1}N$.

As a consequence, we can define a group structure on the set G/\sim of equivalence classes modulo \sim , by setting

$$(g_1N)(g_2N) = (g_1g_2)N.$$

This group is denoted G/N and called the *quotient of G by N* . The equivalence class gN of an element $g \in G$ is also denoted \bar{g} (or $[g]$). The map $\pi: G \rightarrow G/N$ given by

$$\pi(g) = \bar{g} = gNx$$

is clearly a group homomorphism called the *canonical projection*.

Given a homomorphism of groups $\varphi: G \rightarrow G'$, we easily check that the groups $G/\text{Ker } \varphi$ and $\text{Im } \varphi = \varphi(G)$ are isomorphic. This is often called the *first isomorphism theorem*.

A useful way to construct groups is the *direct product* construction. Given two groups G and H , we let $G \times H$ be the Cartesian product of the sets G and H with the multiplication operation \cdot given by

$$(g_1, h_1) \cdot (g_2, h_2) = (g_1g_2, h_1h_2).$$

It is immediately verified that $G \times H$ is a group. Similarly, given any n groups G_1, \dots, G_n , we can define the direct product $G_1 \times \dots \times G_n$ in a similar way.

If G is an abelian group and H_1, \dots, H_n are subgroups of G , the situation is simpler. Consider the map

$$a: H_1 \times \dots \times H_n \rightarrow G$$

given by

$$a(h_1, \dots, h_n) = h_1 + \dots + h_n,$$

using $+$ for the operation of the group G . It is easy to verify that a is a group homomorphism, so its image is a subgroup of G denoted by $H_1 + \dots + H_n$, and called the *sum* of the groups H_i . The following proposition will be needed.

Proposition 2.4. *Given an abelian group G , if H_1 and H_2 are any subgroups of G such that $H_1 \cap H_2 = \{0\}$, then the map a is an isomorphism*

$$a: H_1 \times H_2 \rightarrow H_1 + H_2.$$

Proof. The map is surjective by definition, so we just have to check that it is injective. For this, we show that $\text{Ker } a = \{(0, 0)\}$. We have $a(a_1, a_2) = 0$ iff $a_1 + a_2 = 0$ iff $a_1 = -a_2$. Since $a_1 \in H_1$ and $a_2 \in H_2$, we see that $a_1, a_2 \in H_1 \cap H_2 = \{0\}$, so $a_1 = a_2 = 0$, which proves that $\text{Ker } a = \{(0, 0)\}$. \square

Under the conditions of Proposition 2.4, namely $H_1 \cap H_2 = \{0\}$, the group $H_1 + H_2$ is called the *direct sum* of H_1 and H_2 ; it is denoted by $H_1 \oplus H_2$, and we have an isomorphism $H_1 \times H_2 \cong H_1 \oplus H_2$.

The groups $\mathbb{Z}, \mathbb{Q}, \mathbb{R}, \mathbb{C}$, and $M_n(\mathbb{R})$ are more than an abelian groups, they are also commutative rings. Furthermore, \mathbb{Q}, \mathbb{R} , and \mathbb{C} are fields. We now introduce rings and fields.

Definition 2.6. A *ring* is a set A equipped with two operations $+: A \times A \rightarrow A$ (called *addition*) and $*: A \times A \rightarrow A$ (called *multiplication*) having the following properties:

- (R1) A is an abelian group w.r.t. $+$;
- (R2) $*$ is associative and has an identity element $1 \in A$;
- (R3) $*$ is distributive w.r.t. $+$.

The identity element for addition is denoted 0 , and the additive inverse of $a \in A$ is denoted by $-a$. More explicitly, the axioms of a ring are the following equations which hold for all $a, b, c \in A$:

$$a + (b + c) = (a + b) + c \quad (\text{associativity of } +) \quad (2.1)$$

$$a + b = b + a \quad (\text{commutativity of } +) \quad (2.2)$$

$$a + 0 = 0 + a = a \quad (\text{zero}) \quad (2.3)$$

$$a + (-a) = (-a) + a = 0 \quad (\text{additive inverse}) \quad (2.4)$$

$$a * (b * c) = (a * b) * c \quad (\text{associativity of } *) \quad (2.5)$$

$$a * 1 = 1 * a = a \quad (\text{identity for } *) \quad (2.6)$$

$$(a + b) * c = (a * c) + (b * c) \quad (\text{distributivity}) \quad (2.7)$$

$$a * (b + c) = (a * b) + (a * c) \quad (\text{distributivity}) \quad (2.8)$$

The ring A is *commutative* if

$$a * b = b * a$$

for all $a, b \in A$.

From (2.7) and (2.8), we easily obtain

$$a * 0 = 0 * a = 0 \quad (2.9)$$

$$a * (-b) = (-a) * b = -(a * b). \quad (2.10)$$

Note that (2.9) implies that if $1 = 0$, then $a = 0$ for all $a \in A$, and thus, $A = \{0\}$. The ring $A = \{0\}$ is called the *trivial ring*. A ring for which $1 \neq 0$ is called *nontrivial*. The multiplication $a * b$ of two elements $a, b \in A$ is often denoted by ab .

Example 2.2.

1. The additive groups $\mathbb{Z}, \mathbb{Q}, \mathbb{R}, \mathbb{C}$, are commutative rings.
2. The group $\mathbb{R}[X]$ of polynomials in one variable with real coefficients is a ring under multiplication of polynomials. It is a commutative ring.
3. The group of $n \times n$ matrices $M_n(\mathbb{R})$ is a ring under matrix multiplication. However, it is not a commutative ring.
4. The group $\mathcal{C}(]a, b[)$ of continuous functions $f:]a, b[\rightarrow \mathbb{R}$ is a ring under the operation $f \cdot g$ defined such that

$$(f \cdot g)(x) = f(x)g(x)$$

for all $x \in]a, b[$.

When $ab = 0$ with $b \neq 0$, we say that a is a *zero divisor*. A ring A is an *integral domain* (or an *entire ring*) if $0 \neq 1$, A is commutative, and $ab = 0$ implies that $a = 0$ or $b = 0$, for all $a, b \in A$. In other words, an integral domain is a nontrivial commutative ring with no zero divisors besides 0.

Example 2.3.

1. The rings $\mathbb{Z}, \mathbb{Q}, \mathbb{R}, \mathbb{C}$, are integral domains.
2. The ring $\mathbb{R}[X]$ of polynomials in one variable with real coefficients is an integral domain.
- 3.
4. For any positive integer, $p \in \mathbb{N}$, define a relation on \mathbb{Z} , denoted $m \equiv n \pmod{p}$, as follows:

$$m \equiv n \pmod{p} \quad \text{iff} \quad m - n = kp \quad \text{for some } k \in \mathbb{Z}.$$

The reader will easily check that this is an equivalence relation, and, moreover, it is compatible with respect to addition and multiplication, which means that if $m_1 \equiv n_1 \pmod{p}$ and $m_2 \equiv n_2 \pmod{p}$, then $m_1 + m_2 \equiv n_1 + n_2 \pmod{p}$ and $m_1 m_2 \equiv n_1 n_2 \pmod{p}$. Consequently, we can define an addition operation and a multiplication operation of the set of equivalence classes \pmod{p} :

$$[m] + [n] = [m + n]$$

and

$$[m] \cdot [n] = [mn].$$

Again, the reader will easily check that the ring axioms are satisfied, with $[0]$ as zero and $[1]$ as multiplicative unit. The resulting ring is denoted by $\mathbb{Z}/p\mathbb{Z}$.¹ Observe that if p is composite, then this ring has zero-divisors. For example, if $p = 4$, then we have

$$2 \cdot 2 \equiv 0 \pmod{4}.$$

¹The notation \mathbb{Z}_p is sometimes used instead of $\mathbb{Z}/p\mathbb{Z}$ but it clashes with the notation for the *p-adic integers* so we prefer not to use it.

However, the reader should prove that $\mathbb{Z}/p\mathbb{Z}$ is an integral domain if p is prime (in fact, it is a field).

5. The ring of $n \times n$ matrices $M_n(\mathbb{R})$ is not an integral domain. It has zero divisors.

A homomorphism between rings is a mapping preserving addition and multiplication (and 0 and 1).

Definition 2.7. Given two rings A and B , a *homomorphism between A and B* is a function $h: A \rightarrow B$ satisfying the following conditions for all $x, y \in A$:

$$\begin{aligned} h(x + y) &= h(x) + h(y) \\ h(xy) &= h(x)h(y) \\ h(0) &= 0 \\ h(1) &= 1. \end{aligned}$$

Actually, because B is a group under addition, $h(0) = 0$ follows from

$$h(x + y) = h(x) + h(y).$$

Example 2.4.

1. If A is a ring, for any integer $n \in \mathbb{Z}$, for any $a \in A$, we define $n \cdot a$ by

$$n \cdot a = \underbrace{a + \cdots + a}_n$$

if $n \geq 0$ (with $0 \cdot a = 0$) and

$$n \cdot a = -(-n) \cdot a$$

if $n < 0$. Then, the map $h: \mathbb{Z} \rightarrow A$ given by

$$h(n) = n \cdot 1_A$$

is a ring homomorphism (where 1_A is the multiplicative identity of A).

2. Given any real $\lambda \in \mathbb{R}$, the evaluation map $\eta_\lambda: \mathbb{R}[X] \rightarrow \mathbb{R}$ defined by

$$\eta_\lambda(f(X)) = f(\lambda)$$

for every polynomial $f(X) \in \mathbb{R}[X]$ is a ring homomorphism.

A ring homomorphism $h: A \rightarrow B$ is an *isomorphism* iff there is a homomorphism $g: B \rightarrow A$ such that $g \circ f = \text{id}_A$ and $f \circ g = \text{id}_B$. Then, g is unique and denoted by h^{-1} . It is easy to show that a bijective ring homomorphism $h: A \rightarrow B$ is an isomorphism. An isomorphism from a ring to itself is called an *automorphism*.

Given a ring A , a subset A' of A is a *subring* of A if A' is a subgroup of A (under addition), is closed under multiplication, and contains 1. If $h: A \rightarrow B$ is a homomorphism of rings, then for any subring A' , the image $h(A')$ is a subring of B , and for any subring B' of B , the inverse image $h^{-1}(B')$ is a subring of A .

A field is a commutative ring K for which $K - \{0\}$ is a group under multiplication.

Definition 2.8. A set K is a *field* if it is a ring and the following properties hold:

(F1) $0 \neq 1$;

(F2) $K^* = K - \{0\}$ is a group w.r.t. $*$ (i.e., every $a \neq 0$ has an inverse w.r.t. $*$);

(F3) $*$ is commutative.

If $*$ is not commutative but (F1) and (F2) hold, we say that we have a *skew field* (or *noncommutative field*).

Note that we are assuming that the operation $*$ of a field is commutative. This convention is not universally adopted, but since $*$ will be commutative for most fields we will encounter, we may as well include this condition in the definition.

Example 2.5.

1. The rings \mathbb{Q} , \mathbb{R} , and \mathbb{C} are fields.
2. The set of (formal) fractions $f(X)/g(X)$ of polynomials $f(X), g(X) \in \mathbb{R}[X]$, where $g(X)$ is not the null polynomial, is a field.
3. The ring $\mathcal{C}(]a, b[)$ of continuous functions $f:]a, b[\rightarrow \mathbb{R}$ such that $f(x) \neq 0$ for all $x \in]a, b[$ is a field.
4. The ring $\mathbb{Z}/p\mathbb{Z}$ is a field whenever p is prime.

A homomorphism $h: K_1 \rightarrow K_2$ between two fields K_1 and K_2 is just a homomorphism between the rings K_1 and K_2 . However, because K_1^* and K_2^* are groups under multiplication, a homomorphism of fields must be injective.

First, observe that for any $x \neq 0$,

$$1 = h(1) = h(xx^{-1}) = h(x)h(x^{-1})$$

and

$$1 = h(1) = h(x^{-1}x) = h(x^{-1})h(x),$$

so $h(x) \neq 0$ and

$$h(x^{-1}) = h(x)^{-1}.$$

But then, if $h(x) = 0$, we must have $x = 0$. Consequently, h is injective.

A field homomorphism $h: K_1 \rightarrow K_2$ is an *isomorphism* iff there is a homomorphism $g: K_2 \rightarrow K_1$ such that $g \circ h = \text{id}_{K_1}$ and $h \circ g = \text{id}_{K_2}$. Then, g is unique and denoted by h^{-1} . It is easy to show that a bijective field homomorphism $h: K_1 \rightarrow K_2$ is an isomorphism. An isomorphism from a field to itself is called an *automorphism*.

Since every homomorphism $h: K_1 \rightarrow K_2$ between two fields is injective, the image $h(K_1)$ is a subfield of K_2 . We also say that K_2 is an *extension* of K_1 . A field K is said to be *algebraically closed* if every polynomial $p(X)$ with coefficients in K has some root in K ; that is, there is some $a \in K$ such that $p(a) = 0$. It can be shown that every field K has some minimal extension Ω which is algebraically closed, called an *algebraic closure* of K . For example, \mathbb{C} is the algebraic closure of both \mathbb{Q} and \mathbb{R} .

Given a field K and an automorphism $h: K \rightarrow K$ of K , it is easy to check that the set

$$\text{Fix}(h) = \{a \in K \mid h(a) = a\}$$

of elements of K fixed by h is a subfield of K called the *field fixed by h* .

If K is a field, we have the ring homomorphism $h: \mathbb{Z} \rightarrow K$ given by $h(n) = n \cdot 1$. If h is injective, then K contains a copy of \mathbb{Z} , and since it is a field, it contains a copy of \mathbb{Q} . In this case, we say that K has *characteristic 0*. If h is not injective, then $h(\mathbb{Z})$ is a subring of K , and thus an integral domain, which is isomorphic to $\mathbb{Z}/p\mathbb{Z}$ for some $p \geq 1$. But then, p must be prime since $\mathbb{Z}/p\mathbb{Z}$ is an integral domain iff it is a field iff p is prime. The prime p is called the *characteristic* of K , and we also say that K is of *finite characteristic*.

2.2 Vector Spaces

For every $n \geq 1$, let \mathbb{R}^n be the set of n -tuples $x = (x_1, \dots, x_n)$. Addition can be extended to \mathbb{R}^n as follows:

$$(x_1, \dots, x_n) + (y_1, \dots, y_n) = (x_1 + y_1, \dots, x_n + y_n).$$

We can also define an operation $\cdot: \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ as follows:

$$\lambda \cdot (x_1, \dots, x_n) = (\lambda x_1, \dots, \lambda x_n).$$

The resulting algebraic structure has some interesting properties, those of a vector space.

Before defining vector spaces, we need to discuss a strategic choice which, depending how it is settled, may reduce or increase headaches in dealing with notions such as linear combinations and linear dependence (or independence). The issue has to do with using sets of vectors versus sequences of vectors.

Our experience tells us that *it is preferable to use sequences of vectors*; even better, indexed families of vectors. (We are not alone in having opted for sequences over sets, and we are in good company; for example, Artin [4], Axler [7], and Lang [67] use sequences. Nevertheless, some prominent authors such as Lax [71] use sets. We leave it to the reader to conduct a survey on this issue.)

Given a set A , recall that a *sequence* is an ordered n -tuple $(a_1, \dots, a_n) \in A^n$ of elements from A , for some natural number n . The elements of a sequence need not be distinct and the order is important. For example, (a_1, a_2, a_1) and (a_2, a_1, a_1) are two distinct sequences in A^3 . Their underlying set is $\{a_1, a_2\}$.

What we just defined are *finite* sequences, which can also be viewed as functions from $\{1, 2, \dots, n\}$ to the set A ; the i th element of the sequence (a_1, \dots, a_n) is the image of i under the function. This viewpoint is fruitful, because it allows us to define (countably) infinite sequences as functions $s: \mathbb{N} \rightarrow A$. But then, why limit ourselves to ordered sets such as $\{1, \dots, n\}$ or \mathbb{N} as index sets?

The main role of the index set is to tag each element uniquely, and the order of the tags is not crucial, although convenient. Thus, it is natural to define an *I -indexed family* of elements of A , for short a *family*, as a function $a: I \rightarrow A$ where I is any set viewed as an index set. Since the function a is determined by its graph

$$\{(i, a(i)) \mid i \in I\},$$

the family a can be viewed as the set of pairs $a = \{(i, a(i)) \mid i \in I\}$. For notational simplicity, we write a_i instead of $a(i)$, and denote the family $a = \{(i, a(i)) \mid i \in I\}$ by $(a_i)_{i \in I}$. For example, if $I = \{r, g, b, y\}$ and $A = \mathbb{N}$, the set of pairs

$$a = \{(r, 2), (g, 3), (b, 2), (y, 11)\}$$

is an indexed family. The element 2 appears twice in the family with the two distinct tags r and b .

When the indexed set I is totally ordered, a family $(a_i)_{i \in I}$ often called an *I -sequence*. Interestingly, sets can be viewed as special cases of families. Indeed, a set A can be viewed as the A -indexed family $\{(a, a) \mid a \in A\}$ corresponding to the identity function.

Remark: An indexed family should not be confused with a multiset. Given any set A , a *multiset* is similar to a set, except that elements of A may occur more than once. For example, if $A = \{a, b, c, d\}$, then $\{a, a, a, b, c, c, d, d\}$ is a multiset. Each element appears with a certain multiplicity, but the order of the elements does not matter. For example, a has multiplicity 3. Formally, a multiset is a function $s: A \rightarrow \mathbb{N}$, or equivalently a set of pairs $\{(a, i) \mid a \in A\}$. Thus, a multiset is an A -indexed family of elements from \mathbb{N} , but not a \mathbb{N} -indexed family, since distinct elements may have the same multiplicity (such as c and d in the example above). An indexed family is a generalization of a sequence, but a multiset is a generalization of a set.

We also need to take care of an annoying technicality, which is to define sums of the form $\sum_{i \in I} a_i$, where I is any finite index set and $(a_i)_{i \in I}$ is a family of elements in some set A equipped with a binary operation $+: A \times A \rightarrow A$ which is associative (axiom (G1)) and commutative. This will come up when we define linear combinations.

The issue is that the binary operation $+$ only tells us how to compute $a_1 + a_2$ for two elements of A , but it does not tell us what is the sum of three or more elements. For example, how should $a_1 + a_2 + a_3$ be defined?

What we have to do is to define $a_1 + a_2 + a_3$ by using a sequence of steps each involving two elements, and there are two possible ways to do this: $a_1 + (a_2 + a_3)$ and $(a_1 + a_2) + a_3$. If our operation $+$ is not associative, these are different values. If it is associative, then $a_1 + (a_2 + a_3) = (a_1 + a_2) + a_3$, but then there are still six possible permutations of the indices 1, 2, 3, and if $+$ is not commutative, these values are generally different. If our operation is commutative, then all six permutations have the same value. Thus, if $+$ is associative and commutative, it seems intuitively clear that a sum of the form $\sum_{i \in I} a_i$ does not depend on the order of the operations used to compute it.

This is indeed the case, but a rigorous proof requires induction, and such a proof is surprisingly involved. Readers may accept without proof the fact that sums of the form $\sum_{i \in I} a_i$ are indeed well defined, and jump directly to Definition 2.9. For those who want to see the gory details, here we go.

First, we define sums $\sum_{i \in I} a_i$, where I is a finite sequence of distinct natural numbers, say $I = (i_1, \dots, i_m)$. If $I = (i_1, \dots, i_m)$ with $m \geq 2$, we denote the sequence (i_2, \dots, i_m) by $I - \{i_1\}$. We proceed by induction on the size m of I . Let

$$\begin{aligned} \sum_{i \in I} a_i &= a_{i_1}, \quad \text{if } m = 1, \\ \sum_{i \in I} a_i &= a_{i_1} + \left(\sum_{i \in I - \{i_1\}} a_i \right), \quad \text{if } m > 1. \end{aligned}$$

For example, if $I = (1, 2, 3, 4)$, we have

$$\sum_{i \in I} a_i = a_1 + (a_2 + (a_3 + a_4)).$$

If the operation $+$ is not associative, the grouping of the terms matters. For instance, in general

$$a_1 + (a_2 + (a_3 + a_4)) \neq (a_1 + a_2) + (a_3 + a_4).$$

However, if the operation $+$ is associative, the sum $\sum_{i \in I} a_i$ should not depend on the grouping of the elements in I , as long as their order is preserved. For example, if $I = (1, 2, 3, 4, 5)$, $J_1 = (1, 2)$, and $J_2 = (3, 4, 5)$, we expect that

$$\sum_{i \in I} a_i = \left(\sum_{j \in J_1} a_j \right) + \left(\sum_{j \in J_2} a_j \right).$$

This is indeed the case, as we have the following proposition.

Proposition 2.5. *Given any nonempty set A equipped with an associative binary operation $+$: $A \times A \rightarrow A$, for any nonempty finite sequence I of distinct natural numbers and for any partition of I into p nonempty sequences I_{k_1}, \dots, I_{k_p} , for some nonempty sequence $K = (k_1, \dots, k_p)$ of distinct natural numbers such that $k_i < k_j$ implies that $\alpha < \beta$ for all $\alpha \in I_{k_i}$ and all $\beta \in I_{k_j}$, for every sequence $(a_i)_{i \in I}$ of elements in A , we have*

$$\sum_{\alpha \in I} a_\alpha = \sum_{k \in K} \left(\sum_{\alpha \in I_k} a_\alpha \right).$$

Proof. We proceed by induction on the size n of I .

If $n = 1$, then we must have $p = 1$ and $I_{k_1} = I$, so the proposition holds trivially.

Next, assume $n > 1$. If $p = 1$, then $I_{k_1} = I$ and the formula is trivial, so assume that $p \geq 2$ and write $J = (k_2, \dots, k_p)$. There are two cases.

Case 1. The sequence I_{k_1} has a single element, say β , which is the first element of I . In this case, write C for the sequence obtained from I by deleting its first element β . By definition,

$$\sum_{\alpha \in I} a_\alpha = a_\beta + \left(\sum_{\alpha \in C} a_\alpha \right),$$

and

$$\sum_{k \in K} \left(\sum_{\alpha \in I_k} a_\alpha \right) = a_\beta + \left(\sum_{j \in J} \left(\sum_{\alpha \in I_j} a_\alpha \right) \right).$$

Since $|C| = n - 1$, by the induction hypothesis, we have

$$\left(\sum_{\alpha \in C} a_\alpha \right) = \sum_{j \in J} \left(\sum_{\alpha \in I_j} a_\alpha \right),$$

which yields our identity.

Case 2. The sequence I_{k_1} has at least two elements. In this case, let β be the first element of I (and thus of I_{k_1}), let I' be the sequence obtained from I by deleting its first element β , let I'_{k_1} be the sequence obtained from I_{k_1} by deleting its first element β , and let $I'_{k_i} = I_{k_i}$ for $i = 2, \dots, p$. Recall that $J = (k_2, \dots, k_p)$ and $K = (k_1, \dots, k_p)$. The sequence I' has $n - 1$ elements, so by the induction hypothesis applied to I' and the I'_{k_i} , we get

$$\sum_{\alpha \in I'} a_\alpha = \sum_{k \in K} \left(\sum_{\alpha \in I'_k} a_\alpha \right) = \left(\sum_{\alpha \in I'_{k_1}} a_\alpha \right) + \left(\sum_{j \in J} \left(\sum_{\alpha \in I_j} a_\alpha \right) \right).$$

If we add the lefthand side to a_β , by definition we get

$$\sum_{\alpha \in I} a_\alpha.$$

If we add the righthand side to a_β , using associativity and the definition of an indexed sum, we get

$$\begin{aligned} a_\beta + \left(\left(\sum_{\alpha \in I'_{k_1}} a_\alpha \right) + \left(\sum_{j \in J} \left(\sum_{\alpha \in I_j} a_\alpha \right) \right) \right) &= \left(a_\beta + \left(\sum_{\alpha \in I'_{k_1}} a_\alpha \right) \right) + \left(\sum_{j \in J} \left(\sum_{\alpha \in I_j} a_\alpha \right) \right) \\ &= \left(\sum_{\alpha \in I_{k_1}} a_\alpha \right) + \left(\sum_{j \in J} \left(\sum_{\alpha \in I_j} a_\alpha \right) \right) \\ &= \sum_{k \in K} \left(\sum_{\alpha \in I_k} a_\alpha \right), \end{aligned}$$

as claimed. \square

If $I = (1, \dots, n)$, we also write $\sum_{i=1}^n a_i$ instead of $\sum_{i \in I} a_i$. Since $+$ is associative, Proposition 2.5 shows that the sum $\sum_{i=1}^n a_i$ is independent of the grouping of its elements, which justifies the use of the notation $a_1 + \dots + a_n$ (without any parentheses).

If we also assume that our associative binary operation on A is commutative, then we can show that the sum $\sum_{i \in I} a_i$ does not depend on the ordering of the index set I .

Proposition 2.6. *Given any nonempty set A equipped with an associative and commutative binary operation $+: A \times A \rightarrow A$, for any two nonempty finite sequences I and J of distinct natural numbers such that J is a permutation of I (in other words, the underlying sets of I and J are identical), for every sequence $(a_i)_{i \in I}$ of elements in A , we have*

$$\sum_{\alpha \in I} a_\alpha = \sum_{\alpha \in J} a_\alpha.$$

Proof. We proceed by induction on the number p of elements in I . If $p = 1$, we have $I = J$ and the proposition holds trivially.

If $p > 1$, to simplify notation, assume that $I = (1, \dots, p)$ and that J is a permutation (i_1, \dots, i_p) of I . First, assume that $2 \leq i_1 \leq p-1$, let J' be the sequence obtained from J by deleting i_1 , I' be the sequence obtained from I by deleting i_1 , and let $P = (1, 2, \dots, i_1-1)$ and $Q = (i_1+1, \dots, p-1, p)$. Observe that the sequence I' is the concatenation of the sequences P and Q . By the induction hypothesis applied to J' and I' , and then by Proposition 2.5 applied to I' and its partition (P, Q) , we have

$$\sum_{\alpha \in J'} a_\alpha = \sum_{\alpha \in I'} a_\alpha = \left(\sum_{i=1}^{i_1-1} a_i \right) + \left(\sum_{i=i_1+1}^p a_i \right).$$

If we add the lefthand side to a_{i_1} , by definition we get

$$\sum_{\alpha \in J} a_\alpha.$$

If we add the righthand side to a_{i_1} , we get

$$a_{i_1} + \left(\left(\sum_{i=1}^{i_1-1} a_i \right) + \left(\sum_{i=i_1+1}^p a_i \right) \right).$$

Using associativity, we get

$$a_{i_1} + \left(\left(\sum_{i=1}^{i_1-1} a_i \right) + \left(\sum_{i=i_1+1}^p a_i \right) \right) = \left(a_{i_1} + \left(\sum_{i=1}^{i_1-1} a_i \right) \right) + \left(\sum_{i=i_1+1}^p a_i \right),$$

then using associativity and commutativity several times (more rigorously, using induction on $i_1 - 1$), we get

$$\begin{aligned} \left(a_{i_1} + \left(\sum_{i=1}^{i_1-1} a_i \right) \right) + \left(\sum_{i=i_1+1}^p a_i \right) &= \left(\sum_{i=1}^{i_1-1} a_i \right) + a_{i_1} + \left(\sum_{i=i_1+1}^p a_i \right) \\ &= \sum_{i=1}^p a_i, \end{aligned}$$

as claimed.

The cases where $i_1 = 1$ or $i_1 = p$ are treated similarly, but in a simpler manner since either $P = ()$ or $Q = ()$ (where $()$ denotes the empty sequence). \square

Having done all this, we can now make sense of sums of the form $\sum_{i \in I} a_i$, for any finite indexed set I and any family $a = (a_i)_{i \in I}$ of elements in A , where A is a set equipped with a binary operation $+$ which is associative and commutative.

Indeed, since I is finite, it is in bijection with the set $\{1, \dots, n\}$ for some $n \in \mathbb{N}$, and any total ordering \preceq on I corresponds to a permutation I_{\preceq} of $\{1, \dots, n\}$ (where we identify a permutation with its image). For any total ordering \preceq on I , we define $\sum_{i \in I, \preceq} a_i$ as

$$\sum_{i \in I, \preceq} a_i = \sum_{j \in I_{\preceq}} a_j.$$

Then, for any other total ordering \preceq' on I , we have

$$\sum_{i \in I, \preceq'} a_i = \sum_{j \in I_{\preceq'}} a_j,$$

and since I_{\preceq} and $I_{\preceq'}$ are different permutations of $\{1, \dots, n\}$, by Proposition 2.6, we have

$$\sum_{j \in I_{\preceq}} a_j = \sum_{j \in I_{\preceq'}} a_j.$$

Therefore, the sum $\sum_{i \in I, \preceq} a_i$ does not depend on the total ordering on I . We define *the* sum $\sum_{i \in I} a_i$ as the common value $\sum_{i \in I, \preceq} a_i$ for all total orderings \preceq of I .

Vector spaces are defined as follows.

Definition 2.9. Given a field K (with addition $+$ and multiplication $*$), a *vector space over K* (or *K -vector space*) is a set E (of vectors) together with two operations $+: E \times E \rightarrow E$ (called *vector addition*),² and $\cdot: K \times E \rightarrow E$ (called *scalar multiplication*) satisfying the following conditions for all $\alpha, \beta \in K$ and all $u, v \in E$;

(V0) E is an abelian group w.r.t. $+$, with identity element 0 ;³

(V1) $\alpha \cdot (u + v) = (\alpha \cdot u) + (\alpha \cdot v)$;

(V2) $(\alpha + \beta) \cdot u = (\alpha \cdot u) + (\beta \cdot u)$;

(V3) $(\alpha * \beta) \cdot u = \alpha \cdot (\beta \cdot u)$;

(V4) $1 \cdot u = u$.

In (V3), $*$ denotes multiplication in the field K .

Given $\alpha \in K$ and $v \in E$, the element $\alpha \cdot v$ is also denoted by αv . The field K is often called the field of scalars.

Unless specified otherwise or unless we are dealing with several different fields, in the rest of this chapter, we assume that all K -vector spaces are defined with respect to a fixed field K . Thus, we will refer to a K -vector space simply as a vector space. In most cases, the field K will be the field \mathbb{R} of reals.

From (V0), a vector space always contains the null vector 0 , and thus is nonempty. From (V1), we get $\alpha \cdot 0 = 0$, and $\alpha \cdot (-v) = -(\alpha \cdot v)$. From (V2), we get $0 \cdot v = 0$, and $(-\alpha) \cdot v = -(\alpha \cdot v)$.

Another important consequence of the axioms is the following fact: For any $u \in E$ and any $\lambda \in K$, if $\lambda \neq 0$ and $\lambda \cdot u = 0$, then $u = 0$.

Indeed, since $\lambda \neq 0$, it has a multiplicative inverse λ^{-1} , so from $\lambda \cdot u = 0$, we get

$$\lambda^{-1} \cdot (\lambda \cdot u) = \lambda^{-1} \cdot 0.$$

However, we just observed that $\lambda^{-1} \cdot 0 = 0$, and from (V3) and (V4), we have

$$\lambda^{-1} \cdot (\lambda \cdot u) = (\lambda^{-1} \lambda) \cdot u = 1 \cdot u = u,$$

and we deduce that $u = 0$.

Remark: One may wonder whether axiom (V4) is really needed. Could it be derived from the other axioms? The answer is **no**. For example, one can take $E = \mathbb{R}^n$ and define $\cdot: \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ by

$$\lambda \cdot (x_1, \dots, x_n) = (0, \dots, 0)$$

²The symbol $+$ is overloaded, since it denotes both addition in the field K and addition of vectors in E . It is usually clear from the context which $+$ is intended.

³The symbol 0 is also overloaded, since it represents both the zero in K (a scalar) and the identity element of E (the zero vector). Confusion rarely arises, but one may prefer using $\mathbf{0}$ for the zero vector.

for all $(x_1, \dots, x_n) \in \mathbb{R}^n$ and all $\lambda \in \mathbb{R}$. Axioms (V0)–(V3) are all satisfied, but (V4) fails. Less trivial examples can be given using the notion of a basis, which has not been defined yet.

The field K itself can be viewed as a vector space over itself, addition of vectors being addition in the field, and multiplication by a scalar being multiplication in the field.

Example 2.6.

1. The fields \mathbb{R} and \mathbb{C} are vector spaces over \mathbb{R} .
2. The groups \mathbb{R}^n and \mathbb{C}^n are vector spaces over \mathbb{R} , and \mathbb{C}^n is a vector space over \mathbb{C} .
3. The ring $\mathbb{R}[X]$ of polynomials is a vector space over \mathbb{R} , and $\mathbb{C}[X]$ is a vector space over \mathbb{R} and \mathbb{C} . The ring of $n \times n$ matrices $M_n(\mathbb{R})$ is a vector space over \mathbb{R} .
4. The ring $\mathcal{C}(]a, b[)$ of continuous functions $f:]a, b[\rightarrow \mathbb{R}$ is a vector space over \mathbb{R} .

Let E be a vector space. We would like to define the important notions of linear combination and linear independence. These notions can be defined for sets of vectors in E , but it will turn out to be more convenient to define them for families $(v_i)_{i \in I}$, where I is any arbitrary index set.

2.3 Linear Independence, Subspaces

One of the most useful properties of vector spaces is that they possess bases. What this means is that in every vector space, E , there is some set of vectors, $\{e_1, \dots, e_n\}$, such that *every*, vector, $v \in E$, can be written as a linear combination,

$$v = \lambda_1 e_1 + \dots + \lambda_n e_n,$$

of the e_i , for some scalars, $\lambda_1, \dots, \lambda_n \in K$. Furthermore, the n -tuple, $(\lambda_1, \dots, \lambda_n)$, as above is unique.

This description is fine when E has a finite basis, $\{e_1, \dots, e_n\}$, but this is not always the case! For example, the vector space of real polynomials, $\mathbb{R}[X]$, does not have a finite basis but instead it has an infinite basis, namely

$$1, X, X^2, \dots, X^n, \dots$$

One might wonder if it is possible for a vector space to have bases of different sizes, or even to have a finite basis as well as an infinite basis. We will see later on that this is not possible; all bases of a vector space have the same number of elements (cardinality), which is called the *dimension* of the space. However, we have the following problem: If a vector space has

an infinite basis, $\{e_1, e_2, \dots\}$, how do we define linear combinations? Do we allow linear combinations

$$\lambda_1 e_1 + \lambda_2 e_2 + \dots$$

with infinitely many nonzero coefficients?

If we allow linear combinations with infinitely many nonzero coefficients, then we have to make sense of these sums and this can only be done reasonably if we define such a sum as the limit of the sequence of vectors, $s_1, s_2, \dots, s_n, \dots$, with $s_1 = \lambda_1 e_1$ and

$$s_{n+1} = s_n + \lambda_{n+1} e_{n+1}.$$

But then, how do we define such limits? Well, we have to define some topology on our space, by means of a norm, a metric or some other mechanism. This can indeed be done and this is what Banach spaces and Hilbert spaces are all about but this seems to require a lot of machinery.

A way to avoid limits is to restrict our attention to linear combinations involving only *finitely many* vectors. We may have an infinite supply of vectors but we only form linear combinations involving finitely many nonzero coefficients. Technically, this can be done by introducing *families of finite support*. This gives us the ability to manipulate families of scalars indexed by some fixed infinite set and yet to treat these families as if they were finite.

With these motivations in mind, given a set A , recall that an I -indexed family $(a_i)_{i \in I}$ of elements of A (for short, a *family*) is a function $a: I \rightarrow A$, or equivalently a set of pairs $\{(i, a_i) \mid i \in I\}$. We agree that when $I = \emptyset$, $(a_i)_{i \in I} = \emptyset$. A family $(a_i)_{i \in I}$ is finite if I is finite.

Remark: When considering a family $(a_i)_{i \in I}$, there is no reason to assume that I is ordered. The crucial point is that every element of the family is uniquely indexed by an element of I . Thus, unless specified otherwise, we do not assume that the elements of an index set are ordered.

If A is an abelian group (usually, when A is a ring or a vector space) with identity 0, we say that a family $(a_i)_{i \in I}$ has *finite support* if $a_i = 0$ for all $i \in I - J$, where J is a finite subset of I (the support of the family).

Given two disjoint sets I and J , the union of two families $(u_i)_{i \in I}$ and $(v_j)_{j \in J}$, denoted as $(u_i)_{i \in I} \cup (v_j)_{j \in J}$, is the family $(w_k)_{k \in (I \cup J)}$ defined such that $w_k = u_k$ if $k \in I$, and $w_k = v_k$ if $k \in J$. Given a family $(u_i)_{i \in I}$ and any element v , we denote by $(u_i)_{i \in I} \cup_k (v)$ the family $(w_i)_{i \in I \cup \{k\}}$ defined such that, $w_i = u_i$ if $i \in I$, and $w_k = v$, where k is any index such that $k \notin I$. Given a family $(u_i)_{i \in I}$, a subfamily of $(u_i)_{i \in I}$ is a family $(u_j)_{j \in J}$ where J is any subset of I .

In this chapter, unless specified otherwise, it is assumed that all families of scalars have finite support.

Definition 2.10. Let E be a vector space. A vector $v \in E$ is a *linear combination* of a family $(u_i)_{i \in I}$ of elements of E if there is a family $(\lambda_i)_{i \in I}$ of scalars in K such that

$$v = \sum_{i \in I} \lambda_i u_i.$$

When $I = \emptyset$, we stipulate that $v = 0$. (By proposition 2.6, sums of the form $\sum_{i \in I} \lambda_i u_i$ are well defined.) We say that a family $(u_i)_{i \in I}$ is *linearly independent* if for every family $(\lambda_i)_{i \in I}$ of scalars in K ,

$$\sum_{i \in I} \lambda_i u_i = 0 \quad \text{implies that} \quad \lambda_i = 0 \text{ for all } i \in I.$$

Equivalently, a family $(u_i)_{i \in I}$ is *linearly dependent* if there is some family $(\lambda_i)_{i \in I}$ of scalars in K such that

$$\sum_{i \in I} \lambda_i u_i = 0 \quad \text{and} \quad \lambda_j \neq 0 \text{ for some } j \in I.$$

We agree that when $I = \emptyset$, the family \emptyset is linearly independent.

Observe that defining linear combinations for families of vectors rather than for sets of vectors has the advantage that the vectors being combined need not be distinct. For example, for $I = \{1, 2, 3\}$ and the families (u, v, u) and $(\lambda_1, \lambda_2, \lambda_1)$, the linear combination

$$\sum_{i \in I} \lambda_i u_i = \lambda_1 u + \lambda_2 v + \lambda_1 u$$

makes sense. Using sets of vectors in the definition of a linear combination does not allow such linear combinations; this is too restrictive.

Unravelling Definition 2.10, a family $(u_i)_{i \in I}$ is linearly dependent iff some u_j in the family can be expressed as a linear combination of the other vectors in the family. Indeed, there is some family $(\lambda_i)_{i \in I}$ of scalars in K such that

$$\sum_{i \in I} \lambda_i u_i = 0 \quad \text{and} \quad \lambda_j \neq 0 \text{ for some } j \in I,$$

which implies that

$$u_j = \sum_{i \in (I - \{j\})} -\lambda_j^{-1} \lambda_i u_i.$$

Observe that one of the reasons for defining linear dependence for families of vectors rather than for sets of vectors is that our definition allows multiple occurrences of a vector. This is important because a matrix may contain identical columns, and we would like to say that these columns are linearly dependent. The definition of linear dependence for sets does not allow us to do that.

The above also shows that a family $(u_i)_{i \in I}$ is linearly independent iff either $I = \emptyset$, or I consists of a single element i and $u_i \neq 0$, or $|I| \geq 2$ and no vector u_j in the family can be expressed as a linear combination of the other vectors in the family.

When I is nonempty, if the family $(u_i)_{i \in I}$ is linearly independent, note that $u_i \neq 0$ for all $i \in I$. Otherwise, if $u_i = 0$ for some $i \in I$, then we get a nontrivial linear dependence $\sum_{i \in I} \lambda_i u_i = 0$ by picking any nonzero λ_i and letting $\lambda_k = 0$ for all $k \in I$ with $k \neq i$, since $\lambda_i 0 = 0$. If $|I| \geq 2$, we must also have $u_i \neq u_j$ for all $i, j \in I$ with $i \neq j$, since otherwise we get a nontrivial linear dependence by picking $\lambda_i = \lambda$ and $\lambda_j = -\lambda$ for any nonzero λ , and letting $\lambda_k = 0$ for all $k \in I$ with $k \neq i, j$.

Thus, the definition of linear independence implies that a nontrivial linearly independent family is actually a set. This explains why certain authors choose to define linear independence for sets of vectors. The problem with this approach is that linear dependence, which is the logical negation of linear independence, is then only defined for sets of vectors. However, as we pointed out earlier, it is really desirable to define linear dependence for families allowing multiple occurrences of the same vector.

Example 2.7.

1. Any two distinct scalars $\lambda, \mu \neq 0$ in K are linearly dependent.
2. In \mathbb{R}^3 , the vectors $(1, 0, 0)$, $(0, 1, 0)$, and $(0, 0, 1)$ are linearly independent.
3. In \mathbb{R}^4 , the vectors $(1, 1, 1, 1)$, $(0, 1, 1, 1)$, $(0, 0, 1, 1)$, and $(0, 0, 0, 1)$ are linearly independent.
4. In \mathbb{R}^2 , the vectors $u = (1, 1)$, $v = (0, 1)$ and $w = (2, 3)$ are linearly dependent, since

$$w = 2u + v.$$

Note that a family $(u_i)_{i \in I}$ is linearly independent iff $(u_j)_{j \in J}$ is linearly independent for every finite subset J of I (even when $I = \emptyset$). Indeed, when $\sum_{i \in I} \lambda_i u_i = 0$, the family $(\lambda_i)_{i \in I}$ of scalars in K has finite support, and thus $\sum_{i \in I} \lambda_i u_i = 0$ really means that $\sum_{j \in J} \lambda_j u_j = 0$ for a finite subset J of I . When I is finite, we often assume that it is the set $I = \{1, 2, \dots, n\}$. In this case, we denote the family $(u_i)_{i \in I}$ as (u_1, \dots, u_n) .

The notion of a subspace of a vector space is defined as follows.

Definition 2.11. Given a vector space E , a subset F of E is a *linear subspace* (or *subspace*) of E if F is nonempty and $\lambda u + \mu v \in F$ for all $u, v \in F$, and all $\lambda, \mu \in K$.

It is easy to see that a subspace F of E is indeed a vector space, since the restriction of $+: E \times E \rightarrow E$ to $F \times F$ is indeed a function $+: F \times F \rightarrow F$, and the restriction of $\cdot: K \times E \rightarrow E$ to $K \times F$ is indeed a function $\cdot: K \times F \rightarrow F$.

It is also easy to see that any intersection of subspaces is a subspace. Since F is nonempty, if we pick any vector $u \in F$ and if we let $\lambda = \mu = 0$, then $\lambda u + \mu u = 0u + 0u = 0$, so every subspace contains the vector 0. For any nonempty finite index set I , one can show by induction on the cardinality of I that if $(u_i)_{i \in I}$ is any family of vectors $u_i \in F$ and $(\lambda_i)_{i \in I}$ is any family of scalars, then $\sum_{i \in I} \lambda_i u_i \in F$.

The subspace $\{0\}$ will be denoted by (0) , or even 0 (with a mild abuse of notation).

Example 2.8.

1. In \mathbb{R}^2 , the set of vectors $u = (x, y)$ such that

$$x + y = 0$$

is a subspace.

2. In \mathbb{R}^3 , the set of vectors $u = (x, y, z)$ such that

$$x + y + z = 0$$

is a subspace.

3. For any $n \geq 0$, the set of polynomials $f(X) \in \mathbb{R}[X]$ of degree at most n is a subspace of $\mathbb{R}[X]$.
4. The set of upper triangular $n \times n$ matrices is a subspace of the space of $n \times n$ matrices.

Proposition 2.7. *Given any vector space E , if S is any nonempty subset of E , then the smallest subspace $\langle S \rangle$ (or $\text{Span}(S)$) of E containing S is the set of all (finite) linear combinations of elements from S .*

Proof. We prove that the set $\text{Span}(S)$ of all linear combinations of elements of S is a subspace of E , leaving as an exercise the verification that every subspace containing S also contains $\text{Span}(S)$.

First, $\text{Span}(S)$ is nonempty since it contains S (which is nonempty). If $u = \sum_{i \in I} \lambda_i u_i$ and $v = \sum_{j \in J} \mu_j v_j$ are any two linear combinations in $\text{Span}(S)$, for any two scalars $\lambda, \mu \in \mathbb{R}$,

$$\begin{aligned} \lambda u + \mu v &= \lambda \sum_{i \in I} \lambda_i u_i + \mu \sum_{j \in J} \mu_j v_j \\ &= \sum_{i \in I} \lambda \lambda_i u_i + \sum_{j \in J} \mu \mu_j v_j \\ &= \sum_{i \in I-J} \lambda \lambda_i u_i + \sum_{i \in I \cap J} (\lambda \lambda_i + \mu \mu_i) u_i + \sum_{j \in J-I} \mu \mu_j v_j, \end{aligned}$$

which is a linear combination with index set $I \cup J$, and thus $\lambda u + \mu v \in \text{Span}(S)$, which proves that $\text{Span}(S)$ is a subspace. \square

One might wonder what happens if we add extra conditions to the coefficients involved in forming linear combinations. Here are three natural restrictions which turn out to be important (as usual, we assume that our index sets are finite):

- (1) Consider combinations $\sum_{i \in I} \lambda_i u_i$ for which

$$\sum_{i \in I} \lambda_i = 1.$$

These are called *affine combinations*. One should realize that every linear combination $\sum_{i \in I} \lambda_i u_i$ can be viewed as an affine combination. For example, if k is an index not in I , if we let $J = I \cup \{k\}$, $u_k = 0$, and $\lambda_k = 1 - \sum_{i \in I} \lambda_i$, then $\sum_{j \in J} \lambda_j u_j$ is an affine combination and

$$\sum_{i \in I} \lambda_i u_i = \sum_{j \in J} \lambda_j u_j.$$

However, we get new spaces. For example, in \mathbb{R}^3 , the set of all affine combinations of the three vectors $e_1 = (1, 0, 0)$, $e_2 = (0, 1, 0)$, and $e_3 = (0, 0, 1)$, is the plane passing through these three points. Since it does not contain $0 = (0, 0, 0)$, it is not a linear subspace.

- (2) Consider combinations $\sum_{i \in I} \lambda_i u_i$ for which

$$\lambda_i \geq 0, \quad \text{for all } i \in I.$$

These are called *positive* (or *conic*) *combinations*. It turns out that positive combinations of families of vectors are *cones*. They show naturally in convex optimization.

- (3) Consider combinations $\sum_{i \in I} \lambda_i u_i$ for which we require (1) *and* (2), that is

$$\sum_{i \in I} \lambda_i = 1, \quad \text{and} \quad \lambda_i \geq 0 \quad \text{for all } i \in I.$$

These are called *convex combinations*. Given any finite family of vectors, the set of all convex combinations of these vectors is a *convex polyhedron*. Convex polyhedra play a very important role in convex optimization.

2.4 Bases of a Vector Space

Given a vector space E , given a family $(v_i)_{i \in I}$, the subset V of E consisting of the null vector 0 and of all linear combinations of $(v_i)_{i \in I}$ is easily seen to be a subspace of E . Subspaces having such a “generating family” play an important role, and motivate the following definition.

Definition 2.12. Given a vector space E and a subspace V of E , a family $(v_i)_{i \in I}$ of vectors $v_i \in V$ *spans* V or *generates* V if for every $v \in V$, there is some family $(\lambda_i)_{i \in I}$ of scalars in K such that

$$v = \sum_{i \in I} \lambda_i v_i.$$

We also say that the elements of $(v_i)_{i \in I}$ are *generators* of V and that V is *spanned by* $(v_i)_{i \in I}$, or *generated by* $(v_i)_{i \in I}$. If a subspace V of E is generated by a finite family $(v_i)_{i \in I}$, we say that V is *finitely generated*. A family $(u_i)_{i \in I}$ that spans V and is linearly independent is called a *basis* of V .

Example 2.9.

1. In \mathbb{R}^3 , the vectors $(1, 0, 0)$, $(0, 1, 0)$, and $(0, 0, 1)$ form a basis.
2. The vectors $(1, 1, 1, 1)$, $(1, 1, -1, -1)$, $(1, -1, 0, 0)$, $(0, 0, 1, -1)$ form a basis of \mathbb{R}^4 known as the *Haar basis*. This basis and its generalization to dimension 2^n are crucial in wavelet theory.
3. In the subspace of polynomials in $\mathbb{R}[X]$ of degree at most n , the polynomials $1, X, X^2, \dots, X^n$ form a basis.
4. The *Bernstein polynomials* $\binom{n}{k} (1 - X)^{n-k} X^k$ for $k = 0, \dots, n$, also form a basis of that space. These polynomials play a major role in the theory of *spline curves*.

It is a standard result of linear algebra that every vector space E has a basis, and that for any two bases $(u_i)_{i \in I}$ and $(v_j)_{j \in J}$, I and J have the same cardinality. In particular, if E has a finite basis of n elements, every basis of E has n elements, and the integer n is called the *dimension* of the vector space E . We begin with a crucial lemma.

Lemma 2.8. *Given a linearly independent family $(u_i)_{i \in I}$ of elements of a vector space E , if $v \in E$ is not a linear combination of $(u_i)_{i \in I}$, then the family $(u_i)_{i \in I} \cup_k (v)$ obtained by adding v to the family $(u_i)_{i \in I}$ is linearly independent (where $k \notin I$).*

Proof. Assume that $\mu v + \sum_{i \in I} \lambda_i u_i = 0$, for any family $(\lambda_i)_{i \in I}$ of scalars in K . If $\mu \neq 0$, then μ has an inverse (because K is a field), and thus we have $v = -\sum_{i \in I} (\mu^{-1} \lambda_i) u_i$, showing that v is a linear combination of $(u_i)_{i \in I}$ and contradicting the hypothesis. Thus, $\mu = 0$. But then, we have $\sum_{i \in I} \lambda_i u_i = 0$, and since the family $(u_i)_{i \in I}$ is linearly independent, we have $\lambda_i = 0$ for all $i \in I$. \square

The next theorem holds in general, but the proof is more sophisticated for vector spaces that do not have a finite set of generators. Thus, in this chapter, we only prove the theorem for finitely generated vector spaces.

Theorem 2.9. *Given any finite family $S = (u_i)_{i \in I}$ generating a vector space E and any linearly independent subfamily $L = (u_j)_{j \in J}$ of S (where $J \subseteq I$), there is a basis B of E such that $L \subseteq B \subseteq S$.*

Proof. Consider the set of linearly independent families B such that $L \subseteq B \subseteq S$. Since this set is nonempty and finite, it has some maximal element, say $B = (u_h)_{h \in H}$. We claim that B generates E . Indeed, if B does not generate E , then there is some $u_p \in S$ that is not a linear combination of vectors in B (since S generates E), with $p \notin H$. Then, by Lemma 2.8, the family $B' = (u_h)_{h \in H \cup \{p\}}$ is linearly independent, and since $L \subseteq B \subset B' \subseteq S$, this contradicts the maximality of B . Thus, B is a basis of E such that $L \subseteq B \subseteq S$. \square

Remark: Theorem 2.9 also holds for vector spaces that are not finitely generated. In this case, the problem is to guarantee the existence of a maximal linearly independent family B such that $L \subseteq B \subseteq S$. The existence of such a maximal family can be shown using Zorn's lemma, see Appendix 31 and the references given there.

A situation where the full generality of Theorem 2.9 is needed is the case of the vector space \mathbb{R} over the field of coefficients \mathbb{Q} . The numbers 1 and $\sqrt{2}$ are linearly independent over \mathbb{Q} , so according to Theorem 2.9, the linearly independent family $L = (1, \sqrt{2})$ can be extended to a basis B of \mathbb{R} . Since \mathbb{R} is uncountable and \mathbb{Q} is countable, such a basis must be uncountable!

Let $(v_i)_{i \in I}$ be a family of vectors in E . We say that $(v_i)_{i \in I}$ is a *maximal linearly independent family* of E if it is linearly independent, and if for any vector $w \in E$, the family $(v_i)_{i \in I} \cup \{w\}$ obtained by adding w to the family $(v_i)_{i \in I}$ is linearly dependent. We say that $(v_i)_{i \in I}$ is a *minimal generating family* of E if it spans E , and if for any index $p \in I$, the family $(v_i)_{i \in I - \{p\}}$ obtained by removing v_p from the family $(v_i)_{i \in I}$ does not span E .

The following proposition giving useful properties characterizing a basis is an immediate consequence of Theorem 2.9.

Proposition 2.10. *Given a vector space E , for any family $B = (v_i)_{i \in I}$ of vectors of E , the following properties are equivalent:*

- (1) B is a basis of E .
- (2) B is a maximal linearly independent family of E .
- (3) B is a minimal generating family of E .

Proof. We prove the equivalence of (1) and (2), leaving the equivalence of (1) and (3) as an exercise.

Assume (1). We claim that B is a maximal linearly independent family. If B is not a maximal linearly independent family, then there is some vector $w \in E$ such that the family B' obtained by adding w to B is linearly independent. However, since B is a basis of E , the

vector w can be expressed as a linear combination of vectors in B , contradicting the fact that B' is linearly independent.

Conversely, assume (2). We claim that B spans E . If B does not span E , then there is some vector $w \in E$ which is not a linear combination of vectors in B . By Lemma 2.8, the family B' obtained by adding w to B is linearly independent. Since B is a proper subfamily of B' , this contradicts the assumption that B is a maximal linearly independent family. Therefore, B must span E , and since B is also linearly independent, it is a basis of E . \square

The following *replacement lemma* due to Steinitz shows the relationship between finite linearly independent families and finite families of generators of a vector space. We begin with a version of the lemma which is a bit informal, but easier to understand than the precise and more formal formulation given in Proposition 2.12. The technical difficulty has to do with the fact that some of the indices need to be renamed.

Proposition 2.11. (*Replacement lemma, version 1*) *Given a vector space E , let (u_1, \dots, u_m) be any finite linearly independent family in E , and let (v_1, \dots, v_n) be any finite family such that every u_i is a linear combination of (v_1, \dots, v_n) . Then, we must have $m \leq n$, and m of the vectors v_j can be replaced by (u_1, \dots, u_m) , such that after renaming some of the indices of the v s, the families $(u_1, \dots, u_m, v_{m+1}, \dots, v_n)$ and (v_1, \dots, v_n) generate the same subspace of E .*

Proof. We proceed by induction on m . When $m = 0$, the family (u_1, \dots, u_m) is empty, and the proposition holds trivially. For the induction step, we have a linearly independent family $(u_1, \dots, u_m, u_{m+1})$. Consider the linearly independent family (u_1, \dots, u_m) . By the induction hypothesis, $m \leq n$, and m of the vectors v_j can be replaced by (u_1, \dots, u_m) , such that after renaming some of the indices of the v s, the families $(u_1, \dots, u_m, v_{m+1}, \dots, v_n)$ and (v_1, \dots, v_n) generate the same subspace of E . The vector u_{m+1} can also be expressed as a linear combination of (v_1, \dots, v_n) , and since $(u_1, \dots, u_m, v_{m+1}, \dots, v_n)$ and (v_1, \dots, v_n) generate the same subspace, u_{m+1} can be expressed as a linear combination of $(u_1, \dots, u_m, v_{m+1}, \dots, v_n)$, say

$$u_{m+1} = \sum_{i=1}^m \lambda_i u_i + \sum_{j=m+1}^n \lambda_j v_j.$$

We claim that $\lambda_j \neq 0$ for some j with $m+1 \leq j \leq n$, which implies that $m+1 \leq n$.

Otherwise, we would have

$$u_{m+1} = \sum_{i=1}^m \lambda_i u_i,$$

a nontrivial linear dependence of the u_i , which is impossible since (u_1, \dots, u_{m+1}) are linearly independent.

Therefore $m + 1 \leq n$, and after renaming indices if necessary, we may assume that $\lambda_{m+1} \neq 0$, so we get

$$v_{m+1} = -\sum_{i=1}^m (\lambda_{m+1}^{-1} \lambda_i) u_i - \lambda_{m+1}^{-1} u_{m+1} - \sum_{j=m+2}^n (\lambda_{m+1}^{-1} \lambda_j) v_j.$$

Observe that the families $(u_1, \dots, u_m, v_{m+1}, \dots, v_n)$ and $(u_1, \dots, u_{m+1}, v_{m+2}, \dots, v_n)$ generate the same subspace, since u_{m+1} is a linear combination of $(u_1, \dots, u_m, v_{m+1}, \dots, v_n)$ and v_{m+1} is a linear combination of $(u_1, \dots, u_{m+1}, v_{m+2}, \dots, v_n)$. Since $(u_1, \dots, u_m, v_{m+1}, \dots, v_n)$ and (v_1, \dots, v_n) generate the same subspace, we conclude that $(u_1, \dots, u_{m+1}, v_{m+2}, \dots, v_n)$ and (v_1, \dots, v_n) generate the same subspace, which concludes the induction hypothesis. \square

For the sake of completeness, here is a more formal statement of the replacement lemma (and its proof).

Proposition 2.12. (*Replacement lemma, version 2*) *Given a vector space E , let $(u_i)_{i \in I}$ be any finite linearly independent family in E , where $|I| = m$, and let $(v_j)_{j \in J}$ be any finite family such that every u_i is a linear combination of $(v_j)_{j \in J}$, where $|J| = n$. Then, there exists a set L and an injection $\rho: L \rightarrow J$ (a relabeling function) such that $L \cap I = \emptyset$, $|L| = n - m$, and the families $(u_i)_{i \in I} \cup (v_{\rho(l)})_{l \in L}$ and $(v_j)_{j \in J}$ generate the same subspace of E . In particular, $m \leq n$.*

Proof. We proceed by induction on $|I| = m$. When $m = 0$, the family $(u_i)_{i \in I}$ is empty, and the proposition holds trivially with $L = J$ (ρ is the identity). Assume $|I| = m + 1$. Consider the linearly independent family $(u_i)_{i \in (I - \{p\})}$, where p is any member of I . By the induction hypothesis, there exists a set L and an injection $\rho: L \rightarrow J$ such that $L \cap (I - \{p\}) = \emptyset$, $|L| = n - m$, and the families $(u_i)_{i \in (I - \{p\})} \cup (v_{\rho(l)})_{l \in L}$ and $(v_j)_{j \in J}$ generate the same subspace of E . If $p \in L$, we can replace L by $(L - \{p\}) \cup \{p'\}$ where p' does not belong to $I \cup L$, and replace ρ by the injection ρ' which agrees with ρ on $L - \{p\}$ and such that $\rho'(p') = \rho(p)$. Thus, we can always assume that $L \cap I = \emptyset$. Since u_p is a linear combination of $(v_j)_{j \in J}$ and the families $(u_i)_{i \in (I - \{p\})} \cup (v_{\rho(l)})_{l \in L}$ and $(v_j)_{j \in J}$ generate the same subspace of E , u_p is a linear combination of $(u_i)_{i \in (I - \{p\})} \cup (v_{\rho(l)})_{l \in L}$. Let

$$u_p = \sum_{i \in (I - \{p\})} \lambda_i u_i + \sum_{l \in L} \lambda_l v_{\rho(l)}. \quad (1)$$

If $\lambda_l = 0$ for all $l \in L$, we have

$$\sum_{i \in (I - \{p\})} \lambda_i u_i - u_p = 0,$$

contradicting the fact that $(u_i)_{i \in I}$ is linearly independent. Thus, $\lambda_l \neq 0$ for some $l \in L$, say $l = q$. Since $\lambda_q \neq 0$, we have

$$v_{\rho(q)} = \sum_{i \in (I - \{p\})} (-\lambda_q^{-1} \lambda_i) u_i + \lambda_q^{-1} u_p + \sum_{l \in (L - \{q\})} (-\lambda_q^{-1} \lambda_l) v_{\rho(l)}. \quad (2)$$

We claim that the families $(u_i)_{i \in (I - \{p\})} \cup (v_{\rho(l)})_{l \in L}$ and $(u_i)_{i \in I} \cup (v_{\rho(l)})_{l \in (L - \{q\})}$ generate the same subset of E . Indeed, the second family is obtained from the first by replacing $v_{\rho(q)}$ by u_p , and vice-versa, and u_p is a linear combination of $(u_i)_{i \in (I - \{p\})} \cup (v_{\rho(l)})_{l \in L}$, by (1), and $v_{\rho(q)}$ is a linear combination of $(u_i)_{i \in I} \cup (v_{\rho(l)})_{l \in (L - \{q\})}$, by (2). Thus, the families $(u_i)_{i \in I} \cup (v_{\rho(l)})_{l \in (L - \{q\})}$ and $(v_j)_{j \in J}$ generate the same subspace of E , and the proposition holds for $L - \{q\}$ and the restriction of the injection $\rho: L \rightarrow J$ to $L - \{q\}$, since $L \cap I = \emptyset$ and $|L| = n - m$ imply that $(L - \{q\}) \cap I = \emptyset$ and $|L - \{q\}| = n - (m + 1)$. \square

The idea is that m of the vectors v_j can be *replaced* by the linearly independent u_i 's in such a way that the same subspace is still generated. The purpose of the function $\rho: L \rightarrow J$ is to pick $n - m$ elements j_1, \dots, j_{n-m} of J and to relabel them l_1, \dots, l_{n-m} in such a way that these new indices do not clash with the indices in I ; this way, the vectors $v_{j_1}, \dots, v_{j_{n-m}}$ who “survive” (i.e. are not replaced) are relabeled $v_{l_1}, \dots, v_{l_{n-m}}$, and the other m vectors v_j with $j \in J - \{j_1, \dots, j_{n-m}\}$ are replaced by the u_i . The index set of this new family is $I \cup L$.

Actually, one can prove that Proposition 2.12 implies Theorem 2.9 when the vector space is finitely generated. Putting Theorem 2.9 and Proposition 2.12 together, we obtain the following fundamental theorem.

Theorem 2.13. *Let E be a finitely generated vector space. Any family $(u_i)_{i \in I}$ generating E contains a subfamily $(u_j)_{j \in J}$ which is a basis of E . Any linearly independent family $(u_i)_{i \in I}$ can be extended to a family $(u_j)_{j \in J}$ which is a basis of E (with $I \subseteq J$). Furthermore, for every two bases $(u_i)_{i \in I}$ and $(v_j)_{j \in J}$ of E , we have $|I| = |J| = n$ for some fixed integer $n \geq 0$.*

Proof. The first part follows immediately by applying Theorem 2.9 with $L = \emptyset$ and $S = (u_i)_{i \in I}$. For the second part, consider the family $S' = (u_i)_{i \in I} \cup (v_h)_{h \in H}$, where $(v_h)_{h \in H}$ is any finitely generated family generating E , and with $I \cap H = \emptyset$. Then, apply Theorem 2.9 to $L = (u_i)_{i \in I}$ and to S' . For the last statement, assume that $(u_i)_{i \in I}$ and $(v_j)_{j \in J}$ are bases of E . Since $(u_i)_{i \in I}$ is linearly independent and $(v_j)_{j \in J}$ spans E , proposition 2.12 implies that $|I| \leq |J|$. A symmetric argument yields $|J| \leq |I|$. \square

Remark: Theorem 2.13 also holds for vector spaces that are not finitely generated. This can be shown as follows. Let $(u_i)_{i \in I}$ be a basis of E , let $(v_j)_{j \in J}$ be a generating family of E , and assume that I is infinite. For every $j \in J$, let $L_j \subseteq I$ be the finite set

$$L_j = \{i \in I \mid v_j = \sum_{i \in I} \lambda_i u_i, \lambda_i \neq 0\}.$$

Let $L = \bigcup_{j \in J} L_j$. By definition $L \subseteq I$, and since $(u_i)_{i \in I}$ is a basis of E , we must have $I = L$, since otherwise $(u_i)_{i \in L}$ would be another basis of E , and this would contradict the fact that $(u_i)_{i \in I}$ is linearly independent. Furthermore, J must be infinite, since otherwise, because the L_j are finite, I would be finite. But then, since $I = \bigcup_{j \in J} L_j$ with J infinite and the L_j finite, by a standard result of set theory, $|I| \leq |J|$. If $(v_j)_{j \in J}$ is also a basis, by a symmetric argument, we obtain $|J| \leq |I|$, and thus, $|I| = |J|$ for any two bases $(u_i)_{i \in I}$ and $(v_j)_{j \in J}$ of E .

When E is not finitely generated, we say that E is of infinite dimension. The *dimension* of a vector space E is the common cardinality of all of its bases and is denoted by $\dim(E)$. Clearly, if the field K itself is viewed as a vector space, then every family (a) where $a \in K$ and $a \neq 0$ is a basis. Thus $\dim(K) = 1$. Note that $\dim(\{0\}) = 0$.

If E is a vector space, for any subspace U of E , if $\dim(U) = 1$, then U is called a *line*; if $\dim(U) = 2$, then U is called a *plane*. If $\dim(U) = k$, then U is sometimes called a *k-plane*.

Let $(u_i)_{i \in I}$ be a basis of a vector space E . For any vector $v \in E$, since the family $(u_i)_{i \in I}$ generates E , there is a family $(\lambda_i)_{i \in I}$ of scalars in K , such that

$$v = \sum_{i \in I} \lambda_i u_i.$$

A very important fact is that the family $(\lambda_i)_{i \in I}$ is **unique**.

Proposition 2.14. *Given a vector space E , let $(u_i)_{i \in I}$ be a family of vectors in E . Let $v \in E$, and assume that $v = \sum_{i \in I} \lambda_i u_i$. Then, the family $(\lambda_i)_{i \in I}$ of scalars such that $v = \sum_{i \in I} \lambda_i u_i$ is unique iff $(u_i)_{i \in I}$ is linearly independent.*

Proof. First, assume that $(u_i)_{i \in I}$ is linearly independent. If $(\mu_i)_{i \in I}$ is another family of scalars in K such that $v = \sum_{i \in I} \mu_i u_i$, then we have

$$\sum_{i \in I} (\lambda_i - \mu_i) u_i = 0,$$

and since $(u_i)_{i \in I}$ is linearly independent, we must have $\lambda_i - \mu_i = 0$ for all $i \in I$, that is, $\lambda_i = \mu_i$ for all $i \in I$. The converse is shown by contradiction. If $(u_i)_{i \in I}$ was linearly dependent, there would be a family $(\mu_i)_{i \in I}$ of scalars not all null such that

$$\sum_{i \in I} \mu_i u_i = 0$$

and $\mu_j \neq 0$ for some $j \in I$. But then,

$$v = \sum_{i \in I} \lambda_i u_i + 0 = \sum_{i \in I} \lambda_i u_i + \sum_{i \in I} \mu_i u_i = \sum_{i \in I} (\lambda_i + \mu_i) u_i,$$

with $\lambda_j \neq \lambda_j + \mu_j$ since $\mu_j \neq 0$, contradicting the assumption that $(\lambda_i)_{i \in I}$ is the unique family such that $v = \sum_{i \in I} \lambda_i u_i$. \square

If $(u_i)_{i \in I}$ is a basis of a vector space E , for any vector $v \in E$, if $(x_i)_{i \in I}$ is the unique family of scalars in K such that

$$v = \sum_{i \in I} x_i u_i,$$

each x_i is called the *component* (or *coordinate*) of index i of v with respect to the basis $(u_i)_{i \in I}$.

Given a field K and any (nonempty) set I , we can form a vector space $K^{(I)}$ which, in some sense, is the standard vector space of dimension $|I|$.

Definition 2.13. Given a field K and any (nonempty) set I , let $K^{(I)}$ be the subset of the cartesian product K^I consisting of all families $(\lambda_i)_{i \in I}$ with finite support of scalars in K .⁴ We define addition and multiplication by a scalar as follows:

$$(\lambda_i)_{i \in I} + (\mu_i)_{i \in I} = (\lambda_i + \mu_i)_{i \in I},$$

and

$$\lambda \cdot (\mu_i)_{i \in I} = (\lambda \mu_i)_{i \in I}.$$

It is immediately verified that addition and multiplication by a scalar are well defined. Thus, $K^{(I)}$ is a vector space. Furthermore, because families with finite support are considered, the family $(e_i)_{i \in I}$ of vectors e_i , defined such that $(e_i)_j = 0$ if $j \neq i$ and $(e_i)_i = 1$, is clearly a basis of the vector space $K^{(I)}$. When $I = \{1, \dots, n\}$, we denote $K^{(I)}$ by K^n . The function $\iota: I \rightarrow K^{(I)}$, such that $\iota(i) = e_i$ for every $i \in I$, is clearly an injection.



When I is a finite set, $K^{(I)} = K^I$, but this is false when I is infinite. In fact, $\dim(K^{(I)}) = |I|$, but $\dim(K^I)$ is strictly greater when I is infinite.

Many interesting mathematical structures are vector spaces. A very important example is the set of linear maps between two vector spaces to be defined in the next section. Here is an example that will prepare us for the vector space of linear maps.

Example 2.10. Let X be any nonempty set and let E be a vector space. The set of all functions $f: X \rightarrow E$ can be made into a vector space as follows: Given any two functions $f: X \rightarrow E$ and $g: X \rightarrow E$, let $(f + g): X \rightarrow E$ be defined such that

$$(f + g)(x) = f(x) + g(x)$$

for all $x \in X$, and for every $\lambda \in K$, let $\lambda f: X \rightarrow E$ be defined such that

$$(\lambda f)(x) = \lambda f(x)$$

for all $x \in X$. The axioms of a vector space are easily verified. Now, let $E = K$, and let I be the set of all nonempty subsets of X . For every $S \in I$, let $f_S: X \rightarrow E$ be the function such that $f_S(x) = 1$ iff $x \in S$, and $f_S(x) = 0$ iff $x \notin S$. We leave as an exercise to show that $(f_S)_{S \in I}$ is linearly independent.

2.5 Linear Maps

A function between two vector spaces that preserves the vector space structure is called a homomorphism of vector spaces, or linear map. Linear maps formalize the concept of linearity of a function. In the rest of this section, we assume that all vector spaces are over a given field K (say \mathbb{R}).

⁴Where K^I denotes the set of all functions from I to K .

Definition 2.14. Given two vector spaces E and F , a *linear map* between E and F is a function $f: E \rightarrow F$ satisfying the following two conditions:

$$\begin{aligned} f(x + y) &= f(x) + f(y) && \text{for all } x, y \in E; \\ f(\lambda x) &= \lambda f(x) && \text{for all } \lambda \in K, x \in E. \end{aligned}$$

Setting $x = y = 0$ in the first identity, we get $f(0) = 0$. The basic property of linear maps is that they transform linear combinations into linear combinations. Given a family $(u_i)_{i \in I}$ of vectors in E , given any family $(\lambda_i)_{i \in I}$ of scalars in K , we have

$$f\left(\sum_{i \in I} \lambda_i u_i\right) = \sum_{i \in I} \lambda_i f(u_i).$$

The above identity is shown by induction on the size of the support of the family $(\lambda_i u_i)_{i \in I}$, using the properties of Definition 2.14.

Example 2.11.

1. The map $f: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ defined such that

$$\begin{aligned} x' &= x - y \\ y' &= x + y \end{aligned}$$

is a linear map. The reader should check that it is the composition of a rotation by $\pi/4$ with a magnification of ratio $\sqrt{2}$.

2. For any vector space E , the *identity map* $\text{id}: E \rightarrow E$ given by

$$\text{id}(u) = u \quad \text{for all } u \in E$$

is a linear map. When we want to be more precise, we write id_E instead of id .

3. The map $D: \mathbb{R}[X] \rightarrow \mathbb{R}[X]$ defined such that

$$D(f(X)) = f'(X),$$

where $f'(X)$ is the derivative of the polynomial $f(X)$, is a linear map.

4. The map $\Phi: \mathcal{C}([a, b]) \rightarrow \mathbb{R}$ given by

$$\Phi(f) = \int_a^b f(t) dt,$$

where $\mathcal{C}([a, b])$ is the set of continuous functions defined on the interval $[a, b]$, is a linear map.

5. The function $\langle -, - \rangle: \mathcal{C}([a, b]) \times \mathcal{C}([a, b]) \rightarrow \mathbb{R}$ given by

$$\langle f, g \rangle = \int_a^b f(t)g(t)dt,$$

is linear in each of the variable f, g . It also satisfies the properties $\langle f, g \rangle = \langle g, f \rangle$ and $\langle f, f \rangle = 0$ iff $f = 0$. It is an example of an *inner product*.

Definition 2.15. Given a linear map $f: E \rightarrow F$, we define its *image (or range)* $\text{Im } f = f(E)$, as the set

$$\text{Im } f = \{y \in F \mid (\exists x \in E)(y = f(x))\},$$

and its *Kernel (or nullspace)* $\text{Ker } f = f^{-1}(0)$, as the set

$$\text{Ker } f = \{x \in E \mid f(x) = 0\}.$$

Proposition 2.15. *Given a linear map $f: E \rightarrow F$, the set $\text{Im } f$ is a subspace of F and the set $\text{Ker } f$ is a subspace of E . The linear map $f: E \rightarrow F$ is injective iff $\text{Ker } f = 0$ (where 0 is the trivial subspace $\{0\}$).*

Proof. Given any $x, y \in \text{Im } f$, there are some $u, v \in E$ such that $x = f(u)$ and $y = f(v)$, and for all $\lambda, \mu \in K$, we have

$$f(\lambda u + \mu v) = \lambda f(u) + \mu f(v) = \lambda x + \mu y,$$

and thus, $\lambda x + \mu y \in \text{Im } f$, showing that $\text{Im } f$ is a subspace of F .

Given any $x, y \in \text{Ker } f$, we have $f(x) = 0$ and $f(y) = 0$, and thus,

$$f(\lambda x + \mu y) = \lambda f(x) + \mu f(y) = 0,$$

that is, $\lambda x + \mu y \in \text{Ker } f$, showing that $\text{Ker } f$ is a subspace of E .

First, assume that $\text{Ker } f = 0$. We need to prove that $f(x) = f(y)$ implies that $x = y$. However, if $f(x) = f(y)$, then $f(x) - f(y) = 0$, and by linearity of f we get $f(x - y) = 0$. Because $\text{Ker } f = 0$, we must have $x - y = 0$, that is $x = y$, so f is injective. Conversely, assume that f is injective. If $x \in \text{Ker } f$, that is $f(x) = 0$, since $f(0) = 0$ we have $f(x) = f(0)$, and by injectivity, $x = 0$, which proves that $\text{Ker } f = 0$. Therefore, f is injective iff $\text{Ker } f = 0$. \square

Since by Proposition 2.15, the image $\text{Im } f$ of a linear map f is a subspace of F , we can define the *rank* $\text{rk}(f)$ of f as the dimension of $\text{Im } f$.

A fundamental property of bases in a vector space is that they allow the definition of linear maps as unique homomorphic extensions, as shown in the following proposition.

Proposition 2.16. *Given any two vector spaces E and F , given any basis $(u_i)_{i \in I}$ of E , given any other family of vectors $(v_i)_{i \in I}$ in F , there is a unique linear map $f: E \rightarrow F$ such that $f(u_i) = v_i$ for all $i \in I$. Furthermore, f is injective iff $(v_i)_{i \in I}$ is linearly independent, and f is surjective iff $(v_i)_{i \in I}$ generates F .*

Proof. If such a linear map $f: E \rightarrow F$ exists, since $(u_i)_{i \in I}$ is a basis of E , every vector $x \in E$ can be written uniquely as a linear combination

$$x = \sum_{i \in I} x_i u_i,$$

and by linearity, we must have

$$f(x) = \sum_{i \in I} x_i f(u_i) = \sum_{i \in I} x_i v_i.$$

Define the function $f: E \rightarrow F$, by letting

$$f(x) = \sum_{i \in I} x_i v_i$$

for every $x = \sum_{i \in I} x_i u_i$. It is easy to verify that f is indeed linear, it is unique by the previous reasoning, and obviously, $f(u_i) = v_i$.

Now, assume that f is injective. Let $(\lambda_i)_{i \in I}$ be any family of scalars, and assume that

$$\sum_{i \in I} \lambda_i v_i = 0.$$

Since $v_i = f(u_i)$ for every $i \in I$, we have

$$f\left(\sum_{i \in I} \lambda_i u_i\right) = \sum_{i \in I} \lambda_i f(u_i) = \sum_{i \in I} \lambda_i v_i = 0.$$

Since f is injective iff $\text{Ker } f = 0$, we have

$$\sum_{i \in I} \lambda_i u_i = 0,$$

and since $(u_i)_{i \in I}$ is a basis, we have $\lambda_i = 0$ for all $i \in I$, which shows that $(v_i)_{i \in I}$ is linearly independent. Conversely, assume that $(v_i)_{i \in I}$ is linearly independent. Since $(u_i)_{i \in I}$ is a basis of E , every vector $x \in E$ is a linear combination $x = \sum_{i \in I} \lambda_i u_i$ of $(u_i)_{i \in I}$. If

$$f(x) = f\left(\sum_{i \in I} \lambda_i u_i\right) = 0,$$

then

$$\sum_{i \in I} \lambda_i v_i = \sum_{i \in I} \lambda_i f(u_i) = f\left(\sum_{i \in I} \lambda_i u_i\right) = 0,$$

and $\lambda_i = 0$ for all $i \in I$ because $(v_i)_{i \in I}$ is linearly independent, which means that $x = 0$. Therefore, $\text{Ker } f = 0$, which implies that f is injective. The part where f is surjective is left as a simple exercise. \square

By the second part of Proposition 2.16, an injective linear map $f: E \rightarrow F$ sends a basis $(u_i)_{i \in I}$ to a linearly independent family $(f(u_i))_{i \in I}$ of F , which is also a basis when f is bijective. Also, when E and F have the same finite dimension n , $(u_i)_{i \in I}$ is a basis of E , and $f: E \rightarrow F$ is injective, then $(f(u_i))_{i \in I}$ is a basis of F (by Proposition 2.10).

We can now show that the vector space $K^{(I)}$ of Definition 2.13 has a universal property that amounts to saying that $K^{(I)}$ is the vector space freely generated by I . Recall that $\iota: I \rightarrow K^{(I)}$, such that $\iota(i) = e_i$ for every $i \in I$, is an injection from I to $K^{(I)}$.

Proposition 2.17. *Given any set I , for any vector space F , and for any function $f: I \rightarrow F$, there is a unique linear map $\bar{f}: K^{(I)} \rightarrow F$, such that*

$$f = \bar{f} \circ \iota,$$

as in the following diagram:

$$\begin{array}{ccc} I & \xrightarrow{\iota} & K^{(I)} \\ & \searrow f & \downarrow \bar{f} \\ & & F \end{array}$$

Proof. If such a linear map $\bar{f}: K^{(I)} \rightarrow F$ exists, since $f = \bar{f} \circ \iota$, we must have

$$f(i) = \bar{f}(\iota(i)) = \bar{f}(e_i),$$

for every $i \in I$. However, the family $(e_i)_{i \in I}$ is a basis of $K^{(I)}$, and $(f(i))_{i \in I}$ is a family of vectors in F , and by Proposition 2.16, there is a unique linear map $\bar{f}: K^{(I)} \rightarrow F$ such that $\bar{f}(e_i) = f(i)$ for every $i \in I$, which proves the existence and uniqueness of a linear map \bar{f} such that $f = \bar{f} \circ \iota$. \square

The following simple proposition is also useful.

Proposition 2.18. *Given any two vector spaces E and F , with F nontrivial, given any family $(u_i)_{i \in I}$ of vectors in E , the following properties hold:*

- (1) *The family $(u_i)_{i \in I}$ generates E iff for every family of vectors $(v_i)_{i \in I}$ in F , there is at most one linear map $f: E \rightarrow F$ such that $f(u_i) = v_i$ for all $i \in I$.*
- (2) *The family $(u_i)_{i \in I}$ is linearly independent iff for every family of vectors $(v_i)_{i \in I}$ in F , there is some linear map $f: E \rightarrow F$ such that $f(u_i) = v_i$ for all $i \in I$.*

Proof. (1) If there is any linear map $f: E \rightarrow F$ such that $f(u_i) = v_i$ for all $i \in I$, since $(u_i)_{i \in I}$ generates E , every vector $x \in E$ can be written as some linear combination

$$x = \sum_{i \in I} x_i u_i,$$

and by linearity, we must have

$$f(x) = \sum_{i \in I} x_i f(u_i) = \sum_{i \in I} x_i v_i.$$

This shows that f is unique if it exists. Conversely, assume that $(u_i)_{i \in I}$ does not generate E . Since F is nontrivial, there is some vector $y \in F$ such that $y \neq 0$. Since $(u_i)_{i \in I}$ does not generate E , there is some vector $w \in E$ that is not in the subspace generated by $(u_i)_{i \in I}$. By Theorem 2.13, there is a linearly independent subfamily $(u_i)_{i \in I_0}$ of $(u_i)_{i \in I}$ generating the same subspace. Since by hypothesis, $w \in E$ is not in the subspace generated by $(u_i)_{i \in I_0}$, by Lemma 2.8 and by Theorem 2.13 again, there is a basis $(e_j)_{j \in I_0 \cup J}$ of E , such that $e_i = u_i$, for all $i \in I_0$, and $w = e_{j_0}$, for some $j_0 \in J$. Letting $(v_i)_{i \in I}$ be the family in F such that $v_i = 0$ for all $i \in I$, defining $f: E \rightarrow F$ to be the constant linear map with value 0, we have a linear map such that $f(u_i) = 0$ for all $i \in I$. By Proposition 2.16, there is a unique linear map $g: E \rightarrow F$ such that $g(w) = y$, and $g(e_j) = 0$, for all $j \in (I_0 \cup J) - \{j_0\}$. By definition of the basis $(e_j)_{j \in I_0 \cup J}$ of E , we have, $g(u_i) = 0$ for all $i \in I$, and since $f \neq g$, this contradicts the fact that there is at most one such map.

(2) If the family $(u_i)_{i \in I}$ is linearly independent, then by Theorem 2.13, $(u_i)_{i \in I}$ can be extended to a basis of E , and the conclusion follows by Proposition 2.16. Conversely, assume that $(u_i)_{i \in I}$ is linearly dependent. Then, there is some family $(\lambda_i)_{i \in I}$ of scalars (not all zero) such that

$$\sum_{i \in I} \lambda_i u_i = 0.$$

By the assumption, for any nonzero vector, $y \in F$, for every $i \in I$, there is some linear map $f_i: E \rightarrow F$, such that $f_i(u_i) = y$, and $f_i(u_j) = 0$, for $j \in I - \{i\}$. Then, we would get

$$0 = f_i\left(\sum_{i \in I} \lambda_i u_i\right) = \sum_{i \in I} \lambda_i f_i(u_i) = \lambda_i y,$$

and since $y \neq 0$, this implies $\lambda_i = 0$, for every $i \in I$. Thus, $(u_i)_{i \in I}$ is linearly independent. \square

Given vector spaces E , F , and G , and linear maps $f: E \rightarrow F$ and $g: F \rightarrow G$, it is easily verified that the composition $g \circ f: E \rightarrow G$ of f and g is a linear map.

A linear map $f: E \rightarrow F$ is an *isomorphism* iff there is a linear map $g: F \rightarrow E$, such that

$$g \circ f = \text{id}_E \quad \text{and} \quad f \circ g = \text{id}_F. \quad (*)$$

Such a map g is unique. This is because if g and h both satisfy $g \circ f = \text{id}_E$, $f \circ g = \text{id}_F$, $h \circ f = \text{id}_E$, and $f \circ h = \text{id}_F$, then

$$g = g \circ \text{id}_F = g \circ (f \circ h) = (g \circ f) \circ h = \text{id}_E \circ h = h.$$

The map g satisfying $(*)$ above is called the *inverse* of f and it is also denoted by f^{-1} .

Proposition 2.16 implies that if E and F are two vector spaces, $(u_i)_{i \in I}$ is a basis of E , and $f: E \rightarrow F$ is a linear map which is an isomorphism, then the family $(f(u_i))_{i \in I}$ is a basis of F .

One can verify that if $f: E \rightarrow F$ is a bijective linear map, then its inverse $f^{-1}: F \rightarrow E$ is also a linear map, and thus f is an isomorphism.

Another useful corollary of Proposition 2.16 is this:

Proposition 2.19. *Let E be a vector space of finite dimension $n \geq 1$ and let $f: E \rightarrow E$ be any linear map. The following properties hold:*

- (1) *If f has a left inverse g , that is, if g is a linear map such that $g \circ f = \text{id}$, then f is an isomorphism and $f^{-1} = g$.*
- (2) *If f has a right inverse h , that is, if h is a linear map such that $f \circ h = \text{id}$, then f is an isomorphism and $f^{-1} = h$.*

Proof. (1) The equation $g \circ f = \text{id}$ implies that f is injective; this is a standard result about functions (if $f(x) = f(y)$, then $g(f(x)) = g(f(y))$, which implies that $x = y$ since $g \circ f = \text{id}$). Let (u_1, \dots, u_n) be any basis of E . By Proposition 2.16, since f is injective, $(f(u_1), \dots, f(u_n))$ is linearly independent, and since E has dimension n , it is a basis of E (if $(f(u_1), \dots, f(u_n))$ doesn't span E , then it can be extended to a basis of dimension strictly greater than n , contradicting Theorem 2.13). Then, f is bijective, and by a previous observation its inverse is a linear map. We also have

$$g = g \circ \text{id} = g \circ (f \circ f^{-1}) = (g \circ f) \circ f^{-1} = \text{id} \circ f^{-1} = f^{-1}.$$

(2) The equation $f \circ h = \text{id}$ implies that f is surjective; this is a standard result about functions (for any $y \in E$, we have $f(h(y)) = y$). Let (u_1, \dots, u_n) be any basis of E . By Proposition 2.16, since f is surjective, $(f(u_1), \dots, f(u_n))$ spans E , and since E has dimension n , it is a basis of E (if $(f(u_1), \dots, f(u_n))$ is not linearly independent, then because it spans E , it contains a basis of dimension strictly smaller than n , contradicting Theorem 2.13). Then, f is bijective, and by a previous observation its inverse is a linear map. We also have

$$h = \text{id} \circ h = (f^{-1} \circ f) \circ h = f^{-1} \circ (f \circ h) = f^{-1} \circ \text{id} = f^{-1}.$$

This completes the proof. □

The set of all linear maps between two vector spaces E and F is denoted by $\text{Hom}(E, F)$ or by $\mathcal{L}(E; F)$ (the notation $\mathcal{L}(E; F)$ is usually reserved to the set of continuous linear maps, where E and F are normed vector spaces). When we wish to be more precise and specify the field K over which the vector spaces E and F are defined we write $\text{Hom}_K(E, F)$.

The set $\text{Hom}(E, F)$ is a vector space under the operations defined at the end of Section 2.1, namely

$$(f + g)(x) = f(x) + g(x)$$

for all $x \in E$, and

$$(\lambda f)(x) = \lambda f(x)$$

for all $x \in E$. The point worth checking carefully is that λf is indeed a linear map, which uses the commutativity of $*$ in the field K . Indeed, we have

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

When E and F have finite dimensions, the vector space $\text{Hom}(E, F)$ also has finite dimension, as we shall see shortly. When $E = F$, a linear map $f: E \rightarrow E$ is also called an *endomorphism*. It is also important to note that composition confers to $\text{Hom}(E, E)$ a ring structure. Indeed, composition is an operation $\circ: \text{Hom}(E, E) \times \text{Hom}(E, E) \rightarrow \text{Hom}(E, E)$, which is associative and has an identity id_E , and the distributivity properties hold:

$$\begin{aligned} (g_1 + g_2) \circ f &= g_1 \circ f + g_2 \circ f; \\ g \circ (f_1 + f_2) &= g \circ f_1 + g \circ f_2. \end{aligned}$$

The ring $\text{Hom}(E, E)$ is an example of a noncommutative ring. It is easily seen that the set of bijective linear maps $f: E \rightarrow E$ is a group under composition. Bijective linear maps are also called *automorphisms*. The group of automorphisms of E is called the *general linear group (of E)*, and it is denoted by $\mathbf{GL}(E)$, or by $\text{Aut}(E)$, or when $E = K^n$, by $\mathbf{GL}(n, K)$, or even by $\mathbf{GL}(n)$.

Although in this book, we will not have many occasions to use quotient spaces, they are fundamental in algebra. The next section may be omitted until needed.

2.6 Quotient Spaces

Let E be a vector space, and let M be any subspace of E . The subspace M induces a relation \equiv_M on E , defined as follows: For all $u, v \in E$,

$$u \equiv_M v \text{ iff } u - v \in M.$$

We have the following simple proposition.

Proposition 2.20. *Given any vector space E and any subspace M of E , the relation \equiv_M is an equivalence relation with the following two congruential properties:*

1. *If $u_1 \equiv_M v_1$ and $u_2 \equiv_M v_2$, then $u_1 + u_2 \equiv_M v_1 + v_2$, and*
2. *if $u \equiv_M v$, then $\lambda u \equiv_M \lambda v$.*

Proof. It is obvious that \equiv_M is an equivalence relation. Note that $u_1 \equiv_M v_1$ and $u_2 \equiv_M v_2$ are equivalent to $u_1 - v_1 = w_1$ and $u_2 - v_2 = w_2$, with $w_1, w_2 \in M$, and thus,

$$(u_1 + u_2) - (v_1 + v_2) = w_1 + w_2,$$

and $w_1 + w_2 \in M$, since M is a subspace of E . Thus, we have $u_1 + u_2 \equiv_M v_1 + v_2$. If $u - v = w$, with $w \in M$, then

$$\lambda u - \lambda v = \lambda w,$$

and $\lambda w \in M$, since M is a subspace of E , and thus $\lambda u \equiv_M \lambda v$. \square

Proposition 2.20 shows that we can define addition and multiplication by a scalar on the set E/M of equivalence classes of the equivalence relation \equiv_M .

Definition 2.16. Given any vector space E and any subspace M of E , we define the following operations of addition and multiplication by a scalar on the set E/M of equivalence classes of the equivalence relation \equiv_M as follows: for any two equivalence classes $[u], [v] \in E/M$, we have

$$\begin{aligned} [u] + [v] &= [u + v], \\ \lambda[u] &= [\lambda u]. \end{aligned}$$

By Proposition 2.20, the above operations do not depend on the specific choice of representatives in the equivalence classes $[u], [v] \in E/M$. It is also immediate to verify that E/M is a vector space. The function $\pi: E \rightarrow E/M$, defined such that $\pi(u) = [u]$ for every $u \in E$, is a surjective linear map called the *natural projection of E onto E/M* . The vector space E/M is called the *quotient space of E by the subspace M* .

Given any linear map $f: E \rightarrow F$, we know that $\text{Ker } f$ is a subspace of E , and it is immediately verified that $\text{Im } f$ is isomorphic to the quotient space $E/\text{Ker } f$.

2.7 Summary

The main concepts and results of this chapter are listed below:

- Groups, rings and fields.
- The notion of a *vector space*.
- *Families* of vectors.
- *Linear combinations* of vectors; *linear dependence* and *linear independence* of a family of vectors.
- *Linear subspaces*.
- *Spanning* (or *generating*) family; *generators*, *finitely generated subspace*; *basis of a subspace*.
- *Every linearly independent family can be extended to a basis* (Theorem 2.9).

- A family B of vectors is a basis iff it is a maximal linearly independent family iff it is a minimal generating family (Proposition 2.10).
- The replacement lemma (Proposition 2.12).
- Any two bases in a finitely generated vector space E have the *same number of elements*; this is the *dimension* of E (Theorem 2.13).
- *Hyperplanes*.
- Every vector has a *unique representation* over a basis (in terms of its coordinates).
- The notion of a *linear map*.
- The *image* $\text{Im } f$ (or *range*) of a linear map f .
- The *kernel* $\text{Ker } f$ (or *nullspace*) of a linear map f .
- The *rank* $\text{rk}(f)$ of a linear map f .
- The image and the kernel of a linear map are subspaces. A linear map is injective iff its kernel is the trivial space (0) (Proposition 2.15).
- The *unique homomorphic extension property* of linear maps with respect to bases (Proposition 2.16).
- *Quotient spaces*.

Chapter 3

Matrices and Linear Maps

3.1 Matrices

Proposition 2.16 shows that given two vector spaces E and F and a basis $(u_j)_{j \in J}$ of E , every linear map $f: E \rightarrow F$ is uniquely determined by the family $(f(u_j))_{j \in J}$ of the images under f of the vectors in the basis $(u_j)_{j \in J}$. Thus, in particular, taking $F = K^{(J)}$, we get an isomorphism between any vector space E of dimension $|J|$ and $K^{(J)}$. If $J = \{1, \dots, n\}$, a vector space E of dimension n is isomorphic to the vector space K^n . If we also have a basis $(v_i)_{i \in I}$ of F , then every vector $f(u_j)$ can be written in a unique way as

$$f(u_j) = \sum_{i \in I} a_{ij} v_i,$$

where $j \in J$, for a family of scalars $(a_{ij})_{i \in I}$. Thus, with respect to the two bases $(u_j)_{j \in J}$ of E and $(v_i)_{i \in I}$ of F , the linear map f is completely determined by a possibly infinite “ $I \times J$ -matrix” $M(f) = (a_{ij})_{i \in I, j \in J}$.

Remark: Note that we intentionally assigned the index set J to the basis $(u_j)_{j \in J}$ of E , and the index I to the basis $(v_i)_{i \in I}$ of F , so that the rows of the matrix $M(f)$ associated with $f: E \rightarrow F$ are indexed by I , and the columns of the matrix $M(f)$ are indexed by J . Obviously, this causes a mildly unpleasant reversal. If we had considered the bases $(u_i)_{i \in I}$ of E and $(v_j)_{j \in J}$ of F , we would obtain a $J \times I$ -matrix $M(f) = (a_{ji})_{j \in J, i \in I}$. No matter what we do, there will be a reversal! We decided to stick to the bases $(u_j)_{j \in J}$ of E and $(v_i)_{i \in I}$ of F , so that we get an $I \times J$ -matrix $M(f)$, knowing that we may occasionally suffer from this decision!

When I and J are finite, and say, when $|I| = m$ and $|J| = n$, the linear map f is determined by the matrix $M(f)$ whose entries in the j -th column are the components of the

vector $f(u_j)$ over the basis (v_1, \dots, v_m) , that is, the matrix

$$M(f) = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix}$$

whose entry on row i and column j is a_{ij} ($1 \leq i \leq m$, $1 \leq j \leq n$).

We will now show that when E and F have finite dimension, linear maps can be very conveniently represented by matrices, and that composition of linear maps corresponds to matrix multiplication. We will follow rather closely an elegant presentation method due to Emil Artin.

Let E and F be two vector spaces, and assume that E has a finite basis (u_1, \dots, u_n) and that F has a finite basis (v_1, \dots, v_m) . Recall that we have shown that every vector $x \in E$ can be written in a unique way as

$$x = x_1 u_1 + \dots + x_n u_n,$$

and similarly every vector $y \in F$ can be written in a unique way as

$$y = y_1 v_1 + \dots + y_m v_m.$$

Let $f: E \rightarrow F$ be a linear map between E and F . Then, for every $x = x_1 u_1 + \dots + x_n u_n$ in E , by linearity, we have

$$f(x) = x_1 f(u_1) + \dots + x_n f(u_n).$$

Let

$$f(u_j) = a_{1j} v_1 + \dots + a_{mj} v_m,$$

or more concisely,

$$f(u_j) = \sum_{i=1}^m a_{ij} v_i,$$

for every j , $1 \leq j \leq n$. This can be expressed by writing the coefficients $a_{1j}, a_{2j}, \dots, a_{mj}$ of $f(u_j)$ over the basis (v_1, \dots, v_m) , as the j th column of a matrix, as shown below:

$$\begin{array}{cccc} & f(u_1) & f(u_2) & \dots & f(u_n) \\ \begin{array}{c} v_1 \\ v_2 \\ \vdots \\ v_m \end{array} & \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix} \end{array}$$

Then, substituting the right-hand side of each $f(u_j)$ into the expression for $f(x)$, we get

$$f(x) = x_1 \left(\sum_{i=1}^m a_{i1} v_i \right) + \dots + x_n \left(\sum_{i=1}^m a_{in} v_i \right),$$

which, by regrouping terms to obtain a linear combination of the v_i , yields

$$f(x) = \left(\sum_{j=1}^n a_{1j}x_j\right)v_1 + \cdots + \left(\sum_{j=1}^n a_{mj}x_j\right)v_m.$$

Thus, letting $f(x) = y = y_1v_1 + \cdots + y_mv_m$, we have

$$y_i = \sum_{j=1}^n a_{ij}x_j \tag{1}$$

for all i , $1 \leq i \leq m$.

To make things more concrete, let us treat the case where $n = 3$ and $m = 2$. In this case,

$$\begin{aligned} f(u_1) &= a_{11}v_1 + a_{21}v_2 \\ f(u_2) &= a_{12}v_1 + a_{22}v_2 \\ f(u_3) &= a_{13}v_1 + a_{23}v_2, \end{aligned}$$

which in matrix form is expressed by

$$\begin{matrix} f(u_1) & f(u_2) & f(u_3) \\ v_1 & \begin{pmatrix} a_{11} & a_{12} & a_{13} \end{pmatrix} \\ v_2 & \begin{pmatrix} a_{21} & a_{22} & a_{23} \end{pmatrix} \end{matrix},$$

and for any $x = x_1u_1 + x_2u_2 + x_3u_3$, we have

$$\begin{aligned} f(x) &= f(x_1u_1 + x_2u_2 + x_3u_3) \\ &= x_1f(u_1) + x_2f(u_2) + x_3f(u_3) \\ &= x_1(a_{11}v_1 + a_{21}v_2) + x_2(a_{12}v_1 + a_{22}v_2) + x_3(a_{13}v_1 + a_{23}v_2) \\ &= (a_{11}x_1 + a_{12}x_2 + a_{13}x_3)v_1 + (a_{21}x_1 + a_{22}x_2 + a_{23}x_3)v_2. \end{aligned}$$

Consequently, since

$$y = y_1v_1 + y_2v_2,$$

we have

$$\begin{aligned} y_1 &= a_{11}x_1 + a_{12}x_2 + a_{13}x_3 \\ y_2 &= a_{21}x_1 + a_{22}x_2 + a_{23}x_3. \end{aligned}$$

This agrees with the matrix equation

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}.$$

Let us now consider how the composition of linear maps is expressed in terms of bases.

Let E , F , and G , be three vectors spaces with respective bases (u_1, \dots, u_p) for E , (v_1, \dots, v_n) for F , and (w_1, \dots, w_m) for G . Let $g: E \rightarrow F$ and $f: F \rightarrow G$ be linear maps. As explained earlier, $g: E \rightarrow F$ is determined by the images of the basis vectors u_j , and $f: F \rightarrow G$ is determined by the images of the basis vectors v_k . We would like to understand how $f \circ g: E \rightarrow G$ is determined by the images of the basis vectors u_j .

Remark: Note that we are considering linear maps $g: E \rightarrow F$ and $f: F \rightarrow G$, instead of $f: E \rightarrow F$ and $g: F \rightarrow G$, which yields the composition $f \circ g: E \rightarrow G$ instead of $g \circ f: E \rightarrow G$. Our perhaps unusual choice is motivated by the fact that if f is represented by a matrix $M(f) = (a_{ik})$ and g is represented by a matrix $M(g) = (b_{kj})$, then $f \circ g: E \rightarrow G$ is represented by the product AB of the matrices A and B . If we had adopted the other choice where $f: E \rightarrow F$ and $g: F \rightarrow G$, then $g \circ f: E \rightarrow G$ would be represented by the product BA . Personally, we find it easier to remember the formula for the entry in row i and column of j of the product of two matrices when this product is written by AB , rather than BA . Obviously, this is a matter of taste! We will have to live with our perhaps unorthodox choice.

Thus, let

$$f(v_k) = \sum_{i=1}^m a_{ik} w_i,$$

for every k , $1 \leq k \leq n$, and let

$$g(u_j) = \sum_{k=1}^n b_{kj} v_k,$$

for every j , $1 \leq j \leq p$; in matrix form, we have

$$\begin{array}{cccc} & f(v_1) & f(v_2) & \dots & f(v_n) \\ \begin{array}{c} w_1 \\ w_2 \\ \vdots \\ w_m \end{array} & \left(\begin{array}{cccc} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{array} \right) \end{array}$$

and

$$\begin{array}{cccc} & g(u_1) & g(u_2) & \dots & g(u_p) \\ \begin{array}{c} v_1 \\ v_2 \\ \vdots \\ v_n \end{array} & \left(\begin{array}{cccc} b_{11} & b_{12} & \dots & b_{1p} \\ b_{21} & b_{22} & \dots & b_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ b_{n1} & b_{n2} & \dots & b_{np} \end{array} \right) \end{array}$$

By previous considerations, for every

$$x = x_1u_1 + \cdots + x_pu_p,$$

letting $g(x) = y = y_1v_1 + \cdots + y_nv_n$, we have

$$y_k = \sum_{j=1}^p b_{kj}x_j \quad (2)$$

for all k , $1 \leq k \leq n$, and for every

$$y = y_1v_1 + \cdots + y_nv_n,$$

letting $f(y) = z = z_1w_1 + \cdots + z_mw_m$, we have

$$z_i = \sum_{k=1}^n a_{ik}y_k \quad (3)$$

for all i , $1 \leq i \leq m$. Then, if $y = g(x)$ and $z = f(y)$, we have $z = f(g(x))$, and in view of (2) and (3), we have

$$\begin{aligned} z_i &= \sum_{k=1}^n a_{ik} \left(\sum_{j=1}^p b_{kj}x_j \right) \\ &= \sum_{k=1}^n \sum_{j=1}^p a_{ik}b_{kj}x_j \\ &= \sum_{j=1}^p \sum_{k=1}^n a_{ik}b_{kj}x_j \\ &= \sum_{j=1}^p \left(\sum_{k=1}^n a_{ik}b_{kj} \right) x_j. \end{aligned}$$

Thus, defining c_{ij} such that

$$c_{ij} = \sum_{k=1}^n a_{ik}b_{kj},$$

for $1 \leq i \leq m$, and $1 \leq j \leq p$, we have

$$z_i = \sum_{j=1}^p c_{ij}x_j \quad (4)$$

Identity (4) suggests defining a multiplication operation on matrices, and we proceed to do so. We have the following definitions.

Definition 3.1. Given a field K , an $m \times n$ -matrix is a family $(a_{ij})_{1 \leq i \leq m, 1 \leq j \leq n}$ of scalars in K , represented as an array

$$\begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix}$$

In the special case where $m = 1$, we have a *row vector*, represented as

$$(a_{11} \cdots a_{1n})$$

and in the special case where $n = 1$, we have a *column vector*, represented as

$$\begin{pmatrix} a_{11} \\ \vdots \\ a_{m1} \end{pmatrix}$$

In these last two cases, we usually omit the constant index 1 (first index in case of a row, second index in case of a column). The set of all $m \times n$ -matrices is denoted by $M_{m,n}(K)$ or $M_{m,n}$. An $n \times n$ -matrix is called a *square matrix of dimension n* . The set of all square matrices of dimension n is denoted by $M_n(K)$, or M_n .

Remark: As defined, a matrix $A = (a_{ij})_{1 \leq i \leq m, 1 \leq j \leq n}$ is a *family*, that is, a function from $\{1, 2, \dots, m\} \times \{1, 2, \dots, n\}$ to K . As such, there is no reason to assume an ordering on the indices. Thus, the matrix A can be represented in many different ways as an array, by adopting different orders for the rows or the columns. However, it is customary (and usually convenient) to assume the natural ordering on the sets $\{1, 2, \dots, m\}$ and $\{1, 2, \dots, n\}$, and to represent A as an array according to this ordering of the rows and columns.

We also define some operations on matrices as follows.

Definition 3.2. Given two $m \times n$ matrices $A = (a_{ij})$ and $B = (b_{ij})$, we define their *sum* $A + B$ as the matrix $C = (c_{ij})$ such that $c_{ij} = a_{ij} + b_{ij}$; that is,

$$\begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix} + \begin{pmatrix} b_{11} & b_{12} & \cdots & b_{1n} \\ b_{21} & b_{22} & \cdots & b_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ b_{m1} & b_{m2} & \cdots & b_{mn} \end{pmatrix} = \begin{pmatrix} a_{11} + b_{11} & a_{12} + b_{12} & \cdots & a_{1n} + b_{1n} \\ a_{21} + b_{21} & a_{22} + b_{22} & \cdots & a_{2n} + b_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} + b_{m1} & a_{m2} + b_{m2} & \cdots & a_{mn} + b_{mn} \end{pmatrix}.$$

For any matrix $A = (a_{ij})$, we let $-A$ be the matrix $(-a_{ij})$. Given a scalar $\lambda \in K$, we define the matrix λA as the matrix $C = (c_{ij})$ such that $c_{ij} = \lambda a_{ij}$; that is

$$\lambda \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix} = \begin{pmatrix} \lambda a_{11} & \lambda a_{12} & \dots & \lambda a_{1n} \\ \lambda a_{21} & \lambda a_{22} & \dots & \lambda a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda a_{m1} & \lambda a_{m2} & \dots & \lambda a_{mn} \end{pmatrix}.$$

Given an $m \times n$ matrices $A = (a_{ik})$ and an $n \times p$ matrices $B = (b_{kj})$, we define their *product* AB as the $m \times p$ matrix $C = (c_{ij})$ such that

$$c_{ij} = \sum_{k=1}^n a_{ik} b_{kj},$$

for $1 \leq i \leq m$, and $1 \leq j \leq p$. In the product $AB = C$ shown below

$$\begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix} \begin{pmatrix} b_{11} & b_{12} & \dots & b_{1p} \\ b_{21} & b_{22} & \dots & b_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ b_{n1} & b_{n2} & \dots & b_{np} \end{pmatrix} = \begin{pmatrix} c_{11} & c_{12} & \dots & c_{1p} \\ c_{21} & c_{22} & \dots & c_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ c_{m1} & c_{m2} & \dots & c_{mp} \end{pmatrix}$$

note that the entry of index i and j of the matrix AB obtained by multiplying the matrices A and B can be identified with the product of the row matrix corresponding to the i -th row of A with the column matrix corresponding to the j -column of B :

$$(a_{i1} \dots a_{in}) \begin{pmatrix} b_{1j} \\ \vdots \\ b_{nj} \end{pmatrix} = \sum_{k=1}^n a_{ik} b_{kj}.$$

The square matrix I_n of dimension n containing 1 on the diagonal and 0 everywhere else is called the *identity matrix*. It is denoted as

$$\begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix}$$

Given an $m \times n$ matrix $A = (a_{ij})$, its *transpose* $A^\top = (a_{ji}^\top)$, is the $n \times m$ -matrix such that $a_{ji}^\top = a_{ij}$, for all i , $1 \leq i \leq m$, and all j , $1 \leq j \leq n$.

The transpose of a matrix A is sometimes denoted by A^t , or even by tA . Note that the transpose A^\top of a matrix A has the property that the j -th row of A^\top is the j -th column of A . In other words, transposition exchanges the rows and the columns of a matrix.

The following observation will be useful later on when we discuss the SVD. Given any $m \times n$ matrix A and any $n \times p$ matrix B , if we denote the columns of A by A^1, \dots, A^n and the rows of B by B_1, \dots, B_n , then we have

$$AB = A^1 B_1 + \dots + A^n B_n.$$

For every square matrix A of dimension n , it is immediately verified that $AI_n = I_n A = A$. If a matrix B such that $AB = BA = I_n$ exists, then it is unique, and it is called the *inverse* of A . The matrix B is also denoted by A^{-1} . An invertible matrix is also called a *nonsingular* matrix, and a matrix that is not invertible is called a *singular* matrix.

Proposition 2.19 shows that if a square matrix A has a left inverse, that is a matrix B such that $BA = I$, or a right inverse, that is a matrix C such that $AC = I$, then A is actually invertible; so $B = A^{-1}$ and $C = A^{-1}$. These facts also follow from Proposition 4.14.

It is immediately verified that the set $M_{m,n}(K)$ of $m \times n$ matrices is a *vector space* under addition of matrices and multiplication of a matrix by a scalar. Consider the $m \times n$ -matrices $E_{i,j} = (e_{hk})$, defined such that $e_{ij} = 1$, and $e_{hk} = 0$, if $h \neq i$ or $k \neq j$. It is clear that every matrix $A = (a_{ij}) \in M_{m,n}(K)$ can be written in a unique way as

$$A = \sum_{i=1}^m \sum_{j=1}^n a_{ij} E_{i,j}.$$

Thus, the family $(E_{i,j})_{1 \leq i \leq m, 1 \leq j \leq n}$ is a basis of the vector space $M_{m,n}(K)$, which has dimension mn .

Remark: Definition 3.1 and Definition 3.2 also make perfect sense when K is a (commutative) ring rather than a field. In this more general setting, the framework of vector spaces is too narrow, but we can consider structures over a commutative ring A satisfying all the axioms of Definition 2.9. Such structures are called *modules*. The theory of modules is (much) more complicated than that of vector spaces. For example, modules do not always have a basis, and other properties holding for vector spaces usually fail for modules. When a module has a basis, it is called a *free module*. For example, when A is a commutative ring, the structure A^n is a module such that the vectors e_i , with $(e_i)_i = 1$ and $(e_i)_j = 0$ for $j \neq i$, form a basis of A^n . Many properties of vector spaces still hold for A^n . Thus, A^n is a free module. As another example, when A is a commutative ring, $M_{m,n}(A)$ is a free module with basis $(E_{i,j})_{1 \leq i \leq m, 1 \leq j \leq n}$. Polynomials over a commutative ring also form a free module of infinite dimension.

Square matrices provide a natural example of a noncommutative ring with zero divisors.

Example 3.1. For example, letting A, B be the 2×2 -matrices

$$A = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}, \quad B = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix},$$

then

$$AB = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix},$$

and

$$BA = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}.$$

We now formalize the representation of linear maps by matrices.

Definition 3.3. Let E and F be two vector spaces, and let (u_1, \dots, u_n) be a basis for E , and (v_1, \dots, v_m) be a basis for F . Each vector $x \in E$ expressed in the basis (u_1, \dots, u_n) as $x = x_1 u_1 + \dots + x_n u_n$ is represented by the column matrix

$$M(x) = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

and similarly for each vector $y \in F$ expressed in the basis (v_1, \dots, v_m) .

Every linear map $f: E \rightarrow F$ is represented by the matrix $M(f) = (a_{ij})$, where a_{ij} is the i -th component of the vector $f(u_j)$ over the basis (v_1, \dots, v_m) , i.e., where

$$f(u_j) = \sum_{i=1}^m a_{ij} v_i, \quad \text{for every } j, 1 \leq j \leq n.$$

The coefficients $a_{1j}, a_{2j}, \dots, a_{mj}$ of $f(u_j)$ over the basis (v_1, \dots, v_m) form the j th column of the matrix $M(f)$ shown below:

$$\begin{array}{cccc} & f(u_1) & f(u_2) & \dots & f(u_n) \\ \begin{array}{c} v_1 \\ v_2 \\ \vdots \\ v_m \end{array} & \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix} \end{array}.$$

The matrix $M(f)$ associated with the linear map $f: E \rightarrow F$ is called the *matrix of f with respect to the bases (u_1, \dots, u_n) and (v_1, \dots, v_m)* . When $E = F$ and the basis (v_1, \dots, v_m) is identical to the basis (u_1, \dots, u_n) of E , the matrix $M(f)$ associated with $f: E \rightarrow E$ (as above) is called the *matrix of f with respect to the base (u_1, \dots, u_n)* .

Remark: As in the remark after Definition 3.1, there is no reason to assume that the vectors in the bases (u_1, \dots, u_n) and (v_1, \dots, v_m) are ordered in any particular way. However, it is often convenient to assume the natural ordering. When this is so, authors sometimes refer

to the matrix $M(f)$ as the matrix of f with respect to the *ordered bases* (u_1, \dots, u_n) and (v_1, \dots, v_m) .

Then, given a linear map $f: E \rightarrow F$ represented by the matrix $M(f) = (a_{ij})$ w.r.t. the bases (u_1, \dots, u_n) and (v_1, \dots, v_m) , by equations (1) and the definition of matrix multiplication, the equation $y = f(x)$ correspond to the matrix equation $M(y) = M(f)M(x)$, that is,

$$\begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}.$$

Recall that

$$\begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = x_1 \begin{pmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{m1} \end{pmatrix} + x_2 \begin{pmatrix} a_{12} \\ a_{22} \\ \vdots \\ a_{m2} \end{pmatrix} + \dots + x_n \begin{pmatrix} a_{1n} \\ a_{2n} \\ \vdots \\ a_{mn} \end{pmatrix}.$$

Sometimes, it is necessary to incorporate the bases (u_1, \dots, u_n) and (v_1, \dots, v_m) in the notation for the matrix $M(f)$ expressing f with respect to these bases. This turns out to be a messy enterprise!

We propose the following course of action: write $\mathcal{U} = (u_1, \dots, u_n)$ and $\mathcal{V} = (v_1, \dots, v_m)$ for the bases of E and F , and denote by $M_{\mathcal{U}, \mathcal{V}}(f)$ the *matrix of f with respect to the bases \mathcal{U} and \mathcal{V}* . Furthermore, write $x_{\mathcal{U}}$ for the coordinates $M(x) = (x_1, \dots, x_n)$ of $x \in E$ w.r.t. the basis \mathcal{U} and write $y_{\mathcal{V}}$ for the coordinates $M(y) = (y_1, \dots, y_m)$ of $y \in F$ w.r.t. the basis \mathcal{V} . Then,

$$y = f(x)$$

is expressed in matrix form by

$$y_{\mathcal{V}} = M_{\mathcal{U}, \mathcal{V}}(f) x_{\mathcal{U}}.$$

When $\mathcal{U} = \mathcal{V}$, we abbreviate $M_{\mathcal{U}, \mathcal{V}}(f)$ as $M_{\mathcal{U}}(f)$.

The above notation seems reasonable, but it has the slight disadvantage that in the expression $M_{\mathcal{U}, \mathcal{V}}(f)x_{\mathcal{U}}$, the input argument $x_{\mathcal{U}}$ which is fed to the matrix $M_{\mathcal{U}, \mathcal{V}}(f)$ does not appear next to the subscript \mathcal{U} in $M_{\mathcal{U}, \mathcal{V}}(f)$. We could have used the notation $M_{\mathcal{V}, \mathcal{U}}(f)$, and some people do that. But then, we find a bit confusing that \mathcal{V} comes before \mathcal{U} when f maps from the space E with the basis \mathcal{U} to the space F with the basis \mathcal{V} . So, we prefer to use the notation $M_{\mathcal{U}, \mathcal{V}}(f)$.

Be aware that other authors such as Meyer [80] use the notation $[f]_{\mathcal{U}, \mathcal{V}}$, and others such as Dummit and Foote [32] use the notation $M_{\mathcal{U}}^{\mathcal{V}}(f)$, instead of $M_{\mathcal{U}, \mathcal{V}}(f)$. This gets worse! You may find the notation $M_{\mathcal{V}}^{\mathcal{U}}(f)$ (as in Lang [67]), or ${}_{\mathcal{U}}[f]_{\mathcal{V}}$, or other strange notations.

Let us illustrate the representation of a linear map by a matrix in a concrete situation. Let E be the vector space $\mathbb{R}[X]_4$ of polynomials of degree at most 4, let F be the vector

space $\mathbb{R}[X]_3$ of polynomials of degree at most 3, and let the linear map be the derivative map d : that is,

$$\begin{aligned}d(P + Q) &= dP + dQ \\d(\lambda P) &= \lambda dP,\end{aligned}$$

with $\lambda \in \mathbb{R}$. We choose $(1, x, x^2, x^3, x^4)$ as a basis of E and $(1, x, x^2, x^3)$ as a basis of F . Then, the 4×5 matrix D associated with d is obtained by expressing the derivative dx^i of each basis vector for $i = 0, 1, 2, 3, 4$ over the basis $(1, x, x^2, x^3)$. We find

$$D = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 4 \end{pmatrix}.$$

Then, if P denotes the polynomial

$$P = 3x^4 - 5x^3 + x^2 - 7x + 5,$$

we have

$$dP = 12x^3 - 15x^2 + 2x - 7,$$

the polynomial P is represented by the vector $(5, -7, 1, -5, 3)$ and dP is represented by the vector $(-7, 2, -15, 12)$, and we have

$$\begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 4 \end{pmatrix} \begin{pmatrix} 5 \\ -7 \\ 1 \\ -5 \\ 3 \end{pmatrix} = \begin{pmatrix} -7 \\ 2 \\ -15 \\ 12 \end{pmatrix},$$

as expected! The kernel (nullspace) of d consists of the polynomials of degree 0, that is, the constant polynomials. Therefore $\dim(\text{Ker } d) = 1$, and from

$$\dim(E) = \dim(\text{Ker } d) + \dim(\text{Im } d)$$

(see Theorem 4.11), we get $\dim(\text{Im } d) = 4$ (since $\dim(E) = 5$).

For fun, let us figure out the linear map from the vector space $\mathbb{R}[X]_3$ to the vector space $\mathbb{R}[X]_4$ given by integration (finding the primitive, or anti-derivative) of x^i , for $i = 0, 1, 2, 3$. The 5×4 matrix S representing \int with respect to the same bases as before is

$$S = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1/2 & 0 & 0 \\ 0 & 0 & 1/3 & 0 \\ 0 & 0 & 0 & 1/4 \end{pmatrix}.$$

We verify that $DS = I_4$,

$$\begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 4 \end{pmatrix} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1/2 & 0 & 0 \\ 0 & 0 & 1/3 & 0 \\ 0 & 0 & 0 & 1/4 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix},$$

as it should! The equation $DS = I_4$ show that S is injective and has D as a left inverse. However, $SD \neq I_5$, and instead

$$\begin{pmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1/2 & 0 & 0 \\ 0 & 0 & 1/3 & 0 \\ 0 & 0 & 0 & 1/4 \end{pmatrix} \begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 4 \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix},$$

because constant polynomials (polynomials of degree 0) belong to the kernel of D .

The function that associates to a linear map $f: E \rightarrow F$ the matrix $M(f)$ w.r.t. the bases (u_1, \dots, u_n) and (v_1, \dots, v_m) has the property that matrix multiplication corresponds to composition of linear maps. This allows us to transfer properties of linear maps to matrices. Here is an illustration of this technique:

Proposition 3.1. (1) Given any matrices $A \in M_{m,n}(K)$, $B \in M_{n,p}(K)$, and $C \in M_{p,q}(K)$, we have

$$(AB)C = A(BC);$$

that is, matrix multiplication is associative.

(2) Given any matrices $A, B \in M_{m,n}(K)$, and $C, D \in M_{n,p}(K)$, for all $\lambda \in K$, we have

$$\begin{aligned} (A + B)C &= AC + BC \\ A(C + D) &= AC + AD \\ (\lambda A)C &= \lambda(AC) \\ A(\lambda C) &= \lambda(AC), \end{aligned}$$

so that matrix multiplication $\cdot: M_{m,n}(K) \times M_{n,p}(K) \rightarrow M_{m,p}(K)$ is bilinear.

Proof. (1) Every $m \times n$ matrix $A = (a_{ij})$ defines the function $f_A: K^n \rightarrow K^m$ given by

$$f_A(x) = Ax,$$

for all $x \in K^n$. It is immediately verified that f_A is linear and that the matrix $M(f_A)$ representing f_A over the canonical bases in K^n and K^m is equal to A . Then, formula (4) proves that

$$M(f_A \circ f_B) = M(f_A)M(f_B) = AB,$$

so we get

$$M((f_A \circ f_B) \circ f_C) = M(f_A \circ f_B)M(f_C) = (AB)C$$

and

$$M(f_A \circ (f_B \circ f_C)) = M(f_A)M(f_B \circ f_C) = A(BC),$$

and since composition of functions is associative, we have $(f_A \circ f_B) \circ f_C = f_A \circ (f_B \circ f_C)$, which implies that

$$(AB)C = A(BC).$$

(2) It is immediately verified that if $f_1, f_2 \in \text{Hom}_K(E, F)$, $A, B \in M_{m,n}(K)$, (u_1, \dots, u_n) is any basis of E , and (v_1, \dots, v_m) is any basis of F , then

$$\begin{aligned} M(f_1 + f_2) &= M(f_1) + M(f_2) \\ f_{A+B} &= f_A + f_B. \end{aligned}$$

Then we have

$$\begin{aligned} (A + B)C &= M(f_{A+B})M(f_C) \\ &= M(f_{A+B} \circ f_C) \\ &= M((f_A + f_B) \circ f_C) \\ &= M((f_A \circ f_C) + (f_B \circ f_C)) \\ &= M(f_A \circ f_C) + M(f_B \circ f_C) \\ &= M(f_A)M(f_C) + M(f_B)M(f_C) \\ &= AC + BC. \end{aligned}$$

The equation $A(C + D) = AC + AD$ is proved in a similar fashion, and the last two equations are easily verified. We could also have verified all the identities by making matrix computations. \square

Note that Proposition 3.1 implies that the vector space $M_n(K)$ of square matrices is a (noncommutative) ring with unit I_n . (It even shows that $M_n(K)$ is an associative *algebra*.)

The following proposition states the main properties of the mapping $f \mapsto M(f)$ between $\text{Hom}(E, F)$ and $M_{m,n}$. In short, it is an isomorphism of vector spaces.

Proposition 3.2. *Given three vector spaces E, F, G , with respective bases (u_1, \dots, u_p) , (v_1, \dots, v_n) , and (w_1, \dots, w_m) , the mapping $M: \text{Hom}(E, F) \rightarrow M_{n,p}$ that associates the matrix $M(g)$ to a linear map $g: E \rightarrow F$ satisfies the following properties for all $x \in E$, all $g, h: E \rightarrow F$, and all $f: F \rightarrow G$:*

$$\begin{aligned} M(g(x)) &= M(g)M(x) \\ M(g + h) &= M(g) + M(h) \\ M(\lambda g) &= \lambda M(g) \\ M(f \circ g) &= M(f)M(g), \end{aligned}$$

where $M(x)$ is the column vector associated with the vector x and $M(g(x))$ is the column vector associated with $g(x)$, as explained in Definition 3.3.

Thus, $M: \text{Hom}(E, F) \rightarrow M_{n,p}$ is an isomorphism of vector spaces, and when $p = n$ and the basis (v_1, \dots, v_n) is identical to the basis (u_1, \dots, u_p) , $M: \text{Hom}(E, E) \rightarrow M_n$ is an isomorphism of rings.

Proof. That $M(g(x)) = M(g)M(x)$ was shown just before stating the proposition, using identity (1). The identities $M(g + h) = M(g) + M(h)$ and $M(\lambda g) = \lambda M(g)$ are straightforward, and $M(f \circ g) = M(f)M(g)$ follows from (4) and the definition of matrix multiplication. The mapping $M: \text{Hom}(E, F) \rightarrow M_{n,p}$ is clearly injective, and since every matrix defines a linear map, it is also surjective, and thus bijective. In view of the above identities, it is an isomorphism (and similarly for $M: \text{Hom}(E, E) \rightarrow M_n$). \square

In view of Proposition 3.2, it seems preferable to represent vectors from a vector space of finite dimension as column vectors rather than row vectors. Thus, from now on, we will denote vectors of \mathbb{R}^n (or more generally, of K^n) as column vectors.

It is important to observe that the isomorphism $M: \text{Hom}(E, F) \rightarrow M_{n,p}$ given by Proposition 3.2 depends on the choice of the bases (u_1, \dots, u_p) and (v_1, \dots, v_n) , and similarly for the isomorphism $M: \text{Hom}(E, E) \rightarrow M_n$, which depends on the choice of the basis (u_1, \dots, u_n) . Thus, it would be useful to know how a change of basis affects the representation of a linear map $f: E \rightarrow F$ as a matrix. The following simple proposition is needed.

Proposition 3.3. *Let E be a vector space, and let (u_1, \dots, u_n) be a basis of E . For every family (v_1, \dots, v_n) , let $P = (a_{ij})$ be the matrix defined such that $v_j = \sum_{i=1}^n a_{ij}u_i$. The matrix P is invertible iff (v_1, \dots, v_n) is a basis of E .*

Proof. Note that we have $P = M(f)$, the matrix associated with the unique linear map $f: E \rightarrow E$ such that $f(u_i) = v_i$. By Proposition 2.16, f is bijective iff (v_1, \dots, v_n) is a basis of E . Furthermore, it is obvious that the identity matrix I_n is the matrix associated with the identity $\text{id}: E \rightarrow E$ w.r.t. any basis. If f is an isomorphism, then $f \circ f^{-1} = f^{-1} \circ f = \text{id}$, and by Proposition 3.2, we get $M(f)M(f^{-1}) = M(f^{-1})M(f) = I_n$, showing that P is invertible and that $M(f^{-1}) = P^{-1}$. \square

Proposition 3.3 suggests the following definition.

Definition 3.4. Given a vector space E of dimension n , for any two bases (u_1, \dots, u_n) and (v_1, \dots, v_n) of E , let $P = (a_{ij})$ be the invertible matrix defined such that

$$v_j = \sum_{i=1}^n a_{ij}u_i,$$

which is also the matrix of the identity $\text{id}: E \rightarrow E$ with respect to the bases (v_1, \dots, v_n) and (u_1, \dots, u_n) , in that order. Indeed, we express each $\text{id}(v_j) = v_j$ over the basis (u_1, \dots, u_n) .

The coefficients $a_{1j}, a_{2j}, \dots, a_{nj}$ of v_j over the basis (u_1, \dots, u_n) form the j th column of the matrix P shown below:

$$\begin{array}{cccc} & v_1 & v_2 & \dots & v_n \\ \begin{array}{c} u_1 \\ u_2 \\ \vdots \\ u_n \end{array} & \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix} \end{array}.$$

The matrix P is called the *change of basis matrix* from (u_1, \dots, u_n) to (v_1, \dots, v_n) .

Clearly, the change of basis matrix from (v_1, \dots, v_n) to (u_1, \dots, u_n) is P^{-1} . Since $P = (a_{i,j})$ is the matrix of the identity $\text{id}: E \rightarrow E$ with respect to the bases (v_1, \dots, v_n) and (u_1, \dots, u_n) , given any vector $x \in E$, if $x = x_1 u_1 + \dots + x_n u_n$ over the basis (u_1, \dots, u_n) and $x = x'_1 v_1 + \dots + x'_n v_n$ over the basis (v_1, \dots, v_n) , from Proposition 3.2, we have

$$\begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix} \begin{pmatrix} x'_1 \\ \vdots \\ x'_n \end{pmatrix},$$

showing that the *old* coordinates (x_i) of x (over (u_1, \dots, u_n)) are expressed in terms of the *new* coordinates (x'_i) of x (over (v_1, \dots, v_n)).

Now we face the painful task of assigning a “good” notation incorporating the bases $\mathcal{U} = (u_1, \dots, u_n)$ and $\mathcal{V} = (v_1, \dots, v_n)$ into the notation for the change of basis matrix from \mathcal{U} to \mathcal{V} . Because the change of basis matrix from \mathcal{U} to \mathcal{V} is the matrix of the identity map id_E with respect to the bases \mathcal{V} and \mathcal{U} in that order, we could denote it by $M_{\mathcal{V},\mathcal{U}}(\text{id})$ (Meyer [80] uses the notation $[I]_{\mathcal{V},\mathcal{U}}$), which we abbreviate as

$$P_{\mathcal{V},\mathcal{U}}.$$

Note that

$$P_{\mathcal{U},\mathcal{V}} = P_{\mathcal{V},\mathcal{U}}^{-1}.$$

Then, if we write $x_{\mathcal{U}} = (x_1, \dots, x_n)$ for the old coordinates of x with respect to the basis \mathcal{U} and $x_{\mathcal{V}} = (x'_1, \dots, x'_n)$ for the new coordinates of x with respect to the basis \mathcal{V} , we have

$$x_{\mathcal{U}} = P_{\mathcal{V},\mathcal{U}} x_{\mathcal{V}}, \quad x_{\mathcal{V}} = P_{\mathcal{V},\mathcal{U}}^{-1} x_{\mathcal{U}}.$$

The above may look backward, but remember that the matrix $M_{\mathcal{U},\mathcal{V}}(f)$ takes input expressed over the basis \mathcal{U} to output expressed over the basis \mathcal{V} . Consequently, $P_{\mathcal{V},\mathcal{U}}$ takes input expressed over the basis \mathcal{V} to output expressed over the basis \mathcal{U} , and $x_{\mathcal{U}} = P_{\mathcal{V},\mathcal{U}} x_{\mathcal{V}}$ matches this point of view!



Beware that some authors (such as Artin [4]) define the change of basis matrix from \mathcal{U} to \mathcal{V} as $P_{\mathcal{U},\mathcal{V}} = P_{\mathcal{V},\mathcal{U}}^{-1}$. Under this point of view, the old basis \mathcal{U} is expressed in terms of the new basis \mathcal{V} . We find this a bit unnatural. Also, in practice, it seems that the new basis is often expressed in terms of the old basis, rather than the other way around.

Since the matrix $P = P_{\mathcal{V},\mathcal{U}}$ expresses the *new* basis (v_1, \dots, v_n) in terms of the *old* basis (u_1, \dots, u_n) , we observe that the coordinates (x_i) of a vector x vary in the *opposite direction* of the change of basis. For this reason, vectors are sometimes said to be *contravariant*. However, this expression does not make sense! Indeed, a vector in an intrinsic quantity that does not depend on a specific basis. What makes sense is that the *coordinates* of a vector vary in a contravariant fashion.

Let us consider some concrete examples of change of bases.

Example 3.2. Let $E = F = \mathbb{R}^2$, with $u_1 = (1, 0)$, $u_2 = (0, 1)$, $v_1 = (1, 1)$ and $v_2 = (-1, 1)$. The change of basis matrix P from the basis $\mathcal{U} = (u_1, u_2)$ to the basis $\mathcal{V} = (v_1, v_2)$ is

$$P = \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix}$$

and its inverse is

$$P^{-1} = \begin{pmatrix} 1/2 & 1/2 \\ -1/2 & 1/2 \end{pmatrix}.$$

The old coordinates (x_1, x_2) with respect to (u_1, u_2) are expressed in terms of the new coordinates (x'_1, x'_2) with respect to (v_1, v_2) by

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x'_1 \\ x'_2 \end{pmatrix},$$

and the new coordinates (x'_1, x'_2) with respect to (v_1, v_2) are expressed in terms of the old coordinates (x_1, x_2) with respect to (u_1, u_2) by

$$\begin{pmatrix} x'_1 \\ x'_2 \end{pmatrix} = \begin{pmatrix} 1/2 & 1/2 \\ -1/2 & 1/2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}.$$

Example 3.3. Let $E = F = \mathbb{R}[X]_3$ be the set of polynomials of degree at most 3, and consider the bases $\mathcal{U} = (1, x, x^2, x^3)$ and $\mathcal{V} = (B_0^3(x), B_1^3(x), B_2^3(x), B_3^3(x))$, where $B_0^3(x), B_1^3(x), B_2^3(x), B_3^3(x)$ are the *Bernstein polynomials* of degree 3, given by

$$B_0^3(x) = (1-x)^3 \quad B_1^3(x) = 3(1-x)^2x \quad B_2^3(x) = 3(1-x)x^2 \quad B_3^3(x) = x^3.$$

By expanding the Bernstein polynomials, we find that the change of basis matrix $P_{\mathcal{V},\mathcal{U}}$ is given by

$$P_{\mathcal{V},\mathcal{U}} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ -3 & 3 & 0 & 0 \\ 3 & -6 & 3 & 0 \\ -1 & 3 & -3 & 1 \end{pmatrix}.$$

We also find that the inverse of $P_{\mathcal{V},\mathcal{U}}$ is

$$P_{\mathcal{V},\mathcal{U}}^{-1} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 1/3 & 0 & 0 \\ 1 & 2/3 & 1/3 & 0 \\ 1 & 1 & 1 & 1 \end{pmatrix}.$$

Therefore, the coordinates of the polynomial $2x^3 - x + 1$ over the basis \mathcal{V} are

$$\begin{pmatrix} 1 \\ 2/3 \\ 1/3 \\ 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 1/3 & 0 & 0 \\ 1 & 2/3 & 1/3 & 0 \\ 1 & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \\ 0 \\ 2 \end{pmatrix},$$

and so

$$2x^3 - x + 1 = B_0^3(x) + \frac{2}{3}B_1^3(x) + \frac{1}{3}B_2^3(x) + 2B_3^3(x).$$

Our next example is the Haar wavelets, a fundamental tool in signal processing.

3.2 Haar Basis Vectors and a Glimpse at Wavelets

We begin by considering *Haar wavelets* in \mathbb{R}^4 . Wavelets play an important role in audio and video signal processing, especially for *compressing* long signals into much smaller ones than still retain enough information so that when they are played, we can't see or hear any difference.

Consider the four vectors w_1, w_2, w_3, w_4 given by

$$w_1 = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} \quad w_2 = \begin{pmatrix} 1 \\ 1 \\ -1 \\ -1 \end{pmatrix} \quad w_3 = \begin{pmatrix} 1 \\ -1 \\ 0 \\ 0 \end{pmatrix} \quad w_4 = \begin{pmatrix} 0 \\ 0 \\ 1 \\ -1 \end{pmatrix}.$$

Note that these vectors are pairwise orthogonal, so they are indeed linearly independent (we will see this in a later chapter). Let $\mathcal{W} = \{w_1, w_2, w_3, w_4\}$ be the *Haar basis*, and let $\mathcal{U} = \{e_1, e_2, e_3, e_4\}$ be the canonical basis of \mathbb{R}^4 . The change of basis matrix $W = P_{\mathcal{W},\mathcal{U}}$ from \mathcal{U} to \mathcal{W} is given by

$$W = \begin{pmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & -1 & 0 \\ 1 & -1 & 0 & 1 \\ 1 & -1 & 0 & -1 \end{pmatrix},$$

and we easily find that the inverse of W is given by

$$W^{-1} = \begin{pmatrix} 1/4 & 0 & 0 & 0 \\ 0 & 1/4 & 0 & 0 \\ 0 & 0 & 1/2 & 0 \\ 0 & 0 & 0 & 1/2 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{pmatrix}.$$

So, the vector $v = (6, 4, 5, 1)$ over the basis \mathcal{U} becomes $c = (c_1, c_2, c_3, c_4)$ over the Haar basis \mathcal{W} , with

$$\begin{pmatrix} c_1 \\ c_2 \\ c_3 \\ c_4 \end{pmatrix} = \begin{pmatrix} 1/4 & 0 & 0 & 0 \\ 0 & 1/4 & 0 & 0 \\ 0 & 0 & 1/2 & 0 \\ 0 & 0 & 0 & 1/2 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{pmatrix} \begin{pmatrix} 6 \\ 4 \\ 5 \\ 1 \end{pmatrix} = \begin{pmatrix} 4 \\ 1 \\ 1 \\ 2 \end{pmatrix}.$$

Given a signal $v = (v_1, v_2, v_3, v_4)$, we first *transform* v into its coefficients $c = (c_1, c_2, c_3, c_4)$ over the Haar basis by computing $c = W^{-1}v$. Observe that

$$c_1 = \frac{v_1 + v_2 + v_3 + v_4}{4}$$

is the overall *average* value of the signal v . The coefficient c_1 corresponds to the background of the image (or of the sound). Then, c_2 gives the coarse details of v , whereas, c_3 gives the details in the first part of v , and c_4 gives the details in the second half of v .

Reconstruction of the signal consists in computing $v = Wc$. The trick for good *compression* is to throw away some of the coefficients of c (set them to zero), obtaining a *compressed signal* \hat{c} , and still retain enough crucial information so that the reconstructed signal $\hat{v} = W\hat{c}$ looks almost as good as the original signal v . Thus, the steps are:

$$\text{input } v \longrightarrow \text{coefficients } c = W^{-1}v \longrightarrow \text{compressed } \hat{c} \longrightarrow \text{compressed } \hat{v} = W\hat{c}.$$

This kind of compression scheme makes modern video conferencing possible.

It turns out that there is a faster way to find $c = W^{-1}v$, without actually using W^{-1} . This has to do with the multiscale nature of Haar wavelets.

Given the original signal $v = (6, 4, 5, 1)$ shown in Figure 3.1, we compute averages and half differences obtaining Figure 3.2. We get the coefficients $c_3 = 1$ and $c_4 = 2$. Then, again we compute averages and half differences obtaining Figure 3.3. We get the coefficients $c_1 = 4$ and $c_2 = 1$.

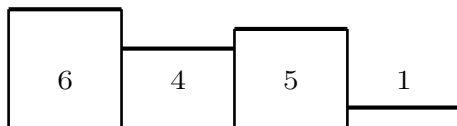


Figure 3.1: The original signal v



Figure 3.2: First averages and first half differences



Figure 3.3: Second averages and second half differences

Note that the original signal v can be reconstructed from the two signals in Figure 3.2, and the signal on the left of Figure 3.2 can be reconstructed from the two signals in Figure 3.3.

This method can be generalized to signals of any length 2^n . The previous case corresponds to $n = 2$. Let us consider the case $n = 3$. The Haar basis $(w_1, w_2, w_3, w_4, w_5, w_6, w_7, w_8)$ is given by the matrix

$$W = \begin{pmatrix} 1 & 1 & 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & -1 & 0 & 0 & 0 \\ 1 & 1 & -1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & -1 & 0 & 0 & -1 & 0 & 0 \\ 1 & -1 & 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & -1 & 0 & 1 & 0 & 0 & -1 & 0 \\ 1 & -1 & 0 & -1 & 0 & 0 & 0 & 1 \\ 1 & -1 & 0 & -1 & 0 & 0 & 0 & -1 \end{pmatrix}$$

The columns of this matrix are orthogonal and it is easy to see that

$$W^{-1} = \text{diag}(1/8, 1/8, 1/4, 1/4, 1/2, 1/2, 1/2, 1/2)W^{\top}.$$

A pattern is beginning to emerge. It looks like the second Haar basis vector w_2 is the “mother” of all the other basis vectors, except the first, whose purpose is to perform averaging. Indeed, in general, given

$$w_2 = (\underbrace{1, \dots, 1, -1, \dots, -1}_{2^n}),$$

the other Haar basis vectors are obtained by a “scaling and shifting process.” Starting from w_2 , the scaling process generates the vectors

$$w_3, w_5, w_9, \dots, w_{2^j+1}, \dots, w_{2^{n-1}+1},$$

such that $w_{2^{j+1}+1}$ is obtained from w_{2^j+1} by forming two consecutive blocks of 1 and -1 of half the size of the blocks in w_{2^j+1} , and setting all other entries to zero. Observe that w_{2^j+1} has 2^j blocks of 2^{n-j} elements. The shifting process, consists in shifting the blocks of 1 and -1 in w_{2^j+1} to the right by inserting a block of $(k-1)2^{n-j}$ zeros from the left, with $0 \leq j \leq n-1$ and $1 \leq k \leq 2^j$. Thus, we obtain the following formula for w_{2^j+k} :

$$w_{2^j+k}(i) = \begin{cases} 0 & 1 \leq i \leq (k-1)2^{n-j} \\ 1 & (k-1)2^{n-j} + 1 \leq i \leq (k-1)2^{n-j} + 2^{n-j-1} \\ -1 & (k-1)2^{n-j} + 2^{n-j-1} + 1 \leq i \leq k2^{n-j} \\ 0 & k2^{n-j} + 1 \leq i \leq 2^n, \end{cases}$$

with $0 \leq j \leq n-1$ and $1 \leq k \leq 2^j$. Of course

$$w_1 = \underbrace{(1, \dots, 1)}_{2^n}.$$

The above formulae look a little better if we change our indexing slightly by letting k vary from 0 to $2^j - 1$ and using the index j instead of 2^j . In this case, the Haar basis is denoted by

$$w_1, h_0^0, h_0^1, h_1^1, h_0^2, h_1^2, h_2^2, h_3^2, \dots, h_k^j, \dots, h_{2^{n-1}-1}^{n-1},$$

and

$$h_k^j(i) = \begin{cases} 0 & 1 \leq i \leq k2^{n-j} \\ 1 & k2^{n-j} + 1 \leq i \leq k2^{n-j} + 2^{n-j-1} \\ -1 & k2^{n-j} + 2^{n-j-1} + 1 \leq i \leq (k+1)2^{n-j} \\ 0 & (k+1)2^{n-j} + 1 \leq i \leq 2^n, \end{cases}$$

with $0 \leq j \leq n-1$ and $0 \leq k \leq 2^j - 1$.

It turns out that there is a way to understand these formulae better if we interpret a vector $u = (u_1, \dots, u_m)$ as a piecewise linear function over the interval $[0, 1)$. We define the function $\text{plf}(u)$ such that

$$\text{plf}(u)(x) = u_i, \quad \frac{i-1}{m} \leq x < \frac{i}{m}, \quad 1 \leq i \leq m.$$

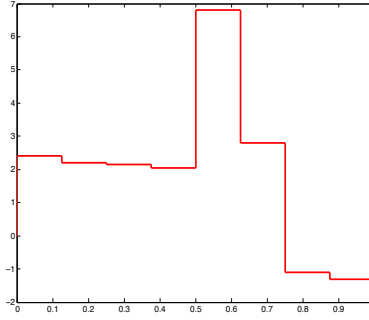
In words, the function $\text{plf}(u)$ has the value u_1 on the interval $[0, 1/m)$, the value u_2 on $[1/m, 2/m)$, etc., and the value u_m on the interval $[(m-1)/m, 1)$. For example, the piecewise linear function associated with the vector

$$u = (2.4, 2.2, 2.15, 2.05, 6.8, 2.8, -1.1, -1.3)$$

is shown in Figure 3.4.

Then, each basis vector h_k^j corresponds to the function

$$\psi_k^j = \text{plf}(h_k^j).$$

Figure 3.4: The piecewise linear function $\text{plf}(u)$

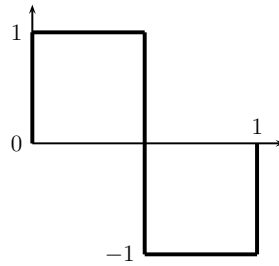
In particular, for all n , the Haar basis vectors

$$h_0^0 = w_2 = \underbrace{(1, \dots, 1, -1, \dots, -1)}_{2^n}$$

yield the same piecewise linear function ψ given by

$$\psi(x) = \begin{cases} 1 & \text{if } 0 \leq x < 1/2 \\ -1 & \text{if } 1/2 \leq x < 1 \\ 0 & \text{otherwise,} \end{cases}$$

whose graph is shown in Figure 3.5. Then, it is easy to see that ψ_k^j is given by the simple

Figure 3.5: The Haar wavelet ψ

expression

$$\psi_k^j(x) = \psi(2^j x - k), \quad 0 \leq j \leq n-1, \quad 0 \leq k \leq 2^j - 1.$$

The above formula makes it clear that ψ_k^j is obtained from ψ by scaling and shifting. The function $\phi_0^0 = \text{plf}(w_1)$ is the piecewise linear function with the constant value 1 on $[0, 1)$, and the functions ψ_k^j together with ϕ_0^0 are known as the *Haar wavelets*.

Rather than using W^{-1} to convert a vector u to a vector c of coefficients over the Haar basis, and the matrix W to reconstruct the vector u from its Haar coefficients c , we can use faster algorithms that use averaging and differencing.

If c is a vector of Haar coefficients of dimension 2^n , we compute the sequence of vectors u^0, u^1, \dots, u^n as follows:

$$\begin{aligned} u^0 &= c \\ u^{j+1} &= u^j \\ u^{j+1}(2i-1) &= u^j(i) + u^j(2^j + i) \\ u^{j+1}(2i) &= u^j(i) - u^j(2^j + i), \end{aligned}$$

for $j = 0, \dots, n-1$ and $i = 1, \dots, 2^j$. The reconstructed vector (signal) is $u = u^n$.

If u is a vector of dimension 2^n , we compute the sequence of vectors c^n, c^{n-1}, \dots, c^0 as follows:

$$\begin{aligned} c^n &= u \\ c^j &= c^{j+1} \\ c^j(i) &= (c^{j+1}(2i-1) + c^{j+1}(2i))/2 \\ c^j(2^j + i) &= (c^{j+1}(2i-1) - c^{j+1}(2i))/2, \end{aligned}$$

for $j = n-1, \dots, 0$ and $i = 1, \dots, 2^j$. The vector over the Haar basis is $c = c^0$.

We leave it as an exercise to implement the above programs in **Matlab** using two variables u and c , and by building iteratively 2^j . Here is an example of the conversion of a vector to its Haar coefficients for $n = 3$.

Given the sequence $u = (31, 29, 23, 17, -6, -8, -2, -4)$, we get the sequence

$$\begin{aligned} c^3 &= (31, 29, 23, 17, -6, -8, -2, -4) \\ c^2 &= (30, 20, -7, -3, 1, 3, 1, 1) \\ c^1 &= (25, -5, 5, -2, 1, 3, 1, 1) \\ c^0 &= (10, 15, 5, -2, 1, 3, 1, 1), \end{aligned}$$

so $c = (10, 15, 5, -2, 1, 3, 1, 1)$. Conversely, given $c = (10, 15, 5, -2, 1, 3, 1, 1)$, we get the sequence

$$\begin{aligned} u^0 &= (10, 15, 5, -2, 1, 3, 1, 1) \\ u^1 &= (25, -5, 5, -2, 1, 3, 1, 1) \\ u^2 &= (30, 20, -7, -3, 1, 3, 1, 1) \\ u^3 &= (31, 29, 23, 17, -6, -8, -2, -4), \end{aligned}$$

which gives back $u = (31, 29, 23, 17, -6, -8, -2, -4)$.

There is another recursive method for constructing the Haar matrix W_n of dimension 2^n that makes it clearer why the above algorithms are indeed correct (which nobody seems to prove!). If we split W_n into two $2^n \times 2^{n-1}$ matrices, then the second matrix containing the last 2^{n-1} columns of W_n has a very simple structure: it consists of the vector

$$\underbrace{(1, -1, 0, \dots, 0)}_{2^n}$$

and $2^{n-1} - 1$ shifted copies of it, as illustrated below for $n = 3$:

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & -1 \end{pmatrix}.$$

Then, we form the $2^n \times 2^{n-2}$ matrix obtained by “doubling” each column of odd index, which means replacing each such column by a column in which the block of 1 is doubled and the block of -1 is doubled. In general, given a current matrix of dimension $2^n \times 2^j$, we form a $2^n \times 2^{j-1}$ matrix by doubling each column of odd index, which means that we replace each such column by a column in which the block of 1 is doubled and the block of -1 is doubled. We repeat this process $n - 1$ times until we get the vector

$$\underbrace{(1, \dots, 1, -1, \dots, -1)}_{2^n}.$$

The first vector is the averaging vector $\underbrace{(1, \dots, 1)}_{2^n}$. This process is illustrated below for $n = 3$:

$$\begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ -1 \\ -1 \\ -1 \\ -1 \end{pmatrix} \leftarrow \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ -1 & 0 \\ -1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & -1 \\ 0 & -1 \end{pmatrix} \leftarrow \begin{pmatrix} 1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & -1 \end{pmatrix}$$

Adding $\underbrace{(1, \dots, 1, 1, \dots, 1)}_{2^n}$ as the first column, we obtain

$$W_3 = \begin{pmatrix} 1 & 1 & 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & -1 & 0 & 0 & 0 \\ 1 & 1 & -1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & -1 & 0 & 0 & -1 & 0 & 0 \\ 1 & -1 & 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & -1 & 0 & 1 & 0 & 0 & -1 & 0 \\ 1 & -1 & 0 & -1 & 0 & 0 & 0 & 1 \\ 1 & -1 & 0 & -1 & 0 & 0 & 0 & -1 \end{pmatrix}.$$

Observe that the right block (of size $2^n \times 2^{n-1}$) shows clearly how the detail coefficients in the second half of the vector c are added and subtracted to the entries in the first half of the partially reconstructed vector after $n - 1$ steps.

An important and attractive feature of the Haar basis is that it provides a *multiresolution analysis* of a signal. Indeed, given a signal u , if $c = (c_1, \dots, c_{2^n})$ is the vector of its Haar coefficients, the coefficients with low index give coarse information about u , and the coefficients with high index represent fine information. For example, if u is an audio signal corresponding to a Mozart concerto played by an orchestra, c_1 corresponds to the “background noise,” c_2 to the bass, c_3 to the first cello, c_4 to the second cello, c_5, c_6, c_7, c_8 to the violas, then the violins, etc. This multiresolution feature of wavelets can be exploited to compress a signal, that is, to use fewer coefficients to represent it. Here is an example.

Consider the signal

$$u = (2.4, 2.2, 2.15, 2.05, 6.8, 2.8, -1.1, -1.3),$$

whose Haar transform is

$$c = (2, 0.2, 0.1, 3, 0.1, 0.05, 2, 0.1).$$

The piecewise-linear curves corresponding to u and c are shown in Figure 3.6. Since some of the coefficients in c are small (smaller than or equal to 0.2) we can compress c by replacing them by 0. We get

$$c_2 = (2, 0, 0, 3, 0, 0, 2, 0),$$

and the reconstructed signal is

$$u_2 = (2, 2, 2, 2, 7, 3, -1, -1).$$

The piecewise-linear curves corresponding to u_2 and c_2 are shown in Figure 3.7.

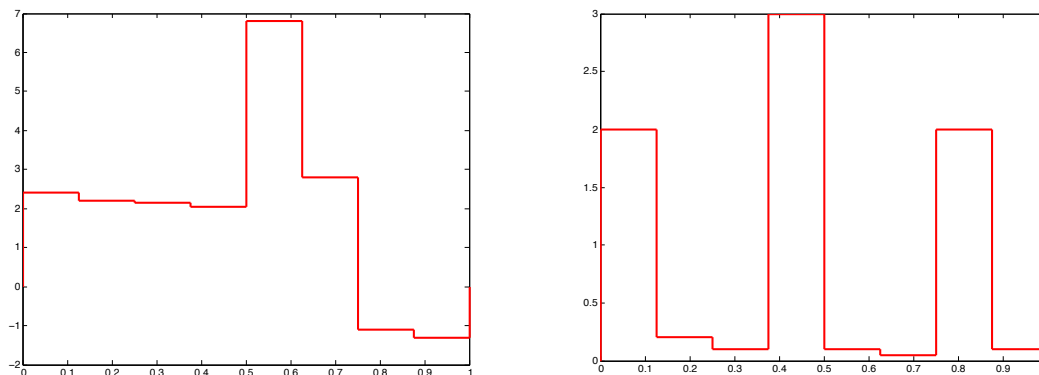


Figure 3.6: A signal and its Haar transform

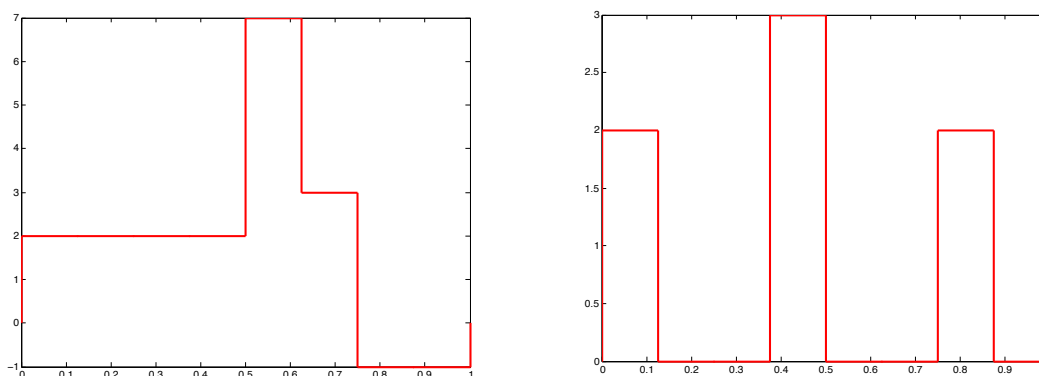


Figure 3.7: A compressed signal and its compressed Haar transform

An interesting (and amusing) application of the Haar wavelets is to the compression of audio signals. It turns out that if you type `load handel` in **Matlab** an audio file will be loaded in a vector denoted by y , and if you type `sound(y)`, the computer will play this piece of music. You can convert y to its vector of Haar coefficients, c . The length of y is 73113, so first truncate the tail of y to get a vector of length $65536 = 2^{16}$. A plot of the signals corresponding to y and c is shown in Figure 3.8. Then, run a program that sets all coefficients of c whose absolute value is less than 0.05 to zero. This sets 37272 coefficients to 0. The resulting vector c_2 is converted to a signal y_2 . A plot of the signals corresponding to y_2 and c_2 is shown in Figure 3.9. When you type `sound(y2)`, you find that the music doesn't differ much from the original, although it sounds less crisp. You should play with other numbers greater than or less than 0.05. You should hear what happens when you type `sound(c)`. It plays the music corresponding to the Haar transform c of y , and it is quite funny.

Another neat property of the Haar transform is that it can be instantly generalized to

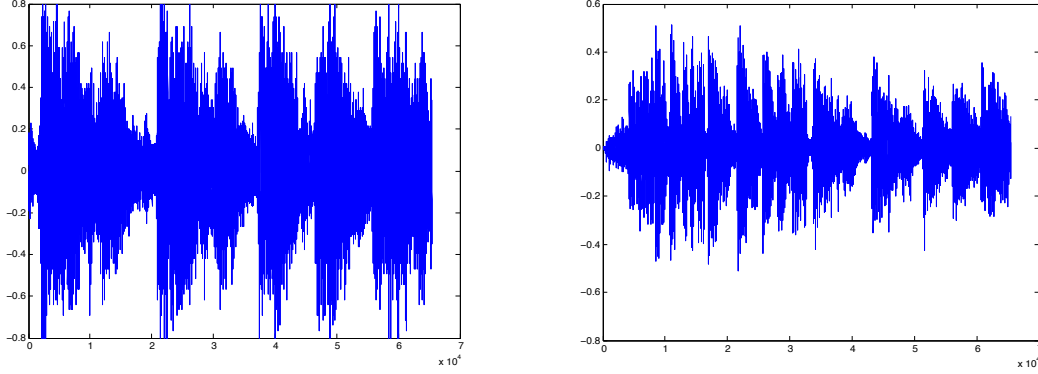


Figure 3.8: The signal “handel” and its Haar transform

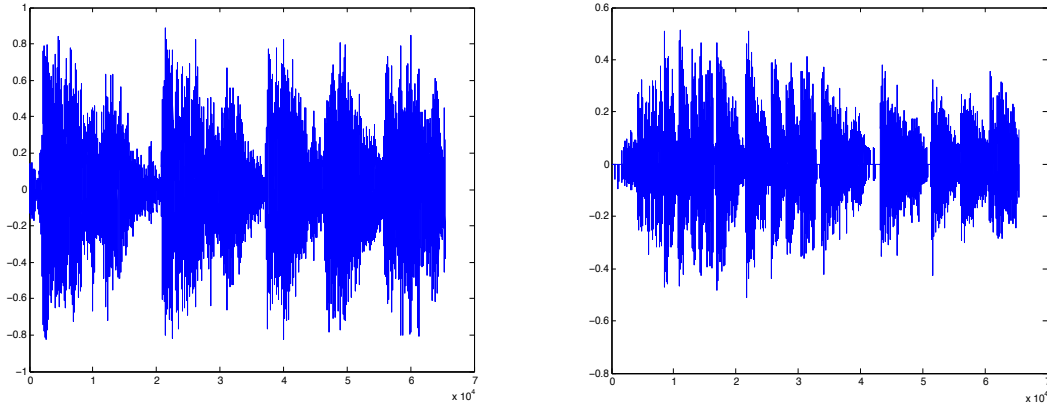


Figure 3.9: The compressed signal “handel” and its Haar transform

matrices (even rectangular) without any extra effort! This allows for the compression of digital images. But first, we address the issue of normalization of the Haar coefficients. As we observed earlier, the $2^n \times 2^n$ matrix W_n of Haar basis vectors has orthogonal columns, but its columns do not have unit length. As a consequence, W_n^\top is not the inverse of W_n , but rather the matrix

$$W_n^{-1} = D_n W_n^\top$$

with $D_n = \text{diag}\left(2^{-n}, \underbrace{2^{-n}}_{2^0}, \underbrace{2^{-(n-1)}, 2^{-(n-1)}}_{2^1}, \underbrace{2^{-(n-2)}, \dots, 2^{-(n-2)}}_{2^2}, \dots, \underbrace{2^{-1}, \dots, 2^{-1}}_{2^{n-1}}\right)$.

Therefore, we define the orthogonal matrix

$$H_n = W_n D_n^{\frac{1}{2}}$$

whose columns are the normalized Haar basis vectors, with

$$D_n^{\frac{1}{2}} = \text{diag}\left(2^{-\frac{n}{2}}, \underbrace{2^{-\frac{n}{2}}}_{2^0}, \underbrace{2^{-\frac{n-1}{2}}, 2^{-\frac{n-1}{2}}}_{2^1}, \underbrace{2^{-\frac{n-2}{2}}, \dots, 2^{-\frac{n-2}{2}}}_{2^2}, \dots, \underbrace{2^{-\frac{1}{2}}, \dots, 2^{-\frac{1}{2}}}_{2^{n-1}}\right).$$

We call H_n the *normalized Haar transform matrix*. Because H_n is orthogonal, $H_n^{-1} = H_n^\top$. Given a vector (signal) u , we call $c = H_n^\top u$ the *normalized Haar coefficients* of u . Then, a moment of reflexion shows that we have to slightly modify the algorithms to compute $H_n^\top u$ and $H_n c$ as follows: When computing the sequence of u^j s, use

$$\begin{aligned} u^{j+1}(2i-1) &= (u^j(i) + u^j(2^j+i))/\sqrt{2} \\ u^{j+1}(2i) &= (u^j(i) - u^j(2^j+i))/\sqrt{2}, \end{aligned}$$

and when computing the sequence of c^j s, use

$$\begin{aligned} c^j(i) &= (c^{j+1}(2i-1) + c^{j+1}(2i))/\sqrt{2} \\ c^j(2^j+i) &= (c^{j+1}(2i-1) - c^{j+1}(2i))/\sqrt{2}. \end{aligned}$$

Note that things are now more symmetric, at the expense of a division by $\sqrt{2}$. However, for long vectors, it turns out that these algorithms are numerically more stable.

Remark: Some authors (for example, Stollnitz, Deroose and Salesin [103]) rescale c by $1/\sqrt{2^n}$ and u by $\sqrt{2^n}$. This is because the norm of the basis functions ψ_k^j is not equal to 1 (under the inner product $\langle f, g \rangle = \int_0^1 f(t)g(t)dt$). The normalized basis functions are the functions $\sqrt{2^j}\psi_k^j$.

Let us now explain the 2D version of the Haar transform. We describe the version using the matrix W_n , the method using H_n being identical (except that $H_n^{-1} = H_n^\top$, but this does not hold for W_n^{-1}). Given a $2^m \times 2^n$ matrix A , we can first convert the *rows* of A to their Haar coefficients using the Haar transform W_n^{-1} , obtaining a matrix B , and then convert the *columns* of B to their Haar coefficients, using the matrix W_m^{-1} . Because columns and rows are exchanged in the first step,

$$B = A(W_n^{-1})^\top,$$

and in the second step $C = W_m^{-1}B$, thus, we have

$$C = W_m^{-1}A(W_n^{-1})^\top = D_m W_m^\top A W_n D_n.$$

In the other direction, given a matrix C of Haar coefficients, we reconstruct the matrix A (the image) by first applying W_m to the columns of C , obtaining B , and then W_n^\top to the rows of B . Therefore

$$A = W_m C W_n^\top.$$

Of course, we don't actually have to invert W_m and W_n and perform matrix multiplications. We just have to use our algorithms using averaging and differencing. Here is an example.

If the data matrix (the image) is the 8×8 matrix

$$A = \begin{pmatrix} 64 & 2 & 3 & 61 & 60 & 6 & 7 & 57 \\ 9 & 55 & 54 & 12 & 13 & 51 & 50 & 16 \\ 17 & 47 & 46 & 20 & 21 & 43 & 42 & 24 \\ 40 & 26 & 27 & 37 & 36 & 30 & 31 & 33 \\ 32 & 34 & 35 & 29 & 28 & 38 & 39 & 25 \\ 41 & 23 & 22 & 44 & 45 & 19 & 18 & 48 \\ 49 & 15 & 14 & 52 & 53 & 11 & 10 & 56 \\ 8 & 58 & 59 & 5 & 4 & 62 & 63 & 1 \end{pmatrix},$$

then applying our algorithms, we find that

$$C = \begin{pmatrix} 32.5 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 4 & -4 & 4 & -4 \\ 0 & 0 & 0 & 0 & 4 & -4 & 4 & -4 \\ 0 & 0 & 0.5 & 0.5 & 27 & -25 & 23 & -21 \\ 0 & 0 & -0.5 & -0.5 & -11 & 9 & -7 & 5 \\ 0 & 0 & 0.5 & 0.5 & -5 & 7 & -9 & 11 \\ 0 & 0 & -0.5 & -0.5 & 21 & -23 & 25 & -27 \end{pmatrix}.$$

As we can see, C has a more zero entries than A ; it is a compressed version of A . We can further compress C by setting to 0 all entries of absolute value at most 0.5. Then, we get

$$C_2 = \begin{pmatrix} 32.5 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 4 & -4 & 4 & -4 \\ 0 & 0 & 0 & 0 & 4 & -4 & 4 & -4 \\ 0 & 0 & 0 & 0 & 27 & -25 & 23 & -21 \\ 0 & 0 & 0 & 0 & -11 & 9 & -7 & 5 \\ 0 & 0 & 0 & 0 & -5 & 7 & -9 & 11 \\ 0 & 0 & 0 & 0 & 21 & -23 & 25 & -27 \end{pmatrix}.$$

We find that the reconstructed image is

$$A_2 = \begin{pmatrix} 63.5 & 1.5 & 3.5 & 61.5 & 59.5 & 5.5 & 7.5 & 57.5 \\ 9.5 & 55.5 & 53.5 & 11.5 & 13.5 & 51.5 & 49.5 & 15.5 \\ 17.5 & 47.5 & 45.5 & 19.5 & 21.5 & 43.5 & 41.5 & 23.5 \\ 39.5 & 25.5 & 27.5 & 37.5 & 35.5 & 29.5 & 31.5 & 33.5 \\ 31.5 & 33.5 & 35.5 & 29.5 & 27.5 & 37.5 & 39.5 & 25.5 \\ 41.5 & 23.5 & 21.5 & 43.5 & 45.5 & 19.5 & 17.5 & 47.5 \\ 49.5 & 15.5 & 13.5 & 51.5 & 53.5 & 11.5 & 9.5 & 55.5 \\ 7.5 & 57.5 & 59.5 & 5.5 & 3.5 & 61.5 & 63.5 & 1.5 \end{pmatrix},$$

which is pretty close to the original image matrix A .

It turns out that **Matlab** has a wonderful command, `image(X)`, which displays the matrix X has an image in which each entry is shown as a little square whose gray level is proportional to the numerical value of that entry (lighter if the value is higher, darker if the value is closer to zero; negative values are treated as zero). The images corresponding to A and C are shown in Figure 3.10. The compressed images corresponding to A_2 and C_2 are shown in

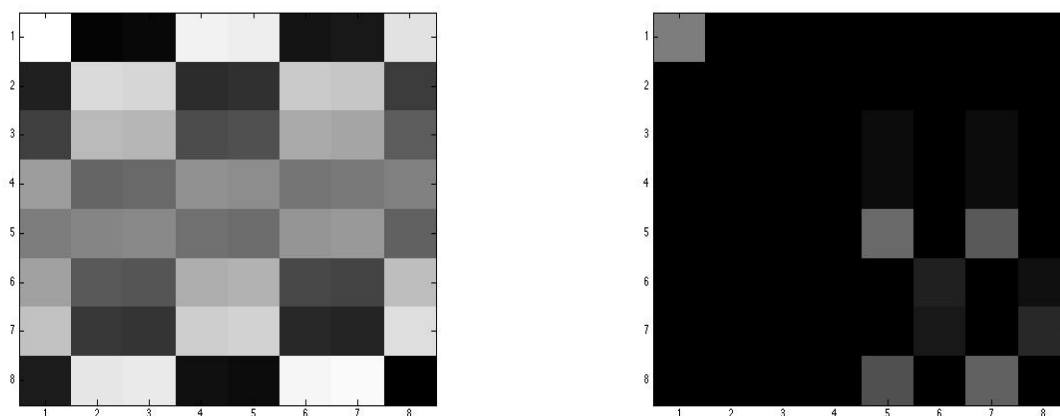


Figure 3.10: An image and its Haar transform

Figure 3.11. The compressed versions appear to be indistinguishable from the originals!

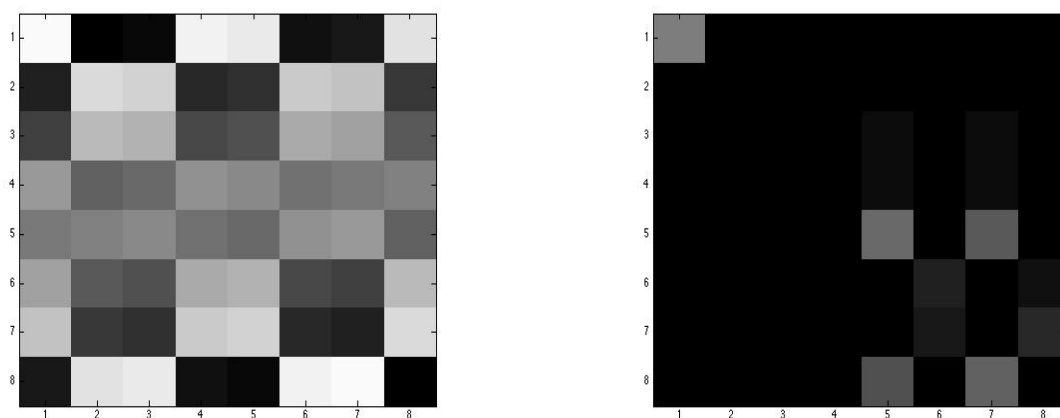


Figure 3.11: Compressed image and its Haar transform

If we use the normalized matrices H_m and H_n , then the equations relating the image

matrix A and its normalized Haar transform C are

$$\begin{aligned} C &= H_m^\top A H_n \\ A &= H_m C H_n^\top. \end{aligned}$$

The Haar transform can also be used to send large images progressively over the internet. Indeed, we can start sending the Haar coefficients of the matrix C starting from the coarsest coefficients (the first column from top down, then the second column, etc.) and at the receiving end we can start reconstructing the image as soon as we have received enough data.

Observe that instead of performing all rounds of averaging and differencing on each row and each column, we can perform partial encoding (and decoding). For example, we can perform a single round of averaging and differencing for each row and each column. The result is an image consisting of four subimages, where the top left quarter is a coarser version of the original, and the rest (consisting of three pieces) contain the finest detail coefficients. We can also perform two rounds of averaging and differencing, or three rounds, *etc.* This process is illustrated on the image shown in Figure 3.12. The result of performing one round, two rounds, three rounds, and nine rounds of averaging is shown in Figure 3.13. Since our images have size 512×512 , nine rounds of averaging yields the Haar transform, displayed as the image on the bottom right. The original image has completely disappeared! We leave it as a fun exercise to modify the algorithms involving averaging and differencing to perform k rounds of averaging/differencing. The reconstruction algorithm is a little tricky.

A nice and easily accessible account of wavelets and their uses in image processing and computer graphics can be found in Stollnitz, Deroose and Salesin [103]. A very detailed account is given in Strang and and Nguyen [106], but this book assumes a fair amount of background in signal processing.

We can find easily a basis of $2^n \times 2^n = 2^{2n}$ vectors w_{ij} ($2^n \times 2^n$ matrices) for the linear map that reconstructs an image from its Haar coefficients, in the sense that for any matrix C of Haar coefficients, the image matrix A is given by

$$A = \sum_{i=1}^{2^n} \sum_{j=1}^{2^n} c_{ij} w_{ij}.$$

Indeed, the matrix w_{ij} is given by the so-called outer product

$$w_{ij} = w_i (w_j)^\top.$$

Similarly, there is a basis of $2^n \times 2^n = 2^{2n}$ vectors h_{ij} ($2^n \times 2^n$ matrices) for the 2D Haar transform, in the sense that for any matrix A , its matrix C of Haar coefficients is given by

$$C = \sum_{i=1}^{2^n} \sum_{j=1}^{2^n} a_{ij} h_{ij}.$$

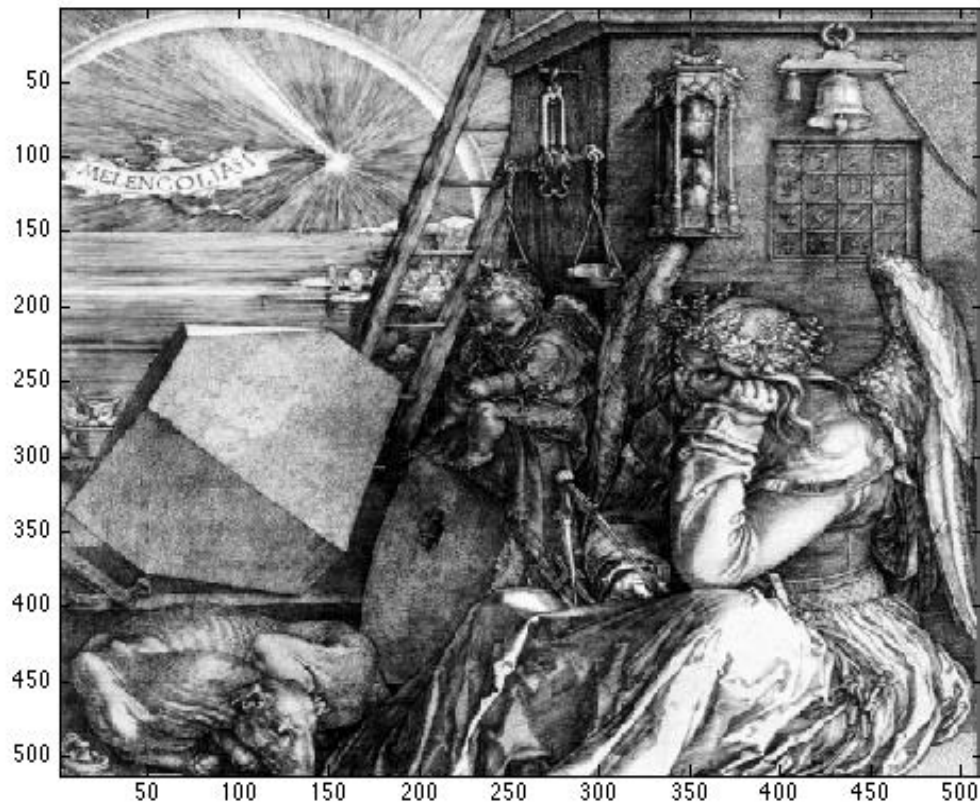


Figure 3.12: Original drawing by Durer

If the columns of W^{-1} are w'_1, \dots, w'_{2n} , then

$$h_{ij} = w'_i(w'_j)^\top.$$

We leave it as exercise to compute the bases (w_{ij}) and (h_{ij}) for $n = 2$, and to display the corresponding images using the command `imagesc`.

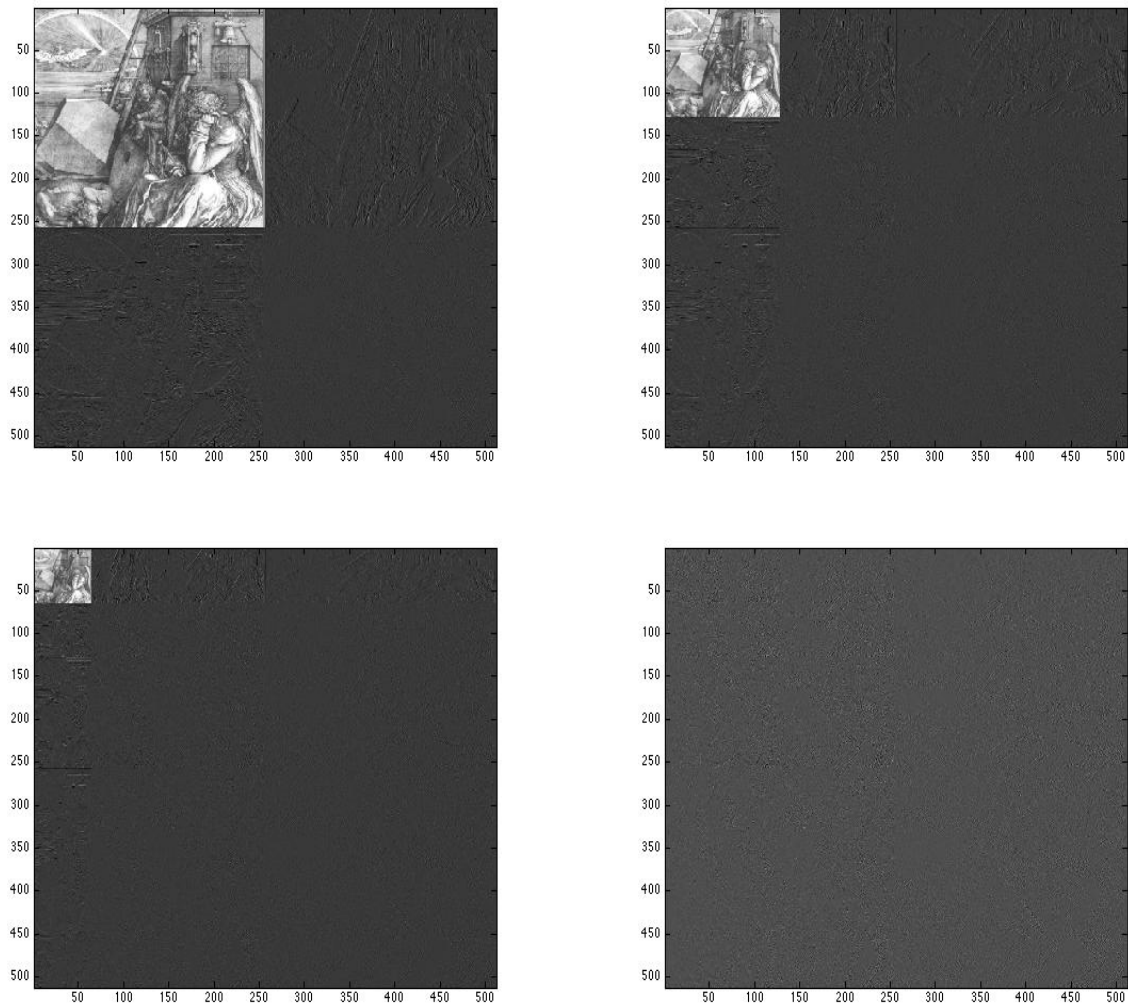


Figure 3.13: Haar tranforms after one, two, three, and nine rounds of averaging

3.3 The Effect of a Change of Bases on Matrices

The effect of a change of bases on the representation of a linear map is described in the following proposition.

Proposition 3.4. *Let E and F be vector spaces, let $\mathcal{U} = (u_1, \dots, u_n)$ and $\mathcal{U}' = (u'_1, \dots, u'_n)$ be two bases of E , and let $\mathcal{V} = (v_1, \dots, v_m)$ and $\mathcal{V}' = (v'_1, \dots, v'_m)$ be two bases of F . Let $P = P_{\mathcal{U}', \mathcal{U}}$ be the change of basis matrix from \mathcal{U} to \mathcal{U}' , and let $Q = P_{\mathcal{V}', \mathcal{V}}$ be the change of basis matrix from \mathcal{V} to \mathcal{V}' . For any linear map $f: E \rightarrow F$, let $M(f) = M_{\mathcal{U}, \mathcal{V}}(f)$ be the matrix associated to f w.r.t. the bases \mathcal{U} and \mathcal{V} , and let $M'(f) = M_{\mathcal{U}', \mathcal{V}'}(f)$ be the matrix associated to f w.r.t. the bases \mathcal{U}' and \mathcal{V}' . We have*

$$M'(f) = Q^{-1}M(f)P,$$

or more explicitly

$$M_{\mathcal{U}', \mathcal{V}'}(f) = P_{\mathcal{V}', \mathcal{V}}^{-1} M_{\mathcal{U}, \mathcal{V}}(f) P_{\mathcal{U}', \mathcal{U}} = P_{\mathcal{V}, \mathcal{V}'} M_{\mathcal{U}, \mathcal{V}}(f) P_{\mathcal{U}', \mathcal{U}}.$$

Proof. Since $f: E \rightarrow F$ can be written as $f = \text{id}_F \circ f \circ \text{id}_E$, since P is the matrix of id_E w.r.t. the bases (u'_1, \dots, u'_n) and (u_1, \dots, u_n) , and Q^{-1} is the matrix of id_F w.r.t. the bases (v_1, \dots, v_m) and (v'_1, \dots, v'_m) , by Proposition 3.2, we have $M'(f) = Q^{-1}M(f)P$. \square

As a corollary, we get the following result.

Corollary 3.5. *Let E be a vector space, and let $\mathcal{U} = (u_1, \dots, u_n)$ and $\mathcal{U}' = (u'_1, \dots, u'_n)$ be two bases of E . Let $P = P_{\mathcal{U}', \mathcal{U}}$ be the change of basis matrix from \mathcal{U} to \mathcal{U}' . For any linear map $f: E \rightarrow E$, let $M(f) = M_{\mathcal{U}}(f)$ be the matrix associated to f w.r.t. the basis \mathcal{U} , and let $M'(f) = M_{\mathcal{U}'}(f)$ be the matrix associated to f w.r.t. the basis \mathcal{U}' . We have*

$$M'(f) = P^{-1}M(f)P,$$

or more explicitly,

$$M_{\mathcal{U}'}(f) = P_{\mathcal{U}', \mathcal{U}}^{-1} M_{\mathcal{U}}(f) P_{\mathcal{U}', \mathcal{U}} = P_{\mathcal{U}, \mathcal{U}'} M_{\mathcal{U}}(f) P_{\mathcal{U}', \mathcal{U}}.$$

Example 3.4. Let $E = \mathbb{R}^2$, $\mathcal{U} = (e_1, e_2)$ where $e_1 = (1, 0)$ and $e_2 = (0, 1)$ are the canonical basis vectors, let $\mathcal{V} = (v_1, v_2) = (e_1, e_1 - e_2)$, and let

$$A = \begin{pmatrix} 2 & 1 \\ 0 & 1 \end{pmatrix}.$$

The change of basis matrix $P = P_{\mathcal{V}, \mathcal{U}}$ from \mathcal{U} to \mathcal{V} is

$$P = \begin{pmatrix} 1 & 1 \\ 0 & -1 \end{pmatrix},$$

and we check that

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

Therefore, in the basis \mathcal{V} , the matrix representing the linear map f defined by A is

$$A' = P^{-1}AP = PAP = \begin{pmatrix} 1 & 1 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} 2 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 0 & -1 \end{pmatrix} = \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix} = D,$$

a diagonal matrix. Therefore, in the basis \mathcal{V} , it is clear what the action of f is: it is a stretch by a factor of 2 in the v_1 direction and it is the identity in the v_2 direction. Observe that v_1 and v_2 are not orthogonal.

What happened is that we *diagonalized* the matrix A . The diagonal entries 2 and 1 are the *eigenvalues* of A (and f) and v_1 and v_2 are corresponding *eigenvectors*. We will come back to eigenvalues and eigenvectors later on.

The above example showed that the same linear map can be represented by different matrices. This suggests making the following definition:

Definition 3.5. Two $n \times n$ matrices A and B are said to be *similar* iff there is some invertible matrix P such that

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

It is easily checked that similarity is an equivalence relation. From our previous considerations, two $n \times n$ matrices A and B are similar iff they represent the same linear map with respect to two different bases. The following surprising fact can be shown: Every square matrix A is similar to its transpose A^\top . The proof requires advanced concepts than we will not discuss in these notes (the Jordan form, or similarity invariants).

If $\mathcal{U} = (u_1, \dots, u_n)$ and $\mathcal{V} = (v_1, \dots, v_n)$ are two bases of E , the change of basis matrix

$$P = P_{\mathcal{V}, \mathcal{U}} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix}$$

from (u_1, \dots, u_n) to (v_1, \dots, v_n) is the matrix whose j th column consists of the coordinates of v_j over the basis (u_1, \dots, u_n) , which means that

$$v_j = \sum_{i=1}^n a_{ij} u_i.$$

It is natural to extend the matrix notation and to express the vector $\begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix}$ in E^n as the product of a matrix times the vector $\begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix}$ in E^n , namely as

$$\begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix} = \begin{pmatrix} a_{11} & a_{21} & \cdots & a_{n1} \\ a_{12} & a_{22} & \cdots & a_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1n} & a_{2n} & \cdots & a_{nn} \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix},$$

but notice that the matrix involved is not P , but its transpose P^\top .

This observation has the following consequence: if $\mathcal{U} = (u_1, \dots, u_n)$ and $\mathcal{V} = (v_1, \dots, v_n)$ are two bases of E and if

$$\begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix} = A \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix},$$

that is,

$$v_i = \sum_{j=1}^n a_{ij} u_j,$$

for any vector $w \in E$, if

$$w = \sum_{i=1}^n x_i u_i = \sum_{k=1}^n y_k v_k,$$

then

$$\begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = A^\top \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix},$$

and so

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = (A^\top)^{-1} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}.$$

It is easy to see that $(A^\top)^{-1} = (A^{-1})^\top$. Also, if $\mathcal{U} = (u_1, \dots, u_n)$, $\mathcal{V} = (v_1, \dots, v_n)$, and $\mathcal{W} = (w_1, \dots, w_n)$ are three bases of E , and if the change of basis matrix from \mathcal{U} to \mathcal{V} is $P = P_{\mathcal{V}, \mathcal{U}}$ and the change of basis matrix from \mathcal{V} to \mathcal{W} is $Q = P_{\mathcal{W}, \mathcal{V}}$, then

$$\begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix} = P^\top \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix}, \quad \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} = Q^\top \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix},$$

so

$$\begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} = Q^\top P^\top \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix} = (PQ)^\top \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix},$$

which means that the change of basis matrix $P_{\mathcal{W},\mathcal{U}}$ from \mathcal{U} to \mathcal{W} is PQ . This proves that

$$P_{\mathcal{W},\mathcal{U}} = P_{\mathcal{V},\mathcal{U}} P_{\mathcal{W},\mathcal{V}}.$$

3.4 Summary

The main concepts and results of this chapter are listed below:

- The representation of linear maps by *matrices*.
- The vector space of linear maps $\text{Hom}_K(E, F)$.
- The vector space $M_{m,n}(K)$ of $m \times n$ matrices over the field K ; The ring $M_n(K)$ of $n \times n$ matrices over the field K .
- *Column vectors, row vectors*.
- *Matrix operations*: addition, scalar multiplication, multiplication.
- The *matrix representation mapping* $M: \text{Hom}(E, F) \rightarrow M_{n,p}$ and the representation isomorphism (Proposition 3.2).
- Haar basis vectors and a glimpse at *Haar wavelets*.
- *Change of basis matrix* and Proposition 3.4.

Chapter 4

Direct Sums, The Dual Space, Duality

4.1 Sums, Direct Sums, Direct Products

Before considering linear forms and hyperplanes, we define the notion of direct sum and prove some simple propositions. There is a subtle point, which is that if we attempt to define the direct sum $E \coprod F$ of two vector spaces using the cartesian product $E \times F$, we don't quite get the right notion because elements of $E \times F$ are ordered pairs, but we want $E \coprod F = F \coprod E$. Thus, we want to think of the elements of $E \coprod F$ as unordered pairs of elements. It is possible to do so by considering the direct sum of a *family* $(E_i)_{i \in \{1,2\}}$, and more generally of a family $(E_i)_{i \in I}$. For simplicity, we begin by considering the case where $I = \{1, 2\}$.

Definition 4.1. Given a family $(E_i)_{i \in \{1,2\}}$ of two vector spaces, we define the (*external*) *direct sum* $E_1 \coprod E_2$ (or *coproduct*) of the family $(E_i)_{i \in \{1,2\}}$ as the set

$$E_1 \coprod E_2 = \{\{\langle 1, u \rangle, \langle 2, v \rangle\} \mid u \in E_1, v \in E_2\},$$

with addition

$$\{\langle 1, u_1 \rangle, \langle 2, v_1 \rangle\} + \{\langle 1, u_2 \rangle, \langle 2, v_2 \rangle\} = \{\langle 1, u_1 + u_2 \rangle, \langle 2, v_1 + v_2 \rangle\},$$

and scalar multiplication

$$\lambda\{\langle 1, u \rangle, \langle 2, v \rangle\} = \{\langle 1, \lambda u \rangle, \langle 2, \lambda v \rangle\}.$$

We define the *injections* $in_1: E_1 \rightarrow E_1 \coprod E_2$ and $in_2: E_2 \rightarrow E_1 \coprod E_2$ as the linear maps defined such that,

$$in_1(u) = \{\langle 1, u \rangle, \langle 2, 0 \rangle\},$$

and

$$in_2(v) = \{\langle 1, 0 \rangle, \langle 2, v \rangle\}.$$

Note that

$$E_2 \coprod E_1 = \{\{\langle 2, v \rangle, \langle 1, u \rangle\} \mid v \in E_2, u \in E_1\} = E_1 \coprod E_2.$$

Thus, every member $\{\langle 1, u \rangle, \langle 2, v \rangle\}$ of $E_1 \coprod E_2$ can be viewed as an *unordered pair* consisting of the two vectors u and v , tagged with the index 1 and 2, respectively.

Remark: In fact, $E_1 \coprod E_2$ is just the product $\prod_{i \in \{1,2\}} E_i$ of the family $(E_i)_{i \in \{1,2\}}$.



This is not to be confused with the cartesian product $E_1 \times E_2$. The vector space $E_1 \times E_2$ is the set of all ordered pairs $\langle u, v \rangle$, where $u \in E_1$, and $v \in E_2$, with addition and multiplication by a scalar defined such that

$$\begin{aligned} \langle u_1, v_1 \rangle + \langle u_2, v_2 \rangle &= \langle u_1 + u_2, v_1 + v_2 \rangle, \\ \lambda \langle u, v \rangle &= \langle \lambda u, \lambda v \rangle. \end{aligned}$$

There is a bijection between $\prod_{i \in \{1,2\}} E_i$ and $E_1 \times E_2$, but as we just saw, elements of $\prod_{i \in \{1,2\}} E_i$ are certain sets. The product $E_1 \times \cdots \times E_n$ of any number of vector spaces can also be defined. We will do this shortly.

The following property holds.

Proposition 4.1. *Given any two vector spaces, E_1 and E_2 , the set $E_1 \coprod E_2$ is a vector space. For every pair of linear maps, $f: E_1 \rightarrow G$ and $g: E_2 \rightarrow G$, there is a unique linear map, $f + g: E_1 \coprod E_2 \rightarrow G$, such that $(f + g) \circ in_1 = f$ and $(f + g) \circ in_2 = g$, as in the following diagram:*

$$\begin{array}{ccc} E_1 & & \\ \downarrow in_1 & \searrow f & \\ E_1 \coprod E_2 & \xrightarrow{f+g} & G \\ \uparrow in_2 & \nearrow g & \\ E_2 & & \end{array}$$

Proof. Define

$$(f + g)(\{\langle 1, u \rangle, \langle 2, v \rangle\}) = f(u) + g(v),$$

for every $u \in E_1$ and $v \in E_2$. It is immediately verified that $f + g$ is the unique linear map with the required properties. \square

We already noted that $E_1 \coprod E_2$ is in bijection with $E_1 \times E_2$. If we define the *projections* $\pi_1: E_1 \coprod E_2 \rightarrow E_1$ and $\pi_2: E_1 \coprod E_2 \rightarrow E_2$, such that

$$\pi_1(\{\langle 1, u \rangle, \langle 2, v \rangle\}) = u,$$

and

$$\pi_2(\{\langle 1, u \rangle, \langle 2, v \rangle\}) = v,$$

we have the following proposition.

Proposition 4.2. *Given any two vector spaces, E_1 and E_2 , for every pair of linear maps, $f: D \rightarrow E_1$ and $g: D \rightarrow E_2$, there is a unique linear map, $f \times g: D \rightarrow E_1 \amalg E_2$, such that $\pi_1 \circ (f \times g) = f$ and $\pi_2 \circ (f \times g) = g$, as in the following diagram:*

$$\begin{array}{ccccc}
 & & E_1 & & \\
 & \nearrow f & \uparrow \pi_1 & & \\
 D & \xrightarrow{f \times g} & E_1 \amalg E_2 & & \\
 & \searrow g & \downarrow \pi_2 & & \\
 & & E_2 & &
 \end{array}$$

Proof. Define

$$(f \times g)(w) = \{\langle 1, f(w) \rangle, \langle 2, g(w) \rangle\},$$

for every $w \in D$. It is immediately verified that $f \times g$ is the unique linear map with the required properties. \square

Remark: It is a peculiarity of linear algebra that direct sums and products of finite families are isomorphic. However, this is no longer true for products and sums of infinite families.

When U, V are subspaces of a vector space E , letting $i_1: U \rightarrow E$ and $i_2: V \rightarrow E$ be the inclusion maps, if $U \amalg V$ is isomorphic to E under the map $i_1 + i_2$ given by Proposition 4.1, we say that E is a *direct sum* of U and V , and we write $E = U \amalg V$ (with a slight abuse of notation, since E and $U \amalg V$ are only isomorphic). It is also convenient to define the sum $U_1 + \cdots + U_p$ and the internal direct sum $U_1 \oplus \cdots \oplus U_p$ of any number of subspaces of E .

Definition 4.2. Given $p \geq 2$ vector spaces E_1, \dots, E_p , the product $F = E_1 \times \cdots \times E_p$ can be made into a vector space by defining addition and scalar multiplication as follows:

$$\begin{aligned}
 (u_1, \dots, u_p) + (v_1, \dots, v_p) &= (u_1 + v_1, \dots, u_p + v_p) \\
 \lambda(u_1, \dots, u_p) &= (\lambda u_1, \dots, \lambda u_p),
 \end{aligned}$$

for all $u_i, v_i \in E_i$ and all $\lambda \in K$. With the above addition and multiplication, the vector space $F = E_1 \times \cdots \times E_p$ is called the *direct product* of the vector spaces E_1, \dots, E_p .

As a special case, when $E_1 = \cdots = E_p = K$, we find again the vector space $F = K^p$. The *projection maps* $pr_i: E_1 \times \cdots \times E_p \rightarrow E_i$ given by

$$pr_i(u_1, \dots, u_p) = u_i$$

are clearly linear. Similarly, the maps $in_i: E_i \rightarrow E_1 \times \cdots \times E_p$ given by

$$in_i(u_i) = (0, \dots, 0, u_i, 0, \dots, 0)$$

are injective and linear. If $\dim(E_i) = n_i$ and if $(e_1^i, \dots, e_{n_i}^i)$ is a basis of E_i for $i = 1, \dots, p$, then it is easy to see that the $n_1 + \dots + n_p$ vectors

$$\begin{array}{ccc} (e_1^1, 0, \dots, 0), & \dots, & (e_{n_1}^1, 0, \dots, 0), \\ \vdots & & \vdots \\ (0, \dots, 0, e_1^i, 0, \dots, 0), & \dots, & (0, \dots, 0, e_{n_i}^i, 0, \dots, 0), \\ \vdots & & \vdots \\ (0, \dots, 0, e_1^p), & \dots, & (0, \dots, 0, e_{n_p}^p) \end{array}$$

form a basis of $E_1 \times \dots \times E_p$, and so

$$\dim(E_1 \times \dots \times E_p) = \dim(E_1) + \dots + \dim(E_p).$$

Let us now consider a vector space E and p subspaces U_1, \dots, U_p of E . We have a map

$$a: U_1 \times \dots \times U_p \rightarrow E$$

given by

$$a(u_1, \dots, u_p) = u_1 + \dots + u_p,$$

with $u_i \in U_i$ for $i = 1, \dots, p$. It is clear that this map is linear, and so its image is a subspace of E denoted by

$$U_1 + \dots + U_p$$

and called the *sum* of the subspaces U_1, \dots, U_p . It is immediately verified that $U_1 + \dots + U_p$ is the smallest subspace of E containing U_1, \dots, U_p . This also implies that $U_1 + \dots + U_p$ does not depend on the order of the factors U_i ; in particular,

$$U_1 + U_2 = U_2 + U_1.$$

If the map a is injective, then $\text{Ker } a = 0$, which means that if $u_i \in U_i$ for $i = 1, \dots, p$ and if

$$u_1 + \dots + u_p = 0$$

then $u_1 = \dots = u_p = 0$. In this case, every $u \in U_1 + \dots + U_p$ has a *unique* expression as a sum

$$u = u_1 + \dots + u_p,$$

with $u_i \in U_i$, for $i = 1, \dots, p$. It is also clear that for any p nonzero vectors $u_i \in U_i$, u_1, \dots, u_p are linearly independent.

Definition 4.3. For any vector space E and any $p \geq 2$ subspaces U_1, \dots, U_p of E , if the map a defined above is injective, then the sum $U_1 + \dots + U_p$ is called a *direct sum* and it is denoted by

$$U_1 \oplus \dots \oplus U_p.$$

The space E is the *direct sum* of the subspaces U_i if

$$E = U_1 \oplus \dots \oplus U_p.$$

As in the case of a sum, $U_1 \oplus U_2 = U_2 \oplus U_1$. Observe that when the map a is injective, then it is a linear isomorphism between $U_1 \times \cdots \times U_p$ and $U_1 \oplus \cdots \oplus U_p$. The difference is that $U_1 \times \cdots \times U_p$ is defined even if the spaces U_i are not assumed to be subspaces of some common space.

Now, if $p = 2$, it is easy to determine the kernel of the map $a: U_1 \times U_2 \rightarrow E$. We have

$$a(u_1, u_2) = u_1 + u_2 = 0 \quad \text{iff} \quad u_1 = -u_2, \quad u_1 \in U_1, u_2 \in U_2,$$

which implies that

$$\text{Ker } a = \{(u, -u) \mid u \in U_1 \cap U_2\}.$$

Now, $U_1 \cap U_2$ is a subspace of E and the linear map $u \mapsto (u, -u)$ is clearly an isomorphism, so $\text{Ker } a$ is isomorphic to $U_1 \cap U_2$. As a result, we get the following result:

Proposition 4.3. *Given any vector space E and any two subspaces U_1 and U_2 , the sum $U_1 + U_2$ is a direct sum iff $U_1 \cap U_2 = (0)$.*

An interesting illustration of the notion of direct sum is the decomposition of a square matrix into its symmetric part and its skew-symmetric part. Recall that an $n \times n$ matrix $A \in M_n$ is *symmetric* if $A^\top = A$, *skew-symmetric* if $A^\top = -A$. It is clear that

$$\mathbf{S}(n) = \{A \in M_n \mid A^\top = A\} \quad \text{and} \quad \mathbf{Skew}(n) = \{A \in M_n \mid A^\top = -A\}$$

are subspaces of M_n , and that $\mathbf{S}(n) \cap \mathbf{Skew}(n) = (0)$. Observe that for any matrix $A \in M_n$, the matrix $H(A) = (A + A^\top)/2$ is symmetric and the matrix $S(A) = (A - A^\top)/2$ is skew-symmetric. Since

$$A = H(A) + S(A) = \frac{A + A^\top}{2} + \frac{A - A^\top}{2},$$

we see that $M_n = \mathbf{S}(n) + \mathbf{Skew}(n)$, and since $\mathbf{S}(n) \cap \mathbf{Skew}(n) = (0)$, we have the direct sum

$$M_n = \mathbf{S}(n) \oplus \mathbf{Skew}(n).$$

Remark: The vector space $\mathbf{Skew}(n)$ of skew-symmetric matrices is also denoted by $\mathfrak{so}(n)$. It is the *Lie algebra* of the group $\mathbf{SO}(n)$.

Proposition 4.3 can be generalized to any $p \geq 2$ subspaces at the expense of notation. The proof of the following proposition is left as an exercise.

Proposition 4.4. *Given any vector space E and any $p \geq 2$ subspaces U_1, \dots, U_p , the following properties are equivalent:*

- (1) *The sum $U_1 + \cdots + U_p$ is a direct sum.*
- (2) *We have*

$$U_i \cap \left(\sum_{j=1, j \neq i}^p U_j \right) = (0), \quad i = 1, \dots, p.$$

(3) We have

$$U_i \cap \left(\sum_{j=1}^{i-1} U_j \right) = (0), \quad i = 2, \dots, p.$$

Because of the isomorphism

$$U_1 \times \cdots \times U_p \approx U_1 \oplus \cdots \oplus U_p,$$

we have

Proposition 4.5. *If E is any vector space, for any (finite-dimensional) subspaces U_1, \dots, U_p of E , we have*

$$\dim(U_1 \oplus \cdots \oplus U_p) = \dim(U_1) + \cdots + \dim(U_p).$$

If E is a direct sum

$$E = U_1 \oplus \cdots \oplus U_p,$$

since every $u \in E$ can be written in a unique way as

$$u = u_1 + \cdots + u_p$$

for some $u_i \in U_i$ for $i = 1, \dots, p$, we can define the maps $\pi_i: E \rightarrow U_i$, called *projections*, by

$$\pi_i(u) = \pi_i(u_1 + \cdots + u_p) = u_i.$$

It is easy to check that these maps are linear and satisfy the following properties:

$$\pi_j \circ \pi_i = \begin{cases} \pi_i & \text{if } i = j \\ 0 & \text{if } i \neq j, \end{cases}$$

$$\pi_1 + \cdots + \pi_p = \text{id}_E.$$

For example, in the case of the direct sum

$$M_n = \mathbf{S}(n) \oplus \mathbf{Skew}(n),$$

the projection onto $\mathbf{S}(n)$ is given by

$$\pi_1(A) = H(A) = \frac{A + A^\top}{2},$$

and the projection onto $\mathbf{Skew}(n)$ is given by

$$\pi_2(A) = S(A) = \frac{A - A^\top}{2}.$$

Clearly, $H(A) + S(A) = A$, $H(H(A)) = H(A)$, $S(S(A)) = S(A)$, and $H(S(A)) = S(H(A)) = 0$.

A function f such that $f \circ f = f$ is said to be *idempotent*. Thus, the projections π_i are idempotent. Conversely, the following proposition can be shown:

Proposition 4.6. *Let E be a vector space. For any $p \geq 2$ linear maps $f_i: E \rightarrow E$, if*

$$f_j \circ f_i = \begin{cases} f_i & \text{if } i = j \\ 0 & \text{if } i \neq j, \end{cases}$$

$$f_1 + \cdots + f_p = \text{id}_E,$$

then if we let $U_i = f_i(E)$, we have a direct sum

$$E = U_1 \oplus \cdots \oplus U_p.$$

We also have the following proposition characterizing idempotent linear maps whose proof is also left as an exercise.

Proposition 4.7. *For every vector space E , if $f: E \rightarrow E$ is an idempotent linear map, i.e., $f \circ f = f$, then we have a direct sum*

$$E = \text{Ker } f \oplus \text{Im } f,$$

so that f is the projection onto its image $\text{Im } f$.

We now give the definition of a direct sum for any arbitrary nonempty index set I . First, let us recall the notion of the product of a family $(E_i)_{i \in I}$. Given a family of sets $(E_i)_{i \in I}$, its product $\prod_{i \in I} E_i$, is the set of all functions $f: I \rightarrow \bigcup_{i \in I} E_i$, such that, $f(i) \in E_i$, for every $i \in I$. It is one of the many versions of the axiom of choice, that, if $E_i \neq \emptyset$ for every $i \in I$, then $\prod_{i \in I} E_i \neq \emptyset$. A member $f \in \prod_{i \in I} E_i$, is often denoted as $(f_i)_{i \in I}$. For every $i \in I$, we have the *projection* $\pi_i: \prod_{i \in I} E_i \rightarrow E_i$, defined such that, $\pi_i((f_i)_{i \in I}) = f_i$. We now define direct sums.

Definition 4.4. Let I be any nonempty set, and let $(E_i)_{i \in I}$ be a family of vector spaces. The *(external) direct sum* $\coprod_{i \in I} E_i$ (or *coproduct*) of the family $(E_i)_{i \in I}$ is defined as follows:

$\coprod_{i \in I} E_i$ consists of all $f \in \prod_{i \in I} E_i$, which have finite support, and addition and multiplication by a scalar are defined as follows:

$$(f_i)_{i \in I} + (g_i)_{i \in I} = (f_i + g_i)_{i \in I},$$

$$\lambda(f_i)_{i \in I} = (\lambda f_i)_{i \in I}.$$

We also have *injection maps* $\text{in}_i: E_i \rightarrow \coprod_{i \in I} E_i$, defined such that, $\text{in}_i(x) = (f_i)_{i \in I}$, where $f_i = x$, and $f_j = 0$, for all $j \in (I - \{i\})$.

The following proposition is an obvious generalization of Proposition 4.1.

Proposition 4.8. *Let I be any nonempty set, let $(E_i)_{i \in I}$ be a family of vector spaces, and let G be any vector space. The direct sum $\coprod_{i \in I} E_i$ is a vector space, and for every family $(h_i)_{i \in I}$ of linear maps $h_i: E_i \rightarrow G$, there is a unique linear map*

$$\left(\sum_{i \in I} h_i \right): \coprod_{i \in I} E_i \rightarrow G,$$

such that, $(\sum_{i \in I} h_i) \circ \text{in}_i = h_i$, for every $i \in I$.

Remark: When $E_i = E$, for all $i \in I$, we denote $\coprod_{i \in I} E_i$ by $E^{(I)}$. In particular, when $E_i = K$, for all $i \in I$, we find the vector space $K^{(I)}$ of Definition 2.13.

We also have the following basic proposition about injective or surjective linear maps.

Proposition 4.9. *Let E and F be vector spaces, and let $f: E \rightarrow F$ be a linear map. If $f: E \rightarrow F$ is injective, then there is a surjective linear map $r: F \rightarrow E$ called a retraction, such that $r \circ f = \text{id}_E$. If $f: E \rightarrow F$ is surjective, then there is an injective linear map $s: F \rightarrow E$ called a section, such that $f \circ s = \text{id}_F$.*

Proof. Let $(u_i)_{i \in I}$ be a basis of E . Since $f: E \rightarrow F$ is an injective linear map, by Proposition 2.16, $(f(u_i))_{i \in I}$ is linearly independent in F . By Theorem 2.9, there is a basis $(v_j)_{j \in J}$ of F , where $I \subseteq J$, and where $v_i = f(u_i)$, for all $i \in I$. By Proposition 2.16, a linear map $r: F \rightarrow E$ can be defined such that $r(v_i) = u_i$, for all $i \in I$, and $r(v_j) = w$ for all $j \in (J - I)$, where w is any given vector in E , say $w = 0$. Since $r(f(u_i)) = u_i$ for all $i \in I$, by Proposition 2.16, we have $r \circ f = \text{id}_E$.

Now, assume that $f: E \rightarrow F$ is surjective. Let $(v_j)_{j \in J}$ be a basis of F . Since $f: E \rightarrow F$ is surjective, for every $v_j \in F$, there is some $u_j \in E$ such that $f(u_j) = v_j$. Since $(v_j)_{j \in J}$ is a basis of F , by Proposition 2.16, there is a unique linear map $s: F \rightarrow E$ such that $s(v_j) = u_j$. Also, since $f(s(v_j)) = v_j$, by Proposition 2.16 (again), we must have $f \circ s = \text{id}_F$. \square

The converse of Proposition 4.9 is obvious. We now have the following fundamental Proposition.

Proposition 4.10. *Let E , F and G , be three vector spaces, $f: E \rightarrow F$ an injective linear map, $g: F \rightarrow G$ a surjective linear map, and assume that $\text{Im } f = \text{Ker } g$. Then, the following properties hold. (a) For any section $s: G \rightarrow F$ of g , we have $F = \text{Ker } g \oplus \text{Im } s$, and the linear map $f + s: E \oplus G \rightarrow F$ is an isomorphism.¹*

(b) For any retraction $r: F \rightarrow E$ of f , we have $F = \text{Im } f \oplus \text{Ker } r$.²

$$E \begin{array}{c} \xrightarrow{f} \\ \xleftarrow{r} \end{array} F \begin{array}{c} \xrightarrow{g} \\ \xleftarrow{s} \end{array} G$$

¹The existence of a section $s: G \rightarrow F$ of g follows from Proposition 4.9.

²The existence of a retraction $r: F \rightarrow E$ of f follows from Proposition 4.9.

Proof. (a) Since $s: G \rightarrow F$ is a section of g , we have $g \circ s = \text{id}_G$, and for every $u \in F$,

$$g(u - s(g(u))) = g(u) - g(s(g(u))) = g(u) - g(u) = 0.$$

Thus, $u - s(g(u)) \in \text{Ker } g$, and we have $F = \text{Ker } g + \text{Im } s$. On the other hand, if $u \in \text{Ker } g \cap \text{Im } s$, then $u = s(v)$ for some $v \in G$ because $u \in \text{Im } s$, $g(u) = 0$ because $u \in \text{Ker } g$, and so,

$$g(u) = g(s(v)) = v = 0,$$

because $g \circ s = \text{id}_G$, which shows that $u = s(v) = 0$. Thus, $F = \text{Ker } g \oplus \text{Im } s$, and since by assumption, $\text{Im } f = \text{Ker } g$, we have $F = \text{Im } f \oplus \text{Im } s$. But then, since f and s are injective, $f + s: E \oplus G \rightarrow F$ is an isomorphism. The proof of (b) is very similar. \square

Note that we can choose a retraction $r: F \rightarrow E$ so that $\text{Ker } r = \text{Im } s$, since $F = \text{Ker } g \oplus \text{Im } s = \text{Im } f \oplus \text{Im } s$ and f is injective so we can set $r \equiv 0$ on $\text{Im } s$.

Given a sequence of linear maps $E \xrightarrow{f} F \xrightarrow{g} G$, when $\text{Im } f = \text{Ker } g$, we say that the sequence $E \xrightarrow{f} F \xrightarrow{g} G$ is *exact at F*. If in addition to being exact at F , f is injective and g is surjective, we say that we have a *short exact sequence*, and this is denoted as

$$0 \longrightarrow E \xrightarrow{f} F \xrightarrow{g} G \longrightarrow 0.$$

The property of a short exact sequence given by Proposition 4.10 is often described by saying that $0 \longrightarrow E \xrightarrow{f} F \xrightarrow{g} G \longrightarrow 0$ is a (short) *split exact sequence*.

As a corollary of Proposition 4.10, we have the following result.

Theorem 4.11. *Let E and F be vector spaces, and let $f: E \rightarrow F$ be a linear map. Then, E is isomorphic to $\text{Ker } f \oplus \text{Im } f$, and thus,*

$$\dim(E) = \dim(\text{Ker } f) + \dim(\text{Im } f) = \dim(\text{Ker } f) + \text{rk}(f).$$

Proof. Consider

$$\text{Ker } f \xrightarrow{i} E \xrightarrow{f'} \text{Im } f,$$

where $\text{Ker } f \xrightarrow{i} E$ is the inclusion map, and $E \xrightarrow{f'} \text{Im } f$ is the surjection associated with $E \xrightarrow{f} F$. Then, we apply Proposition 4.10 to any section $\text{Im } f \xrightarrow{s} E$ of f' to get an isomorphism between E and $\text{Ker } f \oplus \text{Im } f$, and Proposition 4.5, to get $\dim(E) = \dim(\text{Ker } f) + \dim(\text{Im } f)$. \square

Remark: The dimension $\dim(\text{Ker } f)$ of the kernel of a linear map f is often called the *nullity* of f .

We now derive some important results using Theorem 4.11.

Proposition 4.12. *Given a vector space E , if U and V are any two subspaces of E , then*

$$\dim(U) + \dim(V) = \dim(U + V) + \dim(U \cap V),$$

an equation known as Grassmann's relation.

Proof. Recall that $U + V$ is the image of the linear map

$$a: U \times V \rightarrow E$$

given by

$$a(u, v) = u + v,$$

and that we proved earlier that the kernel $\text{Ker } a$ of a is isomorphic to $U \cap V$. By Theorem 4.11,

$$\dim(U \times V) = \dim(\text{Ker } a) + \dim(\text{Im } a),$$

but $\dim(U \times V) = \dim(U) + \dim(V)$, $\dim(\text{Ker } a) = \dim(U \cap V)$, and $\text{Im } a = U + V$, so the Grassmann relation holds. \square

The Grassmann relation can be very useful to figure out whether two subspaces have a nontrivial intersection in spaces of dimension > 3 . For example, it is easy to see that in \mathbb{R}^5 , there are subspaces U and V with $\dim(U) = 3$ and $\dim(V) = 2$ such that $U \cap V = 0$; for example, let U be generated by the vectors $(1, 0, 0, 0, 0)$, $(0, 1, 0, 0, 0)$, $(0, 0, 1, 0, 0)$, and V be generated by the vectors $(0, 0, 0, 1, 0)$ and $(0, 0, 0, 0, 1)$. However, we claim that if $\dim(U) = 3$ and $\dim(V) = 3$, then $\dim(U \cap V) \geq 1$. Indeed, by the Grassmann relation, we have

$$\dim(U) + \dim(V) = \dim(U + V) + \dim(U \cap V),$$

namely

$$3 + 3 = 6 = \dim(U + V) + \dim(U \cap V),$$

and since $U + V$ is a subspace of \mathbb{R}^5 , $\dim(U + V) \leq 5$, which implies

$$6 \leq 5 + \dim(U \cap V),$$

that is $1 \leq \dim(U \cap V)$.

As another consequence of Proposition 4.12, if U and V are two hyperplanes in a vector space of dimension n , so that $\dim(U) = n - 1$ and $\dim(V) = n - 1$, the reader should show that

$$\dim(U \cap V) \geq n - 2,$$

and so, if $U \neq V$, then

$$\dim(U \cap V) = n - 2.$$

Here is a characterization of direct sums that follows directly from Theorem 4.11.

Proposition 4.13. *If U_1, \dots, U_p are any subspaces of a finite dimensional vector space E , then*

$$\dim(U_1 + \dots + U_p) \leq \dim(U_1) + \dots + \dim(U_p),$$

and

$$\dim(U_1 + \dots + U_p) = \dim(U_1) + \dots + \dim(U_p)$$

iff the U_i s form a direct sum $U_1 \oplus \dots \oplus U_p$.

Proof. If we apply Theorem 4.11 to the linear map

$$a: U_1 \times \dots \times U_p \rightarrow U_1 + \dots + U_p$$

given by $a(u_1, \dots, u_p) = u_1 + \dots + u_p$, we get

$$\begin{aligned} \dim(U_1 + \dots + U_p) &= \dim(U_1 \times \dots \times U_p) - \dim(\text{Ker } a) \\ &= \dim(U_1) + \dots + \dim(U_p) - \dim(\text{Ker } a), \end{aligned}$$

so the inequality follows. Since a is injective iff $\text{Ker } a = (0)$, the U_i s form a direct sum iff the second equation holds. \square

Another important corollary of Theorem 4.11 is the following result:

Proposition 4.14. *Let E and F be two vector spaces with the same finite dimension $\dim(E) = \dim(F) = n$. For every linear map $f: E \rightarrow F$, the following properties are equivalent:*

- (a) f is bijective.
- (b) f is surjective.
- (c) f is injective.
- (d) $\text{Ker } f = 0$.

Proof. Obviously, (a) implies (b).

If f is surjective, then $\text{Im } f = F$, and so $\dim(\text{Im } f) = n$. By Theorem 4.11,

$$\dim(E) = \dim(\text{Ker } f) + \dim(\text{Im } f),$$

and since $\dim(E) = n$ and $\dim(\text{Im } f) = n$, we get $\dim(\text{Ker } f) = 0$, which means that $\text{Ker } f = 0$, and so f is injective (see Proposition 2.15). This proves that (b) implies (c).

If f is injective, then by Proposition 2.15, $\text{Ker } f = 0$, so (c) implies (d).

Finally, assume that $\text{Ker } f = 0$, so that $\dim(\text{Ker } f) = 0$ and f is injective (by Proposition 2.15). By Theorem 4.11,

$$\dim(E) = \dim(\text{Ker } f) + \dim(\text{Im } f),$$

and since $\dim(\text{Ker } f) = 0$, we get

$$\dim(\text{Im } f) = \dim(E) = \dim(F),$$

which proves that f is also surjective, and thus bijective. This proves that (d) implies (a) and concludes the proof. \square

One should be warned that Proposition 4.14 fails in infinite dimension.

The following Proposition will also be useful.

Proposition 4.15. *Let E be a vector space. If $E = U \oplus V$ and $E = U \oplus W$, then there is an isomorphism $f: V \rightarrow W$ between V and W .*

Proof. Let R be the relation between V and W , defined such that

$$\langle v, w \rangle \in R \quad \text{iff} \quad w - v \in U.$$

We claim that R is a functional relation that defines a linear isomorphism $f: V \rightarrow W$ between V and W , where $f(v) = w$ iff $\langle v, w \rangle \in R$ (R is the graph of f). If $w - v \in U$ and $w' - v \in U$, then $w' - w \in U$, and since $U \oplus W$ is a direct sum, $U \cap W = 0$, and thus $w' - w = 0$, that is $w' = w$. Thus, R is functional. Similarly, if $w - v \in U$ and $w - v' \in U$, then $v' - v \in U$, and since $U \oplus V$ is a direct sum, $U \cap V = 0$, and $v' = v$. Thus, f is injective. Since $E = U \oplus V$, for every $w \in W$, there exists a unique pair $\langle u, v \rangle \in U \times V$, such that $w = u + v$. Then, $w - v \in U$, and f is surjective. We also need to verify that f is linear. If

$$w - v = u$$

and

$$w' - v' = u',$$

where $u, u' \in U$, then, we have

$$(w + w') - (v + v') = (u + u'),$$

where $u + u' \in U$. Similarly, if

$$w - v = u$$

where $u \in U$, then we have

$$\lambda w - \lambda v = \lambda u,$$

where $\lambda u \in U$. Thus, f is linear. \square

Given a vector space E and any subspace U of E , Proposition 4.15 shows that the dimension of any subspace V such that $E = U \oplus V$ depends only on U . We call $\dim(V)$ the *codimension* of U , and we denote it by $\text{codim}(U)$. A subspace U of codimension 1 is called a *hyperplane*.

The notion of rank of a linear map or of a matrix is an important one, both theoretically and practically, since it is the key to the solvability of linear equations. Recall from Definition 2.15 that the *rank* $\text{rk}(f)$ of a linear map $f: E \rightarrow F$ is the dimension $\dim(\text{Im } f)$ of the image subspace $\text{Im } f$ of F .

We have the following simple proposition.

Proposition 4.16. *Given a linear map $f: E \rightarrow F$, the following properties hold:*

- (i) $\text{rk}(f) = \text{codim}(\text{Ker } f)$.
- (ii) $\text{rk}(f) + \dim(\text{Ker } f) = \dim(E)$.
- (iii) $\text{rk}(f) \leq \min(\dim(E), \dim(F))$.

Proof. Since by Proposition 4.11, $\dim(E) = \dim(\text{Ker } f) + \dim(\text{Im } f)$, and by definition, $\text{rk}(f) = \dim(\text{Im } f)$, we have $\text{rk}(f) = \text{codim}(\text{Ker } f)$. Since $\text{rk}(f) = \dim(\text{Im } f)$, (ii) follows from $\dim(E) = \dim(\text{Ker } f) + \dim(\text{Im } f)$. As for (iii), since $\text{Im } f$ is a subspace of F , we have $\text{rk}(f) \leq \dim(F)$, and since $\text{rk}(f) + \dim(\text{Ker } f) = \dim(E)$, we have $\text{rk}(f) \leq \dim(E)$. \square

The rank of a matrix is defined as follows.

Definition 4.5. Given a $m \times n$ -matrix $A = (a_{ij})$ over the field K , the *rank* $\text{rk}(A)$ of the matrix A is the maximum number of linearly independent columns of A (viewed as vectors in K^m).

In view of Proposition 2.10, the rank of a matrix A is the dimension of the subspace of K^m generated by the columns of A . Let E and F be two vector spaces, and let (u_1, \dots, u_n) be a basis of E , and (v_1, \dots, v_m) a basis of F . Let $f: E \rightarrow F$ be a linear map, and let $M(f)$ be its matrix w.r.t. the bases (u_1, \dots, u_n) and (v_1, \dots, v_m) . Since the rank $\text{rk}(f)$ of f is the dimension of $\text{Im } f$, which is generated by $(f(u_1), \dots, f(u_n))$, the rank of f is the maximum number of linearly independent vectors in $(f(u_1), \dots, f(u_n))$, which is equal to the number of linearly independent columns of $M(f)$, since F and K^m are isomorphic. Thus, we have $\text{rk}(f) = \text{rk}(M(f))$, for every matrix representing f .

We will see later, using duality, that the rank of a matrix A is also equal to the maximal number of linearly independent rows of A .

If U is a hyperplane, then $E = U \oplus V$ for some subspace V of dimension 1. However, a subspace V of dimension 1 is generated by any nonzero vector $v \in V$, and thus we denote V by Kv , and we write $E = U \oplus Kv$. Clearly, $v \notin U$. Conversely, let $x \in E$ be a vector such that $x \notin U$ (and thus, $x \neq 0$). We claim that $E = U \oplus Kx$. Indeed, since U is a hyperplane, we have $E = U \oplus Kv$ for some $v \notin U$ (with $v \neq 0$). Then, $x \in E$ can be written in a unique way as $x = u + \lambda v$, where $u \in U$, and since $x \notin U$, we must have $\lambda \neq 0$, and thus, $v = -\lambda^{-1}u + \lambda^{-1}x$. Since $E = U \oplus Kv$, this shows that $E = U \oplus Kx$. Since $x \notin U$,

we have $U \cap Kx = 0$, and thus $E = U \oplus Kx$. This argument shows that a hyperplane is a maximal proper subspace H of E .

In the next section, we shall see that hyperplanes are precisely the Kernels of nonnull linear maps $f: E \rightarrow K$, called linear forms.

4.2 The Dual Space E^* and Linear Forms

We already observed that the field K itself is a vector space (over itself). The vector space $\text{Hom}(E, K)$ of linear maps from E to the field K , the linear forms, plays a particular role. We take a quick look at the connection between E and $\text{Hom}(E, K)$, its *dual space*. As we will see shortly, every linear map $f: E \rightarrow F$ gives rise to a linear map $f^\top: F^* \rightarrow E^*$, and it turns out that in a suitable basis, the matrix of f^\top is the transpose of the matrix of f . Thus, the notion of dual space provides a conceptual explanation of the phenomena associated with transposition. But it does more, because it allows us to view subspaces as solutions of sets of linear equations and vice-versa.

Consider the following set of two “linear equations” in \mathbb{R}^3 ,

$$x - y + z = 0$$

$$x - y - z = 0,$$

and let us find out what is their set V of common solutions $(x, y, z) \in \mathbb{R}^3$. By subtracting the second equation from the first, we get $2z = 0$, and by adding the two equations, we find that $2(x - y) = 0$, so the set V of solutions is given by

$$y = x$$

$$z = 0.$$

This is a one dimensional subspace of \mathbb{R}^3 . Geometrically, this is the line of equation $y = x$ in the plane $z = 0$.

Now, why did we say that the above equations are linear? This is because, as functions of (x, y, z) , both maps $f_1: (x, y, z) \mapsto x - y + z$ and $f_2: (x, y, z) \mapsto x - y - z$ are linear. The set of all such linear functions from \mathbb{R}^3 to \mathbb{R} is a vector space; we used this fact to form linear combinations of the “equations” f_1 and f_2 . Observe that the dimension of the subspace V is 1. The ambient space has dimension $n = 3$ and there are two “independent” equations f_1, f_2 , so it appears that the dimension $\dim(V)$ of the subspace V defined by m independent equations is

$$\dim(V) = n - m,$$

which is indeed a general fact.

More generally, in \mathbb{R}^n , a linear equation is determined by an n -tuple $(a_1, \dots, a_n) \in \mathbb{R}^n$, and the solutions of this linear equation are given by the n -tuples $(x_1, \dots, x_n) \in \mathbb{R}^n$ such that

$$a_1x_1 + \dots + a_nx_n = 0;$$

these solutions constitute the kernel of the linear map $(x_1, \dots, x_n) \mapsto a_1x_1 + \dots + a_nx_n$. The above considerations assume that we are working in the canonical basis (e_1, \dots, e_n) of \mathbb{R}^n , but we can define “linear equations” independently of bases and in any dimension, by viewing them as elements of the vector space $\text{Hom}(E, K)$ of linear maps from E to the field K .

Definition 4.6. Given a vector space E , the vector space $\text{Hom}(E, K)$ of linear maps from E to the field K is called the *dual space (or dual)* of E . The space $\text{Hom}(E, K)$ is also denoted by E^* , and the linear maps in E^* are called the *linear forms*, or *covectors*. The dual space E^{**} of the space E^* is called the *bidual* of E .

As a matter of notation, linear forms $f: E \rightarrow K$ will also be denoted by starred symbol, such as u^* , x^* , etc.

If E is a vector space of finite dimension n and (u_1, \dots, u_n) is a basis of E , for any linear form $f^* \in E^*$, for every $x = x_1u_1 + \dots + x_nu_n \in E$, we have

$$f^*(x) = \lambda_1x_1 + \dots + \lambda_nx_n,$$

where $\lambda_i = f^*(u_i) \in K$, for every i , $1 \leq i \leq n$. Thus, with respect to the basis (u_1, \dots, u_n) , $f^*(x)$ is a linear combination of the coordinates of x , and we can view a linear form as a *linear equation*, as discussed earlier.

Given a linear form $u^* \in E^*$ and a vector $v \in E$, the result $u^*(v)$ of applying u^* to v is also denoted by $\langle u^*, v \rangle$. This defines a binary operation $\langle -, - \rangle: E^* \times E \rightarrow K$ satisfying the following properties:

$$\begin{aligned} \langle u_1^* + u_2^*, v \rangle &= \langle u_1^*, v \rangle + \langle u_2^*, v \rangle \\ \langle u^*, v_1 + v_2 \rangle &= \langle u^*, v_1 \rangle + \langle u^*, v_2 \rangle \\ \langle \lambda u^*, v \rangle &= \lambda \langle u^*, v \rangle \\ \langle u^*, \lambda v \rangle &= \lambda \langle u^*, v \rangle. \end{aligned}$$

The above identities mean that $\langle -, - \rangle$ is a *bilinear map*, since it is linear in each argument. It is often called the *canonical pairing* between E^* and E . In view of the above identities, given any fixed vector $v \in E$, the map $\text{eval}_v: E^* \rightarrow K$ (*evaluation at v*) defined such that

$$\text{eval}_v(u^*) = \langle u^*, v \rangle = u^*(v) \quad \text{for every } u^* \in E^*$$

is a linear map from E^* to K , that is, eval_v is a linear form in E^{**} . Again, from the above identities, the map $\text{eval}_E: E \rightarrow E^{**}$, defined such that

$$\text{eval}_E(v) = \text{eval}_v \quad \text{for every } v \in E,$$

is a linear map. Observe that

$$\text{eval}_E(v)(u^*) = \langle u^*, v \rangle = u^*(v), \quad \text{for all } v \in E \text{ and all } u^* \in E^*.$$

We shall see that the map eval_E is injective, and that it is an isomorphism when E has finite dimension.

We now formalize the notion of the set V^0 of linear equations vanishing on all vectors in a given subspace $V \subseteq E$, and the notion of the set U^0 of common solutions of a given set $U \subseteq E^*$ of linear equations. The duality theorem (Theorem 4.17) shows that the dimensions of V and V^0 , and the dimensions of U and U^0 , are related in a crucial way. It also shows that, in finite dimension, the maps $V \mapsto V^0$ and $U \mapsto U^0$ are inverse bijections from subspaces of E to subspaces of E^* .

Definition 4.7. Given a vector space E and its dual E^* , we say that a vector $v \in E$ and a linear form $u^* \in E^*$ are *orthogonal* if $\langle u^*, v \rangle = 0$. Given a subspace V of E and a subspace U of E^* , we say that V and U are *orthogonal* if $\langle u^*, v \rangle = 0$ for every $u^* \in U$ and every $v \in V$. Given a subset V of E (resp. a subset U of E^*), the *orthogonal* V^0 of V is the subspace V^0 of E^* defined such that

$$V^0 = \{u^* \in E^* \mid \langle u^*, v \rangle = 0, \text{ for every } v \in V\}$$

(resp. the *orthogonal* U^0 of U is the subspace U^0 of E defined such that

$$U^0 = \{v \in E \mid \langle u^*, v \rangle = 0, \text{ for every } u^* \in U\}.$$

The subspace $V^0 \subseteq E^*$ is also called the *annihilator* of V . The subspace $U^0 \subseteq E$ annihilated by $U \subseteq E^*$ does not have a special name. It seems reasonable to call it the *linear subspace (or linear variety) defined by U* .

Informally, V^0 is the *set of linear equations that vanish on V* , and U^0 is the *set of common zeros of all linear equations in U* .

We can also define V^0 by

$$V^0 = \{u^* \in E^* \mid V \subseteq \text{Ker } u^*\}$$

and U^0 by

$$U^0 = \bigcap_{u^* \in U} \text{Ker } u^*.$$

Observe that $E^0 = 0$, and $\{0\}^0 = E^*$. Furthermore, if $V_1 \subseteq V_2 \subseteq E$, then $V_2^0 \subseteq V_1^0 \subseteq E^*$, and if $U_1 \subseteq U_2 \subseteq E^*$, then $U_2^0 \subseteq U_1^0 \subseteq E$.

Indeed, if $V_1 \subseteq V_2 \subseteq E$, then for any $f^* \in V_2^0$ we have $f^*(v) = 0$ for all $v \in V_2$, and thus $f^*(v) = 0$ for all $v \in V_1$, so $f^* \in V_1^0$. Similarly, if $U_1 \subseteq U_2 \subseteq E^*$, then for any $v \in U_2^0$, we have $f^*(v) = 0$ for all $f^* \in U_2$, so $f^*(v) = 0$ for all $f^* \in U_1$, which means that $v \in U_1^0$.

Here are some examples. Let $E = M_2(\mathbb{R})$, the space of real 2×2 matrices, and let V be the subspace of $M_2(\mathbb{R})$ spanned by the matrices

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}, \quad \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}.$$

We check immediately that the subspace V consists of all matrices of the form

$$\begin{pmatrix} b & a \\ a & c \end{pmatrix},$$

that is, all symmetric matrices. The matrices

$$\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$$

in V satisfy the equation

$$a_{12} - a_{21} = 0,$$

and all scalar multiples of these equations, so V^0 is the subspace of E^* spanned by the linear form given by $u^*(a_{11}, a_{12}, a_{21}, a_{22}) = a_{12} - a_{21}$. We have

$$\dim(V^0) = \dim(E) - \dim(V) = 4 - 3 = 1.$$

The above example generalizes to $E = M_n(\mathbb{R})$ for any $n \geq 1$, but this time, consider the space U of linear forms asserting that a matrix A is symmetric; these are the linear forms spanned by the $n(n-1)/2$ equations

$$a_{ij} - a_{ji} = 0, \quad 1 \leq i < j \leq n;$$

Note there are no constraints on diagonal entries, and half of the equations

$$a_{ij} - a_{ji} = 0, \quad 1 \leq i \neq j \leq n$$

are redundant. It is easy to check that the equations (linear forms) for which $i < j$ are linearly independent. To be more precise, let U be the space of linear forms in E^* spanned by the linear forms

$$u_{ij}^*(a_{11}, \dots, a_{1n}, a_{21}, \dots, a_{2n}, \dots, a_{n1}, \dots, a_{nn}) = a_{ij} - a_{ji}, \quad 1 \leq i < j \leq n.$$

Then, the set U^0 of common solutions of these equations is the space $\mathbf{S}(n)$ of symmetric matrices. This space has dimension

$$\frac{n(n+1)}{2} = n^2 - \frac{n(n-1)}{2}.$$

We leave it as an exercise to find a basis of $\mathbf{S}(n)$.

If $E = M_n(\mathbb{R})$, consider the subspace U of linear forms in E^* spanned by the linear forms

$$\begin{aligned} u_{ij}^*(a_{11}, \dots, a_{1n}, a_{21}, \dots, a_{2n}, \dots, a_{n1}, \dots, a_{nn}) &= a_{ij} + a_{ji}, \quad 1 \leq i < j \leq n \\ u_{ii}^*(a_{11}, \dots, a_{1n}, a_{21}, \dots, a_{2n}, \dots, a_{n1}, \dots, a_{nn}) &= a_{ii}, \quad 1 \leq i \leq n. \end{aligned}$$

It is easy to see that these linear forms are linearly independent, so $\dim(U) = n(n+1)/2$. The space U^0 of matrices $A \in M_n(\mathbb{R})$ satisfying all of the above equations is clearly the space **Skew**(n) of skew-symmetric matrices. The dimension of U^0 is

$$\frac{n(n-1)}{2} = n^2 - \frac{n(n+1)}{2}.$$

We leave it as an exercise to find a basis of **Skew**(n).

For yet another example, with $E = M_n(\mathbb{R})$, for any $A \in M_n(\mathbb{R})$, consider the linear form in E^* given by

$$\text{tr}(A) = a_{11} + a_{22} + \cdots + a_{nn},$$

called the *trace* of A . The subspace U^0 of E consisting of all matrices A such that $\text{tr}(A) = 0$ is a space of dimension $n^2 - 1$. We leave it as an exercise to find a basis of this space.

The dimension equations

$$\begin{aligned}\dim(V) + \dim(V^0) &= \dim(E) \\ \dim(U) + \dim(U^0) &= \dim(E)\end{aligned}$$

are always true (if E is finite-dimensional). This is part of the duality theorem (Theorem 4.17).

In contrast with the previous examples, given a matrix $A \in M_n(\mathbb{R})$, the equations asserting that $A^\top A = I$ are not linear constraints. For example, for $n = 2$, we have

$$\begin{aligned}a_{11}^2 + a_{21}^2 &= 1 \\ a_{21}^2 + a_{22}^2 &= 1 \\ a_{11}a_{12} + a_{21}a_{22} &= 0.\end{aligned}$$

Remarks:

- (1) The notation V^0 (resp. U^0) for the orthogonal of a subspace V of E (resp. a subspace U of E^*) is not universal. Other authors use the notation V^\perp (resp. U^\perp). However, the notation V^\perp is also used to denote the orthogonal complement of a subspace V with respect to an inner product on a space E , in which case V^\perp is a subspace of E and not a subspace of E^* (see Chapter 10). To avoid confusion, we prefer using the notation V^0 .
- (2) Since linear forms can be viewed as linear equations (at least in finite dimension), given a subspace (or even a subset) U of E^* , we can define the set $\mathcal{Z}(U)$ of *common zeros* of the equations in U by

$$\mathcal{Z}(U) = \{v \in E \mid u^*(v) = 0, \text{ for all } u^* \in U\}.$$

Of course $\mathcal{Z}(U) = U^0$, but the notion $\mathcal{Z}(U)$ can be generalized to more general kinds of equations, namely polynomial equations. In this more general setting, U is a set of *polynomials* in n variables with coefficients in K (where $n = \dim(E)$). Sets of the form $\mathcal{Z}(U)$ are called *algebraic varieties*. Linear forms correspond to the special case where homogeneous polynomials of degree 1 are considered.

If V is a subset of E , it is natural to associate with V the *set of polynomials in $K[X_1, \dots, X_n]$ that vanish on V* . This set, usually denoted $\mathcal{I}(V)$, has some special properties that make it an *ideal*. If V is a linear subspace of E , it is natural to restrict our attention to the space V^0 of linear forms that vanish on V , and in this case we identify $\mathcal{I}(V)$ and V^0 (although technically, $\mathcal{I}(V)$ is no longer an ideal).

For any arbitrary set of polynomials $U \subseteq K[X_1, \dots, X_n]$ (resp $V \subseteq E$) the relationship between $\mathcal{I}(\mathcal{Z}(U))$ and U (resp. $\mathcal{Z}(\mathcal{I}(V))$ and V) is generally not simple, even though we always have

$$U \subseteq \mathcal{I}(\mathcal{Z}(U)) \quad (\text{resp.} \quad V \subseteq \mathcal{Z}(\mathcal{I}(V))).$$

However, when the field K is algebraically closed, then $\mathcal{I}(\mathcal{Z}(U))$ is equal to the *radical* of the ideal U , a famous result due to Hilbert known as the *Nullstellensatz* (see Lang [67] or Dummit and Foote [32]). The study of algebraic varieties is the main subject of *algebraic geometry*, a beautiful but formidable subject. For a taste of algebraic geometry, see Lang [67] or Dummit and Foote [32].

The duality theorem (Theorem 4.17) shows that the situation is much simpler if we restrict our attention to linear subspaces; in this case

$$U = \mathcal{I}(\mathcal{Z}(U)) \quad \text{and} \quad V = \mathcal{Z}(\mathcal{I}(V)).$$

We claim that $V \subseteq V^{00}$ for every subspace V of E , and that $U \subseteq U^{00}$ for every subspace U of E^* .

Indeed, for any $v \in V$, to show that $v \in V^{00}$ we need to prove that $u^*(v) = 0$ for all $u^* \in V^0$. However, V^0 consists of all linear forms u^* such that $u^*(y) = 0$ for *all* $y \in V$; in particular, since $v \in V$, $u^*(v) = 0$ for all $u^* \in V^0$, as required.

Similarly, for any $u^* \in U$, to show that $u^* \in U^{00}$ we need to prove that $u^*(v) = 0$ for all $v \in U^0$. However, U^0 consists of all vectors v such that $f^*(v) = 0$ for *all* $f^* \in U$; in particular, since $u^* \in U$, $u^*(v) = 0$ for all $v \in U^0$, as required.

We will see shortly that in finite dimension, we have $V = V^{00}$ and $U = U^{00}$.



However, even though $V = V^{00}$ is always true, when E is of infinite dimension, it is not always true that $U = U^{00}$.

Given a vector space E and any basis $(u_i)_{i \in I}$ for E , we can associate to each u_i a linear form $u_i^* \in E^*$, and the u_i^* have some remarkable properties.

Definition 4.8. Given a vector space E and any basis $(u_i)_{i \in I}$ for E , by Proposition 2.16, for every $i \in I$, there is a unique linear form u_i^* such that

$$u_i^*(u_j) = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j, \end{cases}$$

for every $j \in I$. The linear form u_i^* is called the *coordinate form* of index i w.r.t. the basis $(u_i)_{i \in I}$.

Given an index set I , authors often define the so called “Kronecker symbol” δ_{ij} , such that

$$\delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j, \end{cases}$$

for all $i, j \in I$. Then, $u_i^*(u_j) = \delta_{ij}$.

The reason for the terminology *coordinate form* is as follows: If E has finite dimension and if (u_1, \dots, u_n) is a basis of E , for any vector

$$v = \lambda_1 u_1 + \dots + \lambda_n u_n,$$

we have

$$\begin{aligned} u_i^*(v) &= u_i^*(\lambda_1 u_1 + \dots + \lambda_n u_n) \\ &= \lambda_1 u_i^*(u_1) + \dots + \lambda_i u_i^*(u_i) + \dots + \lambda_n u_i^*(u_n) \\ &= \lambda_i, \end{aligned}$$

since $u_i^*(u_j) = \delta_{ij}$. Therefore, u_i^* is the linear function that returns the i th coordinate of a vector expressed over the basis (u_1, \dots, u_n) .

Given a vector space E and a subspace U of E , by Theorem 2.9, every basis $(u_i)_{i \in I}$ of U can be extended to a basis $(u_j)_{j \in I \cup J}$ of E , where $I \cap J = \emptyset$. We have the following important theorem adapted from E. Artin [3] (Chapter 1).

Theorem 4.17. (*Duality theorem*) *Let E be a vector space. The following properties hold:*

- (a) *For every basis $(u_i)_{i \in I}$ of E , the family $(u_i^*)_{i \in I}$ of coordinate forms is linearly independent.*
- (b) *For every subspace V of E , we have $V^{00} = V$.*
- (c) *For every subspace V of finite codimension m of E , for every subspace W of E such that $E = V \oplus W$ (where W is of finite dimension m), for every basis $(u_i)_{i \in I}$ of E such that (u_1, \dots, u_m) is a basis of W , the family (u_1^*, \dots, u_m^*) is a basis of the orthogonal V^0 of V in E^* , so that*

$$\dim(V^0) = \text{codim}(V).$$

Furthermore, we have $V^{00} = V$.

(d) For every subspace U of finite dimension m of E^* , the orthogonal U^0 of U in E is of finite codimension m , so that

$$\text{codim}(U^0) = \dim(U).$$

Furthermore, $U^{00} = U$.

Proof. (a) Assume that

$$\sum_{i \in I} \lambda_i u_i^* = 0,$$

for a family $(\lambda_i)_{i \in I}$ (of scalars in K). Since $(\lambda_i)_{i \in I}$ has finite support, there is a finite subset J of I such that $\lambda_i = 0$ for all $i \in I - J$, and we have

$$\sum_{j \in J} \lambda_j u_j^* = 0.$$

Applying the linear form $\sum_{j \in J} \lambda_j u_j^*$ to each u_j ($j \in J$), by Definition 4.8, since $u_i^*(u_j) = 1$ if $i = j$ and 0 otherwise, we get $\lambda_j = 0$ for all $j \in J$, that is $\lambda_i = 0$ for all $i \in I$ (by definition of J as the support). Thus, $(u_i^*)_{i \in I}$ is linearly independent.

(b) Clearly, we have $V \subseteq V^{00}$. If $V \neq V^{00}$, then let $(u_i)_{i \in I \cup J}$ be a basis of V^{00} such that $(u_i)_{i \in I}$ is a basis of V (where $I \cap J = \emptyset$). Since $V \neq V^{00}$, $u_{j_0} \in V^{00}$ for some $j_0 \in J$ (and thus, $j_0 \notin I$). Since $u_{j_0} \in V^{00}$, u_{j_0} is orthogonal to every linear form in V^0 . Now, we have $u_{j_0}^*(u_i) = 0$ for all $i \in I$, and thus $u_{j_0}^* \in V^0$. However, $u_{j_0}^*(u_{j_0}) = 1$, contradicting the fact that u_{j_0} is orthogonal to every linear form in V^0 . Thus, $V = V^{00}$.

(c) Let $J = I - \{1, \dots, m\}$. Every linear form $f^* \in V^0$ is orthogonal to every u_j , for $j \in J$, and thus, $f^*(u_j) = 0$, for all $j \in J$. For such a linear form $f^* \in V^0$, let

$$g^* = f^*(u_1)u_1^* + \dots + f^*(u_m)u_m^*.$$

We have $g^*(u_i) = f^*(u_i)$, for every i , $1 \leq i \leq m$. Furthermore, by definition, g^* vanishes on all u_j , where $j \in J$. Thus, f^* and g^* agree on the basis $(u_i)_{i \in I}$ of E , and so, $g^* = f^*$. This shows that (u_1^*, \dots, u_m^*) generates V^0 , and since it is also a linearly independent family, (u_1^*, \dots, u_m^*) is a basis of V^0 . It is then obvious that $\dim(V^0) = \text{codim}(V)$, and by part (b), we have $V^{00} = V$.

(d) Let (u_1^*, \dots, u_m^*) be a basis of U . Note that the map $h: E \rightarrow K^m$ defined such that

$$h(v) = (u_1^*(v), \dots, u_m^*(v))$$

for every $v \in E$, is a linear map, and that its kernel $\text{Ker } h$ is precisely U^0 . Then, by Proposition 4.11,

$$E \approx \text{Ker}(h) \oplus \text{Im } h = U^0 \oplus \text{Im } h,$$

and since $\dim(\text{Im } h) \leq m$, we deduce that U^0 is a subspace of E of finite codimension at most m , and by (c), we have $\dim(U^{00}) = \text{codim}(U^0) \leq m = \dim(U)$. However, it is clear that $U \subseteq U^{00}$, which implies $\dim(U) \leq \dim(U^{00})$, and so $\dim(U^{00}) = \dim(U) = m$, and we must have $U = U^{00}$. \square

Part (a) of Theorem 4.17 shows that

$$\dim(E) \leq \dim(E^*).$$

When E is of finite dimension n and (u_1, \dots, u_n) is a basis of E , by part (c), the family (u_1^*, \dots, u_n^*) is a basis of the dual space E^* , called the *dual basis* of (u_1, \dots, u_n) .

By part (c) and (d) of theorem 4.17, the maps $V \mapsto V^0$ and $U \mapsto U^0$, where V is a subspace of finite codimension of E and U is a subspace of finite dimension of E^* , are inverse bijections. These maps set up a *duality* between subspaces of finite codimension of E , and subspaces of finite dimension of E^* .



One should be careful that this bijection does not extend to subspaces of E^* of infinite dimension.



When E is of infinite dimension, for every basis $(u_i)_{i \in I}$ of E , the family $(u_i^*)_{i \in I}$ of coordinate forms is never a basis of E^* . It is linearly independent, but it is “too small” to generate E^* . For example, if $E = \mathbb{R}^{(\mathbb{N})}$, where $\mathbb{N} = \{0, 1, 2, \dots\}$, the map $f: E \rightarrow \mathbb{R}$ that sums the nonzero coordinates of a vector in E is a linear form, but it is easy to see that it cannot be expressed as a linear combination of coordinate forms. As a consequence, when E is of infinite dimension, E and E^* are not isomorphic.

Here is another example illustrating the power of Theorem 4.17. Let $E = M_n(\mathbb{R})$, and consider the equations asserting that the sum of the entries in every row of a matrix $\in M_n(\mathbb{R})$ is equal to the same number. We have $n - 1$ equations

$$\sum_{j=1}^n (a_{ij} - a_{i+1j}) = 0, \quad 1 \leq i \leq n - 1,$$

and it is easy to see that they are linearly independent. Therefore, the space U of linear forms in E^* spanned by the above linear forms (equations) has dimension $n - 1$, and the space U^0 of matrices satisfying all these equations has dimension $n^2 - n + 1$. It is not so obvious to find a basis for this space.

When E is of finite dimension n and (u_1, \dots, u_n) is a basis of E , we noted that the family (u_1^*, \dots, u_n^*) is a basis of the dual space E^* (called the dual basis of (u_1, \dots, u_n)). Let us see how the coordinates of a linear form φ^* over the dual basis (u_1^*, \dots, u_n^*) vary under a change of basis.

Let (u_1, \dots, u_n) and (v_1, \dots, v_n) be two bases of E , and let $P = (a_{ij})$ be the change of basis matrix from (u_1, \dots, u_n) to (v_1, \dots, v_n) , so that

$$v_j = \sum_{i=1}^n a_{ij} u_i,$$

and let $P^{-1} = (b_{ij})$ be the inverse of P , so that

$$u_i = \sum_{j=1}^n b_{ji} v_j.$$

Since $u_i^*(u_j) = \delta_{ij}$ and $v_i^*(v_j) = \delta_{ij}$, we get

$$v_j^*(u_i) = v_j^*\left(\sum_{k=1}^n b_{ki} v_k\right) = b_{ji},$$

and thus

$$v_j^* = \sum_{i=1}^n b_{ji} u_i^*,$$

and

$$u_i^* = \sum_{j=1}^n a_{ij} v_j^*.$$

This means that the change of basis from the dual basis (u_1^*, \dots, u_n^*) to the dual basis (v_1^*, \dots, v_n^*) is $(P^{-1})^\top$. Since

$$\varphi^* = \sum_{i=1}^n \varphi_i u_i^* = \sum_{i=1}^n \varphi'_i v_i^*,$$

we get

$$\varphi'_j = \sum_{i=1}^n a_{ij} \varphi_i,$$

so the new coordinates φ'_j are expressed in terms of the old coordinates φ_i using the matrix P^\top . If we use the row vectors $(\varphi_1, \dots, \varphi_n)$ and $(\varphi'_1, \dots, \varphi'_n)$, we have

$$(\varphi'_1, \dots, \varphi'_n) = (\varphi_1, \dots, \varphi_n)P.$$

Comparing with the change of basis

$$v_j = \sum_{i=1}^n a_{ij} u_i,$$

we note that this time, the coordinates (φ_i) of the linear form φ^* change in the *same direction* as the change of basis. For this reason, we say that the coordinates of linear forms are *covariant*. By abuse of language, it is often said that linear forms are *covariant*, which explains why the term *covector* is also used for a linear form.

Observe that if (e_1, \dots, e_n) is a basis of the vector space E , then, as a linear map from E to K , every linear form $f \in E^*$ is represented by a $1 \times n$ matrix, that is, by a *row vector*

$$(\lambda_1, \dots, \lambda_n),$$

with respect to the basis (e_1, \dots, e_n) of E , and 1 of K , where $f(e_i) = \lambda_i$. A vector $u = \sum_{i=1}^n u_i e_i \in E$ is represented by a $n \times 1$ matrix, that is, by a *column vector*

$$\begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix},$$

and the action of f on u , namely $f(u)$, is represented by the matrix product

$$\begin{pmatrix} \lambda_1 & \cdots & \lambda_n \end{pmatrix} \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix} = \lambda_1 u_1 + \cdots + \lambda_n u_n.$$

On the other hand, with respect to the dual basis (e_1^*, \dots, e_n^*) of E^* , the linear form f is represented by the column vector

$$\begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_n \end{pmatrix}.$$

Remark: In many texts using tensors, vectors are often indexed with lower indices. If so, it is more convenient to write the coordinates of a vector x over the basis (u_1, \dots, u_n) as (x^i) , using an upper index, so that

$$x = \sum_{i=1}^n x^i u_i,$$

and in a change of basis, we have

$$v_j = \sum_{i=1}^n a_j^i u_i$$

and

$$x^i = \sum_{j=1}^n a_j^i x'^j.$$

Dually, linear forms are indexed with upper indices. Then, it is more convenient to write the coordinates of a covector φ^* over the dual basis (u^{*1}, \dots, u^{*n}) as (φ_i) , using a lower index, so that

$$\varphi^* = \sum_{i=1}^n \varphi_i u^{*i}$$

and in a change of basis, we have

$$u^{*i} = \sum_{j=1}^n a_j^i v^{*j}$$

and

$$\varphi'_j = \sum_{i=1}^n a_j^i \varphi_i.$$

With these conventions, the index of summation appears once in upper position and once in lower position, and the summation sign can be safely omitted, a trick due to *Einstein*. For example, we can write

$$\varphi'_j = a_j^i \varphi_i$$

as an abbreviation for

$$\varphi'_j = \sum_{i=1}^n a_j^i \varphi_i.$$

For another example of the use of Einstein's notation, if the vectors (v_1, \dots, v_n) are linear combinations of the vectors (u_1, \dots, u_n) , with

$$v_i = \sum_{j=1}^n a_{ij} u_j, \quad 1 \leq i \leq n,$$

then the above equations are written as

$$v_i = a_i^j u_j, \quad 1 \leq i \leq n.$$

Thus, in Einstein's notation, the $n \times n$ matrix (a_{ij}) is denoted by (a_i^j) , a $(1, 1)$ -tensor.



Beware that some authors view a matrix as a mapping between *coordinates*, in which case the matrix (a_{ij}) is denoted by (a_j^i) .

We will now pin down the relationship between a vector space E and its bidual E^{**} .

Proposition 4.18. *Let E be a vector space. The following properties hold:*

(a) *The linear map $\text{eval}_E: E \rightarrow E^{**}$ defined such that*

$$\text{eval}_E(v) = \text{eval}_v \quad \text{for all } v \in E,$$

that is, $\text{eval}_E(v)(u^) = \langle u^*, v \rangle = u^*(v)$ for every $u^* \in E^*$, is injective.*

(b) *When E is of finite dimension n , the linear map $\text{eval}_E: E \rightarrow E^{**}$ is an isomorphism (called the canonical isomorphism).*

Proof. (a) Let $(u_i)_{i \in I}$ be a basis of E , and let $v = \sum_{i \in I} v_i u_i$. If $\text{eval}_E(v) = 0$, then in particular, $\text{eval}_E(v)(u_i^*) = 0$ for all u_i^* , and since

$$\text{eval}_E(v)(u_i^*) = \langle u_i^*, v \rangle = v_i,$$

we have $v_i = 0$ for all $i \in I$, that is, $v = 0$, showing that $\text{eval}_E: E \rightarrow E^{**}$ is injective.

If E is of finite dimension n , by Theorem 4.17, for every basis (u_1, \dots, u_n) , the family (u_1^*, \dots, u_n^*) is a basis of the dual space E^* , and thus the family $(u_1^{**}, \dots, u_n^{**})$ is a basis of the bidual E^{**} . This shows that $\dim(E) = \dim(E^{**}) = n$, and since by part (a), we know that $\text{eval}_E: E \rightarrow E^{**}$ is injective, in fact, $\text{eval}_E: E \rightarrow E^{**}$ is bijective (because an injective map carries a linearly independent family to a linearly independent family, and in a vector space of dimension n , a linearly independent family of n vectors is a basis, see Proposition 2.10). \square



When a vector space E has infinite dimension, E and its bidual E^{**} are never isomorphic.

When E is of finite dimension and (u_1, \dots, u_n) is a basis of E , in view of the canonical isomorphism $\text{eval}_E: E \rightarrow E^{**}$, the basis $(u_1^{**}, \dots, u_n^{**})$ of the bidual is identified with (u_1, \dots, u_n) .

Proposition 4.18 can be reformulated very fruitfully in terms of pairings (adapted from E. Artin [3], Chapter 1). Given two vector spaces E and F over a field K , we say that a function $\varphi: E \times F \rightarrow K$ is *bilinear* if for every $v \in V$, the map $u \mapsto \varphi(u, v)$ (from E to K) is linear, and for every $u \in E$, the map $v \mapsto \varphi(u, v)$ (from F to K) is linear.

Definition 4.9. Given two vector spaces E and F over K , a *pairing between E and F* is a bilinear map $\varphi: E \times F \rightarrow K$. Such a pairing is *nondegenerate* iff

- (1) for every $u \in E$, if $\varphi(u, v) = 0$ for all $v \in F$, then $u = 0$, and
- (2) for every $v \in F$, if $\varphi(u, v) = 0$ for all $u \in E$, then $v = 0$.

A pairing $\varphi: E \times F \rightarrow K$ is often denoted by $\langle -, - \rangle: E \times F \rightarrow K$. For example, the map $\langle -, - \rangle: E^* \times E \rightarrow K$ defined earlier is a nondegenerate pairing (use the proof of (a) in Proposition 4.18).

Given a pairing $\varphi: E \times F \rightarrow K$, we can define two maps $l_\varphi: E \rightarrow F^*$ and $r_\varphi: F \rightarrow E^*$ as follows: For every $u \in E$, we define the linear form $l_\varphi(u)$ in F^* such that

$$l_\varphi(u)(y) = \varphi(u, y) \quad \text{for every } y \in F,$$

and for every $v \in F$, we define the linear form $r_\varphi(v)$ in E^* such that

$$r_\varphi(v)(x) = \varphi(x, v) \quad \text{for every } x \in E.$$

We have the following useful proposition.

Proposition 4.19. *Given two vector spaces E and F over K , for every nondegenerate pairing $\varphi: E \times F \rightarrow K$ between E and F , the maps $l_\varphi: E \rightarrow F^*$ and $r_\varphi: F \rightarrow E^*$ are linear and injective. Furthermore, if E and F have finite dimension, then this dimension is the same and $l_\varphi: E \rightarrow F^*$ and $r_\varphi: F \rightarrow E^*$ are bijections.*

Proof. The maps $l_\varphi: E \rightarrow F^*$ and $r_\varphi: F \rightarrow E^*$ are linear because a pairing is bilinear. If $l_\varphi(u) = 0$ (the null form), then

$$l_\varphi(u)(v) = \varphi(u, v) = 0 \quad \text{for every } v \in F,$$

and since φ is nondegenerate, $u = 0$. Thus, $l_\varphi: E \rightarrow F^*$ is injective. Similarly, $r_\varphi: F \rightarrow E^*$ is injective. When F has finite dimension n , we have seen that F and F^* have the same dimension. Since $l_\varphi: E \rightarrow F^*$ is injective, we have $m = \dim(E) \leq \dim(F) = n$. The same argument applies to E , and thus $n = \dim(F) \leq \dim(E) = m$. But then, $\dim(E) = \dim(F)$, and $l_\varphi: E \rightarrow F^*$ and $r_\varphi: F \rightarrow E^*$ are bijections. \square

When E has finite dimension, the nondegenerate pairing $\langle -, - \rangle: E^* \times E \rightarrow K$ yields another proof of the existence of a natural isomorphism between E and E^{**} . Interesting nondegenerate pairings arise in exterior algebra. We now show the relationship between hyperplanes and linear forms.

4.3 Hyperplanes and Linear Forms

Actually, Proposition 4.20 below follows from parts (c) and (d) of Theorem 4.17, but we feel that it is also interesting to give a more direct proof.

Proposition 4.20. *Let E be a vector space. The following properties hold:*

- (a) *Given any nonnull linear form $f^* \in E^*$, its kernel $H = \text{Ker } f^*$ is a hyperplane.*
- (b) *For any hyperplane H in E , there is a (nonnull) linear form $f^* \in E^*$ such that $H = \text{Ker } f^*$.*
- (c) *Given any hyperplane H in E and any (nonnull) linear form $f^* \in E^*$ such that $H = \text{Ker } f^*$, for every linear form $g^* \in E^*$, $H = \text{Ker } g^*$ iff $g^* = \lambda f^*$ for some $\lambda \neq 0$ in K .*

Proof. (a) If $f^* \in E^*$ is nonnull, there is some vector $v_0 \in E$ such that $f^*(v_0) \neq 0$. Let $H = \text{Ker } f^*$. For every $v \in E$, we have

$$f^* \left(v - \frac{f^*(v)}{f^*(v_0)} v_0 \right) = f^*(v) - \frac{f^*(v)}{f^*(v_0)} f^*(v_0) = f^*(v) - f^*(v) = 0.$$

Thus,

$$v - \frac{f^*(v)}{f^*(v_0)} v_0 = h \in H,$$

and

$$v = h + \frac{f^*(v)}{f^*(v_0)} v_0,$$

that is, $E = H + K v_0$. Also, since $f^*(v_0) \neq 0$, we have $v_0 \notin H$, that is, $H \cap K v_0 = 0$. Thus, $E = H \oplus K v_0$, and H is a hyperplane.

(b) If H is a hyperplane, $E = H \oplus Kv_0$ for some $v_0 \notin H$. Then, every $v \in E$ can be written in a unique way as $v = h + \lambda v_0$. Thus, there is a well-defined function $f^*: E \rightarrow K$, such that, $f^*(v) = \lambda$, for every $v = h + \lambda v_0$. We leave as a simple exercise the verification that f^* is a linear form. Since $f^*(v_0) = 1$, the linear form f^* is nonnull. Also, by definition, it is clear that $\lambda = 0$ iff $v \in H$, that is, $\text{Ker } f^* = H$.

(c) Let H be a hyperplane in E , and let $f^* \in E^*$ be any (nonnull) linear form such that $H = \text{Ker } f^*$. Clearly, if $g^* = \lambda f^*$ for some $\lambda \neq 0$, then $H = \text{Ker } g^*$. Conversely, assume that $H = \text{Ker } g^*$ for some nonnull linear form g^* . From (a), we have $E = H \oplus Kv_0$, for some v_0 such that $f^*(v_0) \neq 0$ and $g^*(v_0) \neq 0$. Then, observe that

$$g^* - \frac{g^*(v_0)}{f^*(v_0)} f^*$$

is a linear form that vanishes on H , since both f^* and g^* vanish on H , but also vanishes on Kv_0 . Thus, $g^* = \lambda f^*$, with

$$\lambda = \frac{g^*(v_0)}{f^*(v_0)}.$$

□

We leave as an exercise the fact that every subspace $V \neq E$ of a vector space E , is the intersection of all hyperplanes that contain V . We now consider the notion of transpose of a linear map and of a matrix.

4.4 Transpose of a Linear Map and of a Matrix

Given a linear map $f: E \rightarrow F$, it is possible to define a map $f^\top: F^* \rightarrow E^*$ which has some interesting properties.

Definition 4.10. Given a linear map $f: E \rightarrow F$, the *transpose* $f^\top: F^* \rightarrow E^*$ of f is the linear map defined such that

$$f^\top(v^*) = v^* \circ f, \quad \text{for every } v^* \in F^*,$$

as shown in the diagram below:

$$\begin{array}{ccc} E & \xrightarrow{f} & F \\ & \searrow f^\top(v^*) & \downarrow v^* \\ & & K. \end{array}$$

Equivalently, the linear map $f^\top: F^* \rightarrow E^*$ is defined such that

$$\langle v^*, f(u) \rangle = \langle f^\top(v^*), u \rangle,$$

for all $u \in E$ and all $v^* \in F^*$.

It is easy to verify that the following properties hold:

$$\begin{aligned}(f + g)^\top &= f^\top + g^\top \\ (g \circ f)^\top &= f^\top \circ g^\top \\ \text{id}_E^\top &= \text{id}_{E^*}.\end{aligned}$$



Note the reversal of composition on the right-hand side of $(g \circ f)^\top = f^\top \circ g^\top$.

The equation $(g \circ f)^\top = f^\top \circ g^\top$ implies the following useful proposition.

Proposition 4.21. *If $f: E \rightarrow F$ is any linear map, then the following properties hold:*

(1) *If f is injective, then f^\top is surjective.*

(2) *If f is surjective, then f^\top is injective.*

Proof. If $f: E \rightarrow F$ is injective, then it has a retraction $r: F \rightarrow E$ such that $r \circ f = \text{id}_E$, and if $f: E \rightarrow F$ is surjective, then it has a section $s: F \rightarrow E$ such that $f \circ s = \text{id}_F$. Now, if $f: E \rightarrow F$ is injective, then we have

$$(r \circ f)^\top = f^\top \circ r^\top = \text{id}_{E^*},$$

which implies that f^\top is surjective, and if f is surjective, then we have

$$(f \circ s)^\top = s^\top \circ f^\top = \text{id}_{F^*},$$

which implies that f^\top is injective. □

We also have the following property showing the naturality of the eval map.

Proposition 4.22. *For any linear map $f: E \rightarrow F$, we have*

$$f^{\top\top} \circ \text{eval}_E = \text{eval}_F \circ f,$$

or equivalently, the following diagram commutes:

$$\begin{array}{ccc} E^{**} & \xrightarrow{f^{\top\top}} & F^{**} \\ \text{eval}_E \uparrow & & \uparrow \text{eval}_F \\ E & \xrightarrow{f} & F. \end{array}$$

Proof. For every $u \in E$ and every $\varphi \in F^{**}$, we have

$$\begin{aligned}(f^{\top\top} \circ \text{eval}_E)(u)(\varphi) &= \langle f^{\top\top}(\text{eval}_E(u)), \varphi \rangle \\ &= \langle \text{eval}_E(u), f^\top(\varphi) \rangle \\ &= \langle f^\top(\varphi), u \rangle \\ &= \langle \varphi, f(u) \rangle \\ &= \langle \text{eval}_F(f(u)), \varphi \rangle \\ &= \langle (\text{eval}_F \circ f)(u), \varphi \rangle \\ &= (\text{eval}_F \circ f)(u)(\varphi),\end{aligned}$$

which proves that $f^{\top\top} \circ \text{eval}_E = \text{eval}_F \circ f$, as claimed. □

If E and F are finite-dimensional, then eval_E and then eval_F are isomorphisms, so Proposition 4.22 shows that if we identify E with its bidual E^{**} and F with its bidual F^{**} then

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

As a corollary of Proposition 4.22, if $\dim(E)$ is finite, then we have

$$\text{Ker}(f^{\top\top}) = \text{eval}_E(\text{Ker}(f)).$$

Indeed, if E is finite-dimensional, the map $\text{eval}_E: E \rightarrow E^{**}$ is an isomorphism, so every $\varphi \in E^{**}$ is of the form $\varphi = \text{eval}_E(u)$ for some $u \in E$, the map $\text{eval}_F: F \rightarrow F^{**}$ is injective, and we have

$$\begin{aligned} f^{\top\top}(\varphi) = 0 & \quad \text{iff} \quad f^{\top\top}(\text{eval}_E(u)) = 0 \\ & \quad \text{iff} \quad \text{eval}_F(f(u)) = 0 \\ & \quad \text{iff} \quad f(u) = 0 \\ & \quad \text{iff} \quad u \in \text{Ker}(f) \\ & \quad \text{iff} \quad \varphi \in \text{eval}_E(\text{Ker}(f)), \end{aligned}$$

which proves that $\text{Ker}(f^{\top\top}) = \text{eval}_E(\text{Ker}(f))$.

The following proposition shows the relationship between orthogonality and transposition.

Proposition 4.23. *Given a linear map $f: E \rightarrow F$, for any subspace V of E , we have*

$$f(V)^0 = (f^\top)^{-1}(V^0) = \{w^* \in F^* \mid f^\top(w^*) \in V^0\}.$$

As a consequence,

$$\text{Ker } f^\top = (\text{Im } f)^0 \quad \text{and} \quad \text{Ker } f = (\text{Im } f^\top)^0.$$

Proof. We have

$$\langle w^*, f(v) \rangle = \langle f^\top(w^*), v \rangle,$$

for all $v \in E$ and all $w^* \in F^*$, and thus, we have $\langle w^*, f(v) \rangle = 0$ for every $v \in V$, i.e. $w^* \in f(V)^0$, iff $\langle f^\top(w^*), v \rangle = 0$ for every $v \in V$, iff $f^\top(w^*) \in V^0$, i.e. $w^* \in (f^\top)^{-1}(V^0)$, proving that

$$f(V)^0 = (f^\top)^{-1}(V^0).$$

Since we already observed that $E^0 = 0$, letting $V = E$ in the above identity, we obtain that

$$\text{Ker } f^\top = (\text{Im } f)^0.$$

From the equation

$$\langle w^*, f(v) \rangle = \langle f^\top(w^*), v \rangle,$$

we deduce that $v \in (\text{Im } f^\top)^0$ iff $\langle f^\top(w^*), v \rangle = 0$ for all $w^* \in F^*$ iff $\langle w^*, f(v) \rangle = 0$ for all $w^* \in F^*$. Assume that $v \in (\text{Im } f^\top)^0$. If we pick a basis $(w_i)_{i \in I}$ of F , then we have the linear forms $w_i^*: F \rightarrow K$ such that $w_i^*(w_j) = \delta_{ij}$, and since we must have $\langle w_i^*, f(v) \rangle = 0$ for all $i \in I$ and $(w_i)_{i \in I}$ is a basis of F , we conclude that $f(v) = 0$, and thus $v \in \text{Ker } f$ (this is because $\langle w_i^*, f(v) \rangle$ is the coefficient of $f(v)$ associated with the basis vector w_i). Conversely, if $v \in \text{Ker } f$, then $\langle w^*, f(v) \rangle = 0$ for all $w^* \in F^*$, so we conclude that $v \in (\text{Im } f^\top)^0$. Therefore, $v \in (\text{Im } f^\top)^0$ iff $v \in \text{Ker } f$; that is,

$$\text{Ker } f = (\text{Im } f^\top)^0,$$

as claimed. □

The following proposition gives a natural interpretation of the dual $(E/U)^*$ of a quotient space E/U .

Proposition 4.24. *For any subspace U of a vector space E , if $p: E \rightarrow E/U$ is the canonical surjection onto E/U , then p^\top is injective and*

$$\text{Im}(p^\top) = U^0 = (\text{Ker}(p))^0.$$

Therefore, p^\top is a linear isomorphism between $(E/U)^*$ and U^0 .

Proof. Since p is surjective, by Proposition 4.21, the map p^\top is injective. Obviously, $U = \text{Ker}(p)$. Observe that $\text{Im}(p^\top)$ consists of all linear forms $\psi \in E^*$ such that $\psi = \varphi \circ p$ for some $\varphi \in (E/U)^*$, and since $\text{Ker}(p) = U$, we have $U \subseteq \text{Ker}(\psi)$. Conversely for any linear form $\psi \in E^*$, if $U \subseteq \text{Ker}(\psi)$, then ψ factors through E/U as $\psi = \bar{\psi} \circ p$ as shown in the following commutative diagram

$$\begin{array}{ccc} E & \xrightarrow{p} & E/U \\ & \searrow \psi & \downarrow \bar{\psi} \\ & & K, \end{array}$$

where $\bar{\psi}: E/U \rightarrow K$ is given by

$$\bar{\psi}(\bar{v}) = \psi(v), \quad v \in E,$$

where $\bar{v} \in E/U$ denotes the equivalence class of $v \in E$. The map $\bar{\psi}$ does not depend on the representative chosen in the equivalence class \bar{v} , since if $\bar{v}' = \bar{v}$, that is $v' - v = u \in U$, then $\psi(v') = \psi(v + u) = \psi(v) + \psi(u) = \psi(v) + 0 = \psi(v)$. Therefore, we have

$$\begin{aligned} \text{Im}(p^\top) &= \{\varphi \circ p \mid \varphi \in (E/U)^*\} \\ &= \{\psi: E \rightarrow K \mid U \subseteq \text{Ker}(\psi)\} \\ &= U^0, \end{aligned}$$

which proves our result. □

Proposition 4.24 yields another proof of part (b) of the duality theorem (theorem 4.17) that does not involve the existence of bases (in infinite dimension).

Proposition 4.25. *For any vector space E and any subspace V of E , we have $V^{00} = V$.*

Proof. We begin by observing that $V^0 = V^{000}$. This is because, for any subspace U of E^* , we have $U \subseteq U^{00}$, so $V^0 \subseteq V^{000}$. Furthermore, $V \subseteq V^{00}$ holds, and for any two subspaces M, N of E , if $M \subseteq N$ then $N^0 \subseteq M^0$, so we get $V^{000} \subseteq V^0$. Write $V_1 = V^{00}$, so that $V_1^0 = V^{000} = V^0$. We wish to prove that $V_1 = V$.

Since $V \subseteq V_1 = V^{00}$, the canonical projection $p_1: E \rightarrow E/V_1$ factors as $p_1 = f \circ p$ as in the diagram below,

$$\begin{array}{ccc} E & \xrightarrow{p} & E/V \\ & \searrow p_1 & \downarrow f \\ & & E/V_1 \end{array}$$

where $p: E \rightarrow E/V$ is the canonical projection onto E/V and $f: E/V \rightarrow E/V_1$ is the quotient map induced by p_1 , with $f(\bar{u}_{E/V}) = p_1(u) = \bar{u}_{E/V_1}$, for all $u \in E$ (since $V \subseteq V_1$, if $u - u' = v \in V$, then $u - u' = v \in V_1$, so $p_1(u) = p_1(u')$). Since p_1 is surjective, so is f . We wish to prove that f is actually an isomorphism, and for this, it is enough to show that f is injective. By transposing all the maps, we get the commutative diagram

$$\begin{array}{ccc} E^* & \xleftarrow{p^\top} & (E/V)^* \\ & \swarrow p_1^\top & \uparrow f^\top \\ & & (E/V_1)^* \end{array}$$

but by Proposition 4.24, the maps $p^\top: (E/V)^* \rightarrow V^0$ and $p_1^\top: (E/V_1)^* \rightarrow V_1^0$ are isomorphism, and since $V^0 = V_1^0$, we have the following diagram where both p^\top and p_1^\top are isomorphisms:

$$\begin{array}{ccc} V^0 & \xleftarrow{p^\top} & (E/V)^* \\ & \swarrow p_1^\top & \uparrow f^\top \\ & & (E/V_1)^* \end{array}$$

Therefore, $f^\top = (p^\top)^{-1} \circ p_1^\top$ is an isomorphism. We claim that this implies that f is injective.

If f is not injective, then there is some $x \in E/V$ such that $x \neq 0$ and $f(x) = 0$, so for every $\varphi \in (E/V_1)^*$, we have $f^\top(\varphi)(x) = \varphi(f(x)) = 0$. However, there is linear form $\psi \in (E/V)^*$ such that $\psi(x) = 1$, so $\psi \neq f^\top(\varphi)$ for all $\varphi \in (E/V_1)^*$, contradicting the fact that f^\top is surjective. To find such a linear form ψ , pick any supplement W of Kx in E/V , so that $E/V = Kx \oplus W$ (W is a hyperplane in E/V not containing x), and define ψ to be zero

on W and 1 on x .³ Therefore, f is injective, and since we already know that it is surjective, it is bijective. This means that the canonical map $f: E/V \rightarrow E/V_1$ with $V \subseteq V_1$ is an isomorphism, which implies that $V = V_1 = V^{00}$ (otherwise, if $v \in V_1 - V$, then $p_1(v) = 0$, so $f(p(v)) = p_1(v) = 0$, but $p(v) \neq 0$ since $v \notin V$, and f is not injective). \square

The following theorem shows the relationship between the rank of f and the rank of f^\top .

Theorem 4.26. *Given a linear map $f: E \rightarrow F$, the following properties hold.*

(a) *The dual $(\text{Im } f)^*$ of $\text{Im } f$ is isomorphic to $\text{Im } f^\top = f^\top(F^*)$; that is,*

$$(\text{Im } f)^* \approx \text{Im } f^\top.$$

(b) $\text{rk}(f) \leq \text{rk}(f^\top)$. *If $\text{rk}(f)$ is finite, we have $\text{rk}(f) = \text{rk}(f^\top)$.*

Proof. (a) Consider the linear maps

$$E \xrightarrow{p} \text{Im } f \xrightarrow{j} F,$$

where $E \xrightarrow{p} \text{Im } f$ is the surjective map induced by $E \xrightarrow{f} F$, and $\text{Im } f \xrightarrow{j} F$ is the injective inclusion map of $\text{Im } f$ into F . By definition, $f = j \circ p$. To simplify the notation, let $I = \text{Im } f$. By Proposition 4.21, since $E \xrightarrow{p} I$ is surjective, $I^* \xrightarrow{p^\top} E^*$ is injective, and since $\text{Im } f \xrightarrow{j} F$ is injective, $F^* \xrightarrow{j^\top} I^*$ is surjective. Since $f = j \circ p$, we also have

$$f^\top = (j \circ p)^\top = p^\top \circ j^\top,$$

and since $F^* \xrightarrow{j^\top} I^*$ is surjective, and $I^* \xrightarrow{p^\top} E^*$ is injective, we have an isomorphism between $(\text{Im } f)^*$ and $f^\top(F^*)$.

(b) We already noted that part (a) of Theorem 4.17 shows that $\dim(E) \leq \dim(E^*)$, for every vector space E . Thus, $\dim(\text{Im } f) \leq \dim((\text{Im } f)^*)$, which, by (a), shows that $\text{rk}(f) \leq \text{rk}(f^\top)$. When $\dim(\text{Im } f)$ is finite, we already observed that as a corollary of Theorem 4.17, $\dim(\text{Im } f) = \dim((\text{Im } f)^*)$, and thus, by part (a) we have $\text{rk}(f) = \text{rk}(f^\top)$.

If $\dim(F)$ is finite, then there is also a simple proof of (b) that doesn't use the result of part (a). By Theorem 4.17(c)

$$\dim(\text{Im } f) + \dim((\text{Im } f)^0) = \dim(F),$$

and by Theorem 4.11

$$\dim(\text{Ker } f^\top) + \dim(\text{Im } f^\top) = \dim(F^*).$$

³Using Zorn's lemma, we pick W maximal among all subspaces of E/V such that $Kx \cap W = (0)$; then, $E/V = Kx \oplus W$.

Furthermore, by Proposition 4.23, we have

$$\text{Ker } f^\top = (\text{Im } f)^0,$$

and since F is finite-dimensional $\dim(F) = \dim(F^*)$, so we deduce

$$\dim(\text{Im } f) + \dim((\text{Im } f)^0) = \dim((\text{Im } f)^0) + \dim(\text{Im } f^\top),$$

which yields $\dim(\text{Im } f) = \dim(\text{Im } f^\top)$; that is, $\text{rk}(f) = \text{rk}(f^\top)$. \square

Remarks:

1. If $\dim(E)$ is finite, following an argument of Dan Guralnik, we can also prove that $\text{rk}(f) = \text{rk}(f^\top)$ as follows.

We know from Proposition 4.23 applied to $f^\top : F^* \rightarrow E^*$ that

$$\text{Ker } (f^{\top\top}) = (\text{Im } f^\top)^0,$$

and we showed as a consequence of Proposition 4.22 that

$$\text{Ker } (f^{\top\top}) = \text{eval}_E(\text{Ker } (f)).$$

It follows (since eval_E is an isomorphism) that

$$\dim((\text{Im } f^\top)^0) = \dim(\text{Ker } (f^{\top\top})) = \dim(\text{Ker } (f)) = \dim(E) - \dim(\text{Im } f),$$

and since

$$\dim(\text{Im } f^\top) + \dim((\text{Im } f^\top)^0) = \dim(E),$$

we get

$$\dim(\text{Im } f^\top) = \dim(\text{Im } f).$$

2. As indicated by Dan Guralnik, if $\dim(E)$ is finite, the above result can be used to prove that

$$\text{Im } f^\top = (\text{Ker } (f))^0.$$

From

$$\langle f^\top(\varphi), u \rangle = \langle \varphi, f(u) \rangle$$

for all $\varphi \in F^*$ and all $u \in E$, we see that if $u \in \text{Ker } (f)$, then $\langle f^\top(\varphi), u \rangle = \langle \varphi, 0 \rangle = 0$, which means that $f^\top(\varphi) \in (\text{Ker } (f))^0$, and thus, $\text{Im } f^\top \subseteq (\text{Ker } (f))^0$. For the converse, since $\dim(E)$ is finite, we have

$$\dim((\text{Ker } (f))^0) = \dim(E) - \dim(\text{Ker } (f)) = \dim(\text{Im } f),$$

but we just proved that $\dim(\operatorname{Im} f^\top) = \dim(\operatorname{Im} f)$, so we get

$$\dim((\operatorname{Ker}(f))^0) = \dim(\operatorname{Im} f^\top),$$

and since $\operatorname{Im} f^\top \subseteq (\operatorname{Ker}(f))^0$, we obtain

$$\operatorname{Im} f^\top = (\operatorname{Ker}(f))^0,$$

as claimed. Now, since $(\operatorname{Ker}(f))^0 = \operatorname{Ker}(f)$, the above equation yields another proof of the fact that

$$\operatorname{Ker}(f) = (\operatorname{Im} f^\top)^0,$$

when E is finite-dimensional.

3. The equation

$$\operatorname{Im} f^\top = (\operatorname{Ker}(f))^0$$

is actually valid even if when E is infinite-dimensional, as we now prove.

Proposition 4.27. *If $f: E \rightarrow F$ is any linear map, then the following identities hold:*

$$\begin{aligned} \operatorname{Im} f^\top &= (\operatorname{Ker}(f))^0 \\ \operatorname{Ker}(f^\top) &= (\operatorname{Im} f)^0 \\ \operatorname{Im} f &= (\operatorname{Ker}(f^\top))^0 \\ \operatorname{Ker}(f) &= (\operatorname{Im} f^\top)^0. \end{aligned}$$

Proof. The equation $\operatorname{Ker}(f^\top) = (\operatorname{Im} f)^0$ has already been proved in Proposition 4.23.

By the duality theorem $(\operatorname{Ker}(f))^0 = \operatorname{Ker}(f)$, so from $\operatorname{Im} f^\top = (\operatorname{Ker}(f))^0$ we get $\operatorname{Ker}(f) = (\operatorname{Im} f^\top)^0$. Similarly, $(\operatorname{Im} f)^0 = \operatorname{Im} f$, so from $\operatorname{Ker}(f^\top) = (\operatorname{Im} f)^0$ we get $\operatorname{Im} f = (\operatorname{Ker}(f^\top))^0$. Therefore, what is left to be proved is that $\operatorname{Im} f^\top = (\operatorname{Ker}(f))^0$.

Let $p: E \rightarrow E/\operatorname{Ker}(f)$ be the canonical surjection, $\bar{f}: E/\operatorname{Ker}(f) \rightarrow \operatorname{Im} f$ be the isomorphism induced by f , and $j: \operatorname{Im} f \rightarrow F$ be the inclusion map. Then, we have

$$f = j \circ \bar{f} \circ p,$$

which implies that

$$f^\top = p^\top \circ \bar{f}^\top \circ j^\top.$$

Since p is surjective, p^\top is injective, since j is injective, j^\top is surjective, and since \bar{f} is bijective, \bar{f}^\top is also bijective. It follows that $(E/\operatorname{Ker}(f))^* = \operatorname{Im}(\bar{f}^\top \circ j^\top)$, and we have

$$\operatorname{Im} f^\top = \operatorname{Im} p^\top.$$

Since $p: E \rightarrow E/\operatorname{Ker}(f)$ is the canonical surjection, by Proposition 4.24 applied to $U = \operatorname{Ker}(f)$, we get

$$\operatorname{Im} f^\top = \operatorname{Im} p^\top = (\operatorname{Ker}(f))^0,$$

as claimed. □

In summary, the equation

$$\operatorname{Im} f^\top = (\operatorname{Ker}(f))^0$$

applies in any dimension, and it implies that

$$\operatorname{Ker}(f) = (\operatorname{Im} f^\top)^0.$$

The following proposition shows the relationship between the matrix representing a linear map $f: E \rightarrow F$ and the matrix representing its transpose $f^\top: F^* \rightarrow E^*$.

Proposition 4.28. *Let E and F be two vector spaces, and let (u_1, \dots, u_n) be a basis for E , and (v_1, \dots, v_m) be a basis for F . Given any linear map $f: E \rightarrow F$, if $M(f)$ is the $m \times n$ -matrix representing f w.r.t. the bases (u_1, \dots, u_n) and (v_1, \dots, v_m) , the $n \times m$ -matrix $M(f^\top)$ representing $f^\top: F^* \rightarrow E^*$ w.r.t. the dual bases (v_1^*, \dots, v_m^*) and (u_1^*, \dots, u_n^*) is the transpose $M(f)^\top$ of $M(f)$.*

Proof. Recall that the entry a_{ij} in row i and column j of $M(f)$ is the i -th coordinate of $f(u_j)$ over the basis (v_1, \dots, v_m) . By definition of v_i^* , we have $\langle v_i^*, f(u_j) \rangle = a_{ij}$. The entry a_{ji}^\top in row j and column i of $M(f^\top)$ is the j -th coordinate of

$$f^\top(v_i^*) = a_{1i}^\top u_1^* + \dots + a_{ji}^\top u_j^* + \dots + a_{ni}^\top u_n^*$$

over the basis (u_1^*, \dots, u_n^*) , which is just $a_{ji}^\top = f^\top(v_i^*)(u_j) = \langle f^\top(v_i^*), u_j \rangle$. Since

$$\langle v_i^*, f(u_j) \rangle = \langle f^\top(v_i^*), u_j \rangle,$$

we have $a_{ij} = a_{ji}^\top$, proving that $M(f^\top) = M(f)^\top$. □

We now can give a very short proof of the fact that the rank of a matrix is equal to the rank of its transpose.

Proposition 4.29. *Given a $m \times n$ matrix A over a field K , we have $\operatorname{rk}(A) = \operatorname{rk}(A^\top)$.*

Proof. The matrix A corresponds to a linear map $f: K^n \rightarrow K^m$, and by Theorem 4.26, $\operatorname{rk}(f) = \operatorname{rk}(f^\top)$. By Proposition 4.28, the linear map f^\top corresponds to A^\top . Since $\operatorname{rk}(A) = \operatorname{rk}(f)$, and $\operatorname{rk}(A^\top) = \operatorname{rk}(f^\top)$, we conclude that $\operatorname{rk}(A) = \operatorname{rk}(A^\top)$. □

Thus, given an $m \times n$ -matrix A , the maximum number of linearly independent columns is equal to the maximum number of linearly independent rows. There are other ways of proving this fact that do not involve the dual space, but instead some elementary transformations on rows and columns.

Proposition 4.29 immediately yields the following criterion for determining the rank of a matrix:

Proposition 4.30. *Given any $m \times n$ matrix A over a field K (typically $K = \mathbb{R}$ or $K = \mathbb{C}$), the rank of A is the maximum natural number r such that there is an invertible $r \times r$ submatrix of A obtained by selecting r rows and r columns of A .*

For example, the 3×2 matrix

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{pmatrix}$$

has rank 2 iff one of the three 2×2 matrices

$$\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \quad \begin{pmatrix} a_{11} & a_{12} \\ a_{31} & a_{32} \end{pmatrix} \quad \begin{pmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{pmatrix}$$

is invertible. We will see in Chapter 5 that this is equivalent to the fact the determinant of one of the above matrices is nonzero. This is not a very efficient way of finding the rank of a matrix. We will see that there are better ways using various decompositions such as LU, QR, or SVD.

4.5 The Four Fundamental Subspaces

Given a linear map $f: E \rightarrow F$ (where E and F are finite-dimensional), Proposition 4.23 revealed that the four spaces

$$\operatorname{Im} f, \operatorname{Im} f^\top, \operatorname{Ker} f, \operatorname{Ker} f^\top$$

play a special role. They are often called the *fundamental subspaces* associated with f . These spaces are related in an intimate manner, since Proposition 4.23 shows that

$$\begin{aligned} \operatorname{Ker} f &= (\operatorname{Im} f^\top)^\perp \\ \operatorname{Ker} f^\top &= (\operatorname{Im} f)^\perp, \end{aligned}$$

and Theorem 4.26 shows that

$$\operatorname{rk}(f) = \operatorname{rk}(f^\top).$$

It is instructive to translate these relations in terms of matrices (actually, certain linear algebra books make a big deal about this!). If $\dim(E) = n$ and $\dim(F) = m$, given any basis (u_1, \dots, u_n) of E and a basis (v_1, \dots, v_m) of F , we know that f is represented by an $m \times n$ matrix $A = (a_{ij})$, where the j th column of A is equal to $f(u_j)$ over the basis (v_1, \dots, v_m) . Furthermore, the transpose map f^\top is represented by the $n \times m$ matrix A^\top (with respect to the dual bases). Consequently, the four fundamental spaces

$$\operatorname{Im} f, \operatorname{Im} f^\top, \operatorname{Ker} f, \operatorname{Ker} f^\top$$

correspond to

- (1) The *column space* of A , denoted by $\text{Im } A$ or $\mathcal{R}(A)$; this is the subspace of \mathbb{R}^m spanned by the columns of A , which corresponds to the image $\text{Im } f$ of f .
- (2) The *kernel* or *nullspace* of A , denoted by $\text{Ker } A$ or $\mathcal{N}(A)$; this is the subspace of \mathbb{R}^n consisting of all vectors $x \in \mathbb{R}^n$ such that $Ax = 0$.
- (3) The *row space* of A , denoted by $\text{Im } A^\top$ or $\mathcal{R}(A^\top)$; this is the subspace of \mathbb{R}^n spanned by the rows of A , or equivalently, spanned by the columns of A^\top , which corresponds to the image $\text{Im } f^\top$ of f^\top .
- (4) The *left kernel* or *left nullspace* of A denoted by $\text{Ker } A^\top$ or $\mathcal{N}(A^\top)$; this is the kernel (nullspace) of A^\top , the subspace of \mathbb{R}^m consisting of all vectors $y \in \mathbb{R}^m$ such that $A^\top y = 0$, or equivalently, $y^\top A = 0$.

Recall that the dimension r of $\text{Im } f$, which is also equal to the dimension of the column space $\text{Im } A = \mathcal{R}(A)$, is the *rank* of A (and f). Then, some of our previous results can be reformulated as follows:

1. The column space $\mathcal{R}(A)$ of A has dimension r .
2. The nullspace $\mathcal{N}(A)$ of A has dimension $n - r$.
3. The row space $\mathcal{R}(A^\top)$ has dimension r .
4. The left nullspace $\mathcal{N}(A^\top)$ of A has dimension $m - r$.

The above statements constitute what Strang calls the *Fundamental Theorem of Linear Algebra, Part I* (see Strang [105]).

The two statements

$$\begin{aligned}\text{Ker } f &= (\text{Im } f^\top)^\perp \\ \text{Ker } f^\top &= (\text{Im } f)^\perp\end{aligned}$$

translate to

- (1) The nullspace of A is the orthogonal of the row space of A .
- (2) The left nullspace of A is the orthogonal of the column space of A .

The above statements constitute what Strang calls the *Fundamental Theorem of Linear Algebra, Part II* (see Strang [105]).

Since vectors are represented by column vectors and linear forms by row vectors (over a basis in E or F), a vector $x \in \mathbb{R}^n$ is orthogonal to a linear form y if

$$yx = 0.$$

Then, a vector $x \in \mathbb{R}^n$ is orthogonal to the row space of A iff x is orthogonal to every row of A , namely $Ax = 0$, which is equivalent to the fact that x belong to the nullspace of A . Similarly, the column vector $y \in \mathbb{R}^m$ (representing a linear form over the dual basis of F^*) belongs to the nullspace of A^\top iff $A^\top y = 0$, iff $y^\top A = 0$, which means that the linear form given by y^\top (over the basis in F) is orthogonal to the column space of A .

Since (2) is equivalent to the fact that the column space of A is equal to the orthogonal of the left nullspace of A , we get the following criterion for the solvability of an equation of the form $Ax = b$:

The equation $Ax = b$ has a solution iff for all $y \in \mathbb{R}^m$, if $A^\top y = 0$, then $y^\top b = 0$.

Indeed, the condition on the right-hand side says that b is orthogonal to the left nullspace of A , that is, that b belongs to the column space of A .

This criterion can be cheaper to check that checking directly that b is spanned by the columns of A . For example, if we consider the system

$$\begin{aligned}x_1 - x_2 &= b_1 \\x_2 - x_3 &= b_2 \\x_3 - x_1 &= b_3\end{aligned}$$

which, in matrix form, is written $Ax = b$ as below:

$$\begin{pmatrix} 1 & -1 & 0 \\ 0 & 1 & -1 \\ -1 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix},$$

we see that the rows of the matrix A add up to 0. In fact, it is easy to convince ourselves that the left nullspace of A is spanned by $y = (1, 1, 1)$, and so the system is solvable iff $y^\top b = 0$, namely

$$b_1 + b_2 + b_3 = 0.$$

Note that the above criterion can also be stated negatively as follows:

The equation $Ax = b$ has no solution iff there is some $y \in \mathbb{R}^m$ such that $A^\top y = 0$ and $y^\top b \neq 0$.

4.6 Summary

The main concepts and results of this chapter are listed below:

- *Direct products, sums, direct sums.*
- *Projections.*

- The fundamental equation

$$\dim(E) = \dim(\text{Ker } f) + \dim(\text{Im } f) = \dim(\text{Ker } f) + \text{rk}(f)$$

(Proposition 4.11).

- *Grassmann's relation*

$$\dim(U) + \dim(V) = \dim(U + V) + \dim(U \cap V).$$

- Characterizations of a bijective linear map $f: E \rightarrow F$.
- *Rank* of a matrix.
- The *dual space* E^* and *linear forms* (covector). The *bidual* E^{**} .
- The *bilinear pairing* $\langle -, - \rangle: E^* \times E \rightarrow K$ (the *canonical pairing*).
- *Evaluation at v* : $\text{eval}_v: E^* \rightarrow K$.
- The map $\text{eval}_E: E \rightarrow E^{**}$.
- *Orthogonality* between a subspace V of E and a subspace U of E^* ; the *orthogonal* V^0 and the *orthogonal* U^0 .
- *Coordinate forms*.
- The *Duality theorem* (Theorem 4.17).
- The *dual basis* of a basis.
- The isomorphism $\text{eval}_E: E \rightarrow E^{**}$ when $\dim(E)$ is finite.
- *Pairing* between two vector spaces; *nondegenerate pairing*; Proposition 4.19.
- Hyperplanes and linear forms.
- The *transpose* $f^\top: F^* \rightarrow E^*$ of a linear map $f: E \rightarrow F$.
- The fundamental identities:

$$\text{Ker } f^\top = (\text{Im } f)^0 \quad \text{and} \quad \text{Ker } f = (\text{Im } f^\top)^0$$

(Proposition 4.23).

- If F is finite-dimensional, then

$$\text{rk}(f) = \text{rk}(f^\top).$$

(Theorem 4.26).

- The matrix of the transpose map f^\top is equal to the transpose of the matrix of the map f (Proposition 4.28).
- For any $m \times n$ matrix A ,

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

- Characterization of the rank of a matrix in terms of a maximal invertible submatrix (Proposition 4.30).
- The *four fundamental subspaces*:

$$\text{Im } f, \text{Im } f^\top, \text{Ker } f, \text{Ker } f^\top.$$

- The *column space*, the *nullspace*, the *row space*, and the *left nullspace* (of a matrix).
- Criterion for the solvability of an equation of the form $Ax = b$ in terms of the left nullspace.

Chapter 5

Determinants

5.1 Permutations, Signature of a Permutation

This chapter contains a review of determinants and their use in linear algebra. We begin with permutations and the signature of a permutation. Next, we define multilinear maps and alternating multilinear maps. Determinants are introduced as alternating multilinear maps taking the value 1 on the unit matrix (following Emil Artin). It is then shown how to compute a determinant using the Laplace expansion formula, and the connection with the usual definition is made. It is shown how determinants can be used to invert matrices and to solve (at least in theory!) systems of linear equations (the Cramer formulae). The determinant of a linear map is defined. We conclude by defining the characteristic polynomial of a matrix (and of a linear map) and by proving the celebrated Cayley-Hamilton theorem which states that every matrix is a “zero” of its characteristic polynomial (we give two proofs; one computational, the other one more conceptual).

Determinants can be defined in several ways. For example, determinants can be defined in a fancy way in terms of the exterior algebra (or alternating algebra) of a vector space. We will follow a more algorithmic approach due to Emil Artin. No matter which approach is followed, we need a few preliminaries about permutations on a finite set. We need to show that every permutation on n elements is a product of transpositions, and that the parity of the number of transpositions involved is an invariant of the permutation. Let $[n] = \{1, 2, \dots, n\}$, where $n \in \mathbb{N}$, and $n > 0$.

Definition 5.1. A *permutation on n elements* is a bijection $\pi: [n] \rightarrow [n]$. When $n = 1$, the only function from $[1]$ to $[1]$ is the constant map: $1 \mapsto 1$. Thus, we will assume that $n \geq 2$. A *transposition* is a permutation $\tau: [n] \rightarrow [n]$ such that, for some $i < j$ (with $1 \leq i < j \leq n$), $\tau(i) = j$, $\tau(j) = i$, and $\tau(k) = k$, for all $k \in [n] - \{i, j\}$. In other words, a transposition exchanges two distinct elements $i, j \in [n]$. A *cyclic permutation of order k (or k -cycle)* is a permutation $\sigma: [n] \rightarrow [n]$ such that, for some sequence (i_1, i_2, \dots, i_k) of distinct elements of $[n]$ with $2 \leq k \leq n$,

$$\sigma(i_1) = i_2, \sigma(i_2) = i_3, \dots, \sigma(i_{k-1}) = i_k, \sigma(i_k) = i_1,$$

and $\sigma(j) = j$, for $j \in [n] - \{i_1, \dots, i_k\}$. The set $\{i_1, \dots, i_k\}$ is called the *domain* of the cyclic permutation, and the cyclic permutation is sometimes denoted by (i_1, i_2, \dots, i_k) .

If τ is a transposition, clearly, $\tau \circ \tau = \text{id}$. Also, a cyclic permutation of order 2 is a transposition, and for a cyclic permutation σ of order k , we have $\sigma^k = \text{id}$. Clearly, the composition of two permutations is a permutation and every permutation has an inverse which is also a permutation. Therefore, the set of permutations on $[n]$ is a *group* often denoted \mathfrak{S}_n . It is easy to show by induction that the group \mathfrak{S}_n has $n!$ elements. We will also use the terminology product of permutations (or transpositions), as a synonym for composition of permutations.

The following proposition shows the importance of cyclic permutations and transpositions.

Proposition 5.1. *For every $n \geq 2$, for every permutation $\pi: [n] \rightarrow [n]$, there is a partition of $[n]$ into r subsets called the orbits of π , with $1 \leq r \leq n$, where each set J in this partition is either a singleton $\{i\}$, or it is of the form*

$$J = \{i, \pi(i), \pi^2(i), \dots, \pi^{r_i-1}(i)\},$$

where r_i is the smallest integer, such that, $\pi^{r_i}(i) = i$ and $2 \leq r_i \leq n$. If π is not the identity, then it can be written in a unique way (up to the order) as a composition $\pi = \sigma_1 \circ \dots \circ \sigma_s$ of cyclic permutations with disjoint domains, where s is the number of orbits with at least two elements. Every permutation $\pi: [n] \rightarrow [n]$ can be written as a nonempty composition of transpositions.

Proof. Consider the relation R_π defined on $[n]$ as follows: $iR_\pi j$ iff there is some $k \geq 1$ such that $j = \pi^k(i)$. We claim that R_π is an equivalence relation. Transitivity is obvious. We claim that for every $i \in [n]$, there is some least r ($1 \leq r \leq n$) such that $\pi^r(i) = i$.

Indeed, consider the following sequence of $n + 1$ elements:

$$\langle i, \pi(i), \pi^2(i), \dots, \pi^n(i) \rangle.$$

Since $[n]$ only has n distinct elements, there are some h, k with $0 \leq h < k \leq n$ such that

$$\pi^h(i) = \pi^k(i),$$

and since π is a bijection, this implies $\pi^{k-h}(i) = i$, where $0 \leq k - h \leq n$. Thus, we proved that there is some integer $m \geq 1$ such that $\pi^m(i) = i$, so there is such a smallest integer r .

Consequently, R_π is reflexive. It is symmetric, since if $j = \pi^k(i)$, letting r be the least $r \geq 1$ such that $\pi^r(i) = i$, then

$$i = \pi^{kr}(i) = \pi^{k(r-1)}(\pi^k(i)) = \pi^{k(r-1)}(j).$$

Now, for every $i \in [n]$, the equivalence class (orbit) of i is a subset of $[n]$, either the singleton $\{i\}$ or a set of the form

$$J = \{i, \pi(i), \pi^2(i), \dots, \pi^{r_i-1}(i)\},$$

where r_i is the smallest integer such that $\pi^{r_i}(i) = i$ and $2 \leq r_i \leq n$, and in the second case, the restriction of π to J induces a cyclic permutation σ_i , and $\pi = \sigma_1 \circ \dots \circ \sigma_s$, where s is the number of equivalence classes having at least two elements.

For the second part of the proposition, we proceed by induction on n . If $n = 2$, there are exactly two permutations on $[2]$, the transposition τ exchanging 1 and 2, and the identity. However, $\text{id}_2 = \tau^2$. Now, let $n \geq 3$. If $\pi(n) = n$, since by the induction hypothesis, the restriction of π to $[n-1]$ can be written as a product of transpositions, π itself can be written as a product of transpositions. If $\pi(n) = k \neq n$, letting τ be the transposition such that $\tau(n) = k$ and $\tau(k) = n$, it is clear that $\tau \circ \pi$ leaves n invariant, and by the induction hypothesis, we have $\tau \circ \pi = \tau_m \circ \dots \circ \tau_1$ for some transpositions, and thus

$$\pi = \tau \circ \tau_m \circ \dots \circ \tau_1,$$

a product of transpositions (since $\tau \circ \tau = \text{id}_n$). □

Remark: When $\pi = \text{id}_n$ is the identity permutation, we can agree that the composition of 0 transpositions is the identity. The second part of Proposition 5.1 shows that the transpositions generate the group of permutations \mathfrak{S}_n .

In writing a permutation π as a composition $\pi = \sigma_1 \circ \dots \circ \sigma_s$ of cyclic permutations, it is clear that the order of the σ_i does not matter, since their domains are disjoint. Given a permutation written as a product of transpositions, we now show that the parity of the number of transpositions is an invariant.

Definition 5.2. For every $n \geq 2$, since every permutation $\pi: [n] \rightarrow [n]$ defines a partition of r subsets over which π acts either as the identity or as a cyclic permutation, let $\epsilon(\pi)$, called the *signature* of π , be defined by $\epsilon(\pi) = (-1)^{n-r}$, where r is the number of sets in the partition.

If τ is a transposition exchanging i and j , it is clear that the partition associated with τ consists of $n-1$ equivalence classes, the set $\{i, j\}$, and the $n-2$ singleton sets $\{k\}$, for $k \in [n] - \{i, j\}$, and thus, $\epsilon(\tau) = (-1)^{n-(n-1)} = (-1)^1 = -1$.

Proposition 5.2. For every $n \geq 2$, for every permutation $\pi: [n] \rightarrow [n]$, for every transposition τ , we have

$$\epsilon(\tau \circ \pi) = -\epsilon(\pi).$$

Consequently, for every product of transpositions such that $\pi = \tau_m \circ \dots \circ \tau_1$, we have

$$\epsilon(\pi) = (-1)^m,$$

which shows that the parity of the number of transpositions is an invariant.

Proof. Assume that $\tau(i) = j$ and $\tau(j) = i$, where $i < j$. There are two cases, depending whether i and j are in the same equivalence class J_l of R_π , or if they are in distinct equivalence classes. If i and j are in the same class J_l , then if

$$J_l = \{i_1, \dots, i_p, \dots, i_q, \dots, i_k\},$$

where $i_p = i$ and $i_q = j$, since

$$\tau(\pi(\pi^{-1}(i_p))) = \tau(i_p) = \tau(i) = j = i_q$$

and

$$\tau(\pi(i_{q-1})) = \tau(i_q) = \tau(j) = i = i_p,$$

it is clear that J_l splits into two subsets, one of which is $\{i_p, \dots, i_{q-1}\}$, and thus, the number of classes associated with $\tau \circ \pi$ is $r + 1$, and $\epsilon(\tau \circ \pi) = (-1)^{n-r-1} = -(-1)^{n-r} = -\epsilon(\pi)$. If i and j are in distinct equivalence classes J_l and J_m , say

$$\{i_1, \dots, i_p, \dots, i_h\}$$

and

$$\{j_1, \dots, j_q, \dots, j_k\},$$

where $i_p = i$ and $j_q = j$, since

$$\tau(\pi(\pi^{-1}(i_p))) = \tau(i_p) = \tau(i) = j = j_q$$

and

$$\tau(\pi(\pi^{-1}(j_q))) = \tau(j_q) = \tau(j) = i = i_p,$$

we see that the classes J_l and J_m merge into a single class, and thus, the number of classes associated with $\tau \circ \pi$ is $r - 1$, and $\epsilon(\tau \circ \pi) = (-1)^{n-r+1} = -(-1)^{n-r} = -\epsilon(\pi)$.

Now, let $\pi = \tau_m \circ \dots \circ \tau_1$ be any product of transpositions. By the first part of the proposition, we have

$$\epsilon(\pi) = (-1)^{m-1} \epsilon(\tau_1) = (-1)^{m-1} (-1) = (-1)^m,$$

since $\epsilon(\tau_1) = -1$ for a transposition. □

Remark: When $\pi = \text{id}_n$ is the identity permutation, since we agreed that the composition of 0 transpositions is the identity, it is still correct that $(-1)^0 = \epsilon(\text{id}) = +1$. From the proposition, it is immediate that $\epsilon(\pi' \circ \pi) = \epsilon(\pi')\epsilon(\pi)$. In particular, since $\pi^{-1} \circ \pi = \text{id}_n$, we get $\epsilon(\pi^{-1}) = \epsilon(\pi)$.

We can now proceed with the definition of determinants.

5.2 Alternating Multilinear Maps

First, we define multilinear maps, symmetric multilinear maps, and alternating multilinear maps.

Remark: Most of the definitions and results presented in this section also hold when K is a commutative ring, and when we consider modules over K (free modules, when bases are needed).

Let E_1, \dots, E_n , and F , be vector spaces over a field K , where $n \geq 1$.

Definition 5.3. A function $f: E_1 \times \dots \times E_n \rightarrow F$ is a *multilinear map* (or an *n-linear map*) if it is linear in each argument, holding the others fixed. More explicitly, for every i , $1 \leq i \leq n$, for all $x_1 \in E_1, \dots, x_{i-1} \in E_{i-1}, x_{i+1} \in E_{i+1}, \dots, x_n \in E_n$, for all $x, y \in E_i$, for all $\lambda \in K$,

$$\begin{aligned} f(x_1, \dots, x_{i-1}, x + y, x_{i+1}, \dots, x_n) &= f(x_1, \dots, x_{i-1}, x, x_{i+1}, \dots, x_n) \\ &\quad + f(x_1, \dots, x_{i-1}, y, x_{i+1}, \dots, x_n), \\ f(x_1, \dots, x_{i-1}, \lambda x, x_{i+1}, \dots, x_n) &= \lambda f(x_1, \dots, x_{i-1}, x, x_{i+1}, \dots, x_n). \end{aligned}$$

When $F = K$, we call f an *n-linear form* (or *multilinear form*). If $n \geq 2$ and $E_1 = E_2 = \dots = E_n$, an n -linear map $f: E \times \dots \times E \rightarrow F$ is called *symmetric*, if $f(x_1, \dots, x_n) = f(x_{\pi(1)}, \dots, x_{\pi(n)})$, for every permutation π on $\{1, \dots, n\}$. An n -linear map $f: E \times \dots \times E \rightarrow F$ is called *alternating*, if $f(x_1, \dots, x_n) = 0$ whenever $x_i = x_{i+1}$, for some i , $1 \leq i \leq n-1$ (in other words, when two adjacent arguments are equal). It does not harm to agree that when $n = 1$, a linear map is considered to be both symmetric and alternating, and we will do so.

When $n = 2$, a 2-linear map $f: E_1 \times E_2 \rightarrow F$ is called a *bilinear map*. We have already seen several examples of bilinear maps. Multiplication $\cdot: K \times K \rightarrow K$ is a bilinear map, treating K as a vector space over itself. More generally, multiplication $\cdot: A \times A \rightarrow A$ in a ring A is a bilinear map, viewing A as a module over itself.

The operation $\langle -, - \rangle: E^* \times E \rightarrow K$ applying a linear form to a vector is a bilinear map.

Symmetric bilinear maps (and multilinear maps) play an important role in geometry (inner products, quadratic forms), and in differential calculus (partial derivatives).

A bilinear map is symmetric if $f(u, v) = f(v, u)$, for all $u, v \in E$.

Alternating multilinear maps satisfy the following simple but crucial properties.

Proposition 5.3. *Let $f: E \times \dots \times E \rightarrow F$ be an n -linear alternating map, with $n \geq 2$. The following properties hold:*

(1)

$$f(\dots, x_i, x_{i+1}, \dots) = -f(\dots, x_{i+1}, x_i, \dots)$$

(2)

$$f(\dots, x_i, \dots, x_j, \dots) = 0,$$

where $x_i = x_j$, and $1 \leq i < j \leq n$.

(3)

$$f(\dots, x_i, \dots, x_j, \dots) = -f(\dots, x_j, \dots, x_i, \dots),$$

where $1 \leq i < j \leq n$.

(4)

$$f(\dots, x_i, \dots) = f(\dots, x_i + \lambda x_j, \dots),$$

for any $\lambda \in K$, and where $i \neq j$.

Proof. (1) By multilinearity applied twice, we have

$$\begin{aligned} f(\dots, x_i + x_{i+1}, x_i + x_{i+1}, \dots) &= f(\dots, x_i, x_i, \dots) + f(\dots, x_i, x_{i+1}, \dots) \\ &\quad + f(\dots, x_{i+1}, x_i, \dots) + f(\dots, x_{i+1}, x_{i+1}, \dots), \end{aligned}$$

and since f is alternating, this yields

$$0 = f(\dots, x_i, x_{i+1}, \dots) + f(\dots, x_{i+1}, x_i, \dots),$$

that is, $f(\dots, x_i, x_{i+1}, \dots) = -f(\dots, x_{i+1}, x_i, \dots)$.

(2) If $x_i = x_j$ and i and j are not adjacent, we can interchange x_i and x_{i+1} , and then x_i and x_{i+2} , etc, until x_i and x_j become adjacent. By (1),

$$f(\dots, x_i, \dots, x_j, \dots) = \epsilon f(\dots, x_i, x_j, \dots),$$

where $\epsilon = +1$ or -1 , but $f(\dots, x_i, x_j, \dots) = 0$, since $x_i = x_j$, and (2) holds.

(3) follows from (2) as in (1). (4) is an immediate consequence of (2). \square

Proposition 5.3 will now be used to show a fundamental property of alternating multilinear maps. First, we need to extend the matrix notation a little bit. Let E be a vector space over K . Given an $n \times n$ matrix $A = (a_{ij})$ over K , we can define a map $L(A): E^n \rightarrow E^n$ as follows:

$$L(A)_1(u) = a_{11}u_1 + \dots + a_{1n}u_n,$$

...

$$L(A)_n(u) = a_{n1}u_1 + \dots + a_{nn}u_n,$$

for all $u_1, \dots, u_n \in E$, with $u = (u_1, \dots, u_n)$. It is immediately verified that $L(A)$ is linear. Then, given two $n \times n$ matrices $A = (a_{ij})$ and $B = (b_{ij})$, by repeating the calculations establishing the product of matrices (just before Definition 3.1), we can show that

$$L(AB) = L(A) \circ L(B).$$

It is then convenient to use the matrix notation to describe the effect of the linear map $L(A)$, as

$$\begin{pmatrix} L(A)_1(u) \\ L(A)_2(u) \\ \vdots \\ L(A)_n(u) \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix}.$$

Lemma 5.4. *Let $f: E \times \cdots \times E \rightarrow F$ be an n -linear alternating map. Let (u_1, \dots, u_n) and (v_1, \dots, v_n) be two families of n vectors, such that,*

$$\begin{aligned} v_1 &= a_{11}u_1 + \cdots + a_{n1}u_n, \\ &\quad \dots \\ v_n &= a_{1n}u_1 + \cdots + a_{nn}u_n. \end{aligned}$$

Equivalently, letting

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix}$$

assume that we have

$$\begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix} = A^\top \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix}.$$

Then,

$$f(v_1, \dots, v_n) = \left(\sum_{\pi \in \mathfrak{S}_n} \epsilon(\pi) a_{\pi(1)1} \cdots a_{\pi(n)n} \right) f(u_1, \dots, u_n),$$

where the sum ranges over all permutations π on $\{1, \dots, n\}$.

Proof. Expanding $f(v_1, \dots, v_n)$ by multilinearity, we get a sum of terms of the form

$$a_{\pi(1)1} \cdots a_{\pi(n)n} f(u_{\pi(1)}, \dots, u_{\pi(n)}),$$

for all possible functions $\pi: \{1, \dots, n\} \rightarrow \{1, \dots, n\}$. However, because f is alternating, only the terms for which π is a permutation are nonzero. By Proposition 5.1, every permutation π is a product of transpositions, and by Proposition 5.2, the parity $\epsilon(\pi)$ of the number of transpositions only depends on π . Then, applying Proposition 5.3 (3) to each transposition in π , we get

$$a_{\pi(1)1} \cdots a_{\pi(n)n} f(u_{\pi(1)}, \dots, u_{\pi(n)}) = \epsilon(\pi) a_{\pi(1)1} \cdots a_{\pi(n)n} f(u_1, \dots, u_n).$$

Thus, we get the expression of the lemma. □

The quantity

$$\det(A) = \sum_{\pi \in \mathfrak{S}_n} \epsilon(\pi) a_{\pi(1)1} \cdots a_{\pi(n)n}$$

is in fact the value of the determinant of A (which, as we shall see shortly, is also equal to the determinant of A^\top). However, working directly with the above definition is quite awkward, and we will proceed via a slightly indirect route

5.3 Definition of a Determinant

Recall that the set of all square $n \times n$ -matrices with coefficients in a field K is denoted by $M_n(K)$.

Definition 5.4. A determinant is defined as any map

$$D: M_n(K) \rightarrow K,$$

which, when viewed as a map on $(K^n)^n$, i.e., a map of the n columns of a matrix, is n -linear alternating and such that $D(I_n) = 1$ for the identity matrix I_n . Equivalently, we can consider a vector space E of dimension n , some fixed basis (e_1, \dots, e_n) , and define

$$D: E^n \rightarrow K$$

as an n -linear alternating map such that $D(e_1, \dots, e_n) = 1$.

First, we will show that such maps D exist, using an inductive definition that also gives a recursive method for computing determinants. Actually, we will define a family $(\mathcal{D}_n)_{n \geq 1}$ of (finite) sets of maps $D: M_n(K) \rightarrow K$. Second, we will show that determinants are in fact uniquely defined, that is, we will show that each \mathcal{D}_n consists of a single map. This will show the equivalence of the direct definition $\det(A)$ of Lemma 5.4 with the inductive definition $D(A)$. Finally, we will prove some basic properties of determinants, using the uniqueness theorem.

Given a matrix $A \in M_n(K)$, we denote its n columns by A^1, \dots, A^n .

Definition 5.5. For every $n \geq 1$, we define a finite set \mathcal{D}_n of maps $D: M_n(K) \rightarrow K$ inductively as follows:

When $n = 1$, \mathcal{D}_1 consists of the single map D such that, $D(A) = a$, where $A = (a)$, with $a \in K$.

Assume that \mathcal{D}_{n-1} has been defined, where $n \geq 2$. We define the set \mathcal{D}_n as follows. For every matrix $A \in M_n(K)$, let A_{ij} be the $(n-1) \times (n-1)$ -matrix obtained from $A = (a_{ij})$ by deleting row i and column j . Then, \mathcal{D}_n consists of all the maps D such that, for some i , $1 \leq i \leq n$,

$$D(A) = (-1)^{i+1} a_{i1} D(A_{i1}) + \cdots + (-1)^{i+n} a_{in} D(A_{in}),$$

where for every j , $1 \leq j \leq n$, $D(A_{ij})$ is the result of applying any D in \mathcal{D}_{n-1} to A_{ij} .



We confess that the use of the same letter D for the member of \mathcal{D}_n being defined, and for members of \mathcal{D}_{n-1} , may be slightly confusing. We considered using subscripts to distinguish, but this seems to complicate things unnecessarily. One should not worry too much anyway, since it will turn out that each \mathcal{D}_n contains just one map.

Each $(-1)^{i+j}D(A_{ij})$ is called the *cofactor* of a_{ij} , and the inductive expression for $D(A)$ is called a *Laplace expansion of D according to the i -th row*. Given a matrix $A \in M_n(K)$, each $D(A)$ is called a *determinant* of A .

We can think of each member of \mathcal{D}_n as an *algorithm* to evaluate “the” determinant of A . The main point is that these algorithms, which recursively evaluate a determinant using all possible Laplace row expansions, all yield the same result, $\det(A)$.

We will prove shortly that $D(A)$ is uniquely defined (at the moment, it is not clear that \mathcal{D}_n consists of a single map). Assuming this fact, given a $n \times n$ -matrix $A = (a_{ij})$,

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix}$$

its determinant is denoted by $D(A)$ or $\det(A)$, or more explicitly by

$$\det(A) = \begin{vmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{vmatrix}$$

First, let us first consider some examples.

Example 5.1.

1. When $n = 2$, if

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

expanding according to any row, we have

$$D(A) = ad - bc.$$

2. When $n = 3$, if

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}$$

expanding according to the first row, we have

$$D(A) = a_{11} \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix} - a_{12} \begin{vmatrix} a_{21} & a_{23} \\ a_{31} & a_{33} \end{vmatrix} + a_{13} \begin{vmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{vmatrix}$$

that is,

$$D(A) = a_{11}(a_{22}a_{33} - a_{32}a_{23}) - a_{12}(a_{21}a_{33} - a_{31}a_{23}) + a_{13}(a_{21}a_{32} - a_{31}a_{22}),$$

which gives the explicit formula

$$D(A) = a_{11}a_{22}a_{33} + a_{21}a_{32}a_{13} + a_{31}a_{12}a_{23} - a_{11}a_{32}a_{23} - a_{21}a_{12}a_{33} - a_{31}a_{22}a_{13}.$$

We now show that each $D \in \mathcal{D}_n$ is a determinant (map).

Lemma 5.5. *For every $n \geq 1$, for every $D \in \mathcal{D}_n$ as defined in Definition 5.5, D is an alternating multilinear map such that $D(I_n) = 1$.*

Proof. By induction on n , it is obvious that $D(I_n) = 1$. Let us now prove that D is multilinear. Let us show that D is linear in each column. Consider any column k . Since

$$D(A) = (-1)^{i+1}a_{i1}D(A_{i1}) + \cdots + (-1)^{i+j}a_{ij}D(A_{ij}) + \cdots + (-1)^{i+n}a_{in}D(A_{in}),$$

if $j \neq k$, then by induction, $D(A_{ij})$ is linear in column k , and a_{ij} does not belong to column k , so $(-1)^{i+j}a_{ij}D(A_{ij})$ is linear in column k . If $j = k$, then $D(A_{ij})$ does not depend on column $k = j$, since A_{ij} is obtained from A by deleting row i and column $j = k$, and a_{ij} belongs to column $j = k$. Thus, $(-1)^{i+j}a_{ij}D(A_{ij})$ is linear in column k . Consequently, in all cases, $(-1)^{i+j}a_{ij}D(A_{ij})$ is linear in column k , and thus, $D(A)$ is linear in column k .

Let us now prove that D is alternating. Assume that two adjacent rows of A are equal, say $A^k = A^{k+1}$. First, let $j \neq k$ and $j \neq k+1$. Then, the matrix A_{ij} has two identical adjacent columns, and by the induction hypothesis, $D(A_{ij}) = 0$. The remaining terms of $D(A)$ are

$$(-1)^{i+k}a_{ik}D(A_{ik}) + (-1)^{i+k+1}a_{i,k+1}D(A_{i,k+1}).$$

However, the two matrices A_{ik} and $A_{i,k+1}$ are equal, since we are assuming that columns k and $k+1$ of A are identical, and since A_{ik} is obtained from A by deleting row i and column k , and $A_{i,k+1}$ is obtained from A by deleting row i and column $k+1$. Similarly, $a_{ik} = a_{i,k+1}$, since columns k and $k+1$ of A are equal. But then,

$$(-1)^{i+k}a_{ik}D(A_{ik}) + (-1)^{i+k+1}a_{i,k+1}D(A_{i,k+1}) = (-1)^{i+k}a_{ik}D(A_{ik}) - (-1)^{i+k}a_{ik}D(A_{ik}) = 0.$$

This shows that D is alternating, and completes the proof. \square

Lemma 5.5 shows the existence of determinants. We now prove their uniqueness.

Theorem 5.6. *For every $n \geq 1$, for every $D \in \mathcal{D}_n$, for every matrix $A \in M_n(K)$, we have*

$$D(A) = \sum_{\pi \in \mathfrak{S}_n} \epsilon(\pi) a_{\pi(1)1} \cdots a_{\pi(n)n},$$

where the sum ranges over all permutations π on $\{1, \dots, n\}$. As a consequence, \mathcal{D}_n consists of a single map for every $n \geq 1$, and this map is given by the above explicit formula.

Proof. Consider the standard basis (e_1, \dots, e_n) of K^n , where $(e_i)_i = 1$ and $(e_i)_j = 0$, for $j \neq i$. Then, each column A^j of A corresponds to a vector v_j whose coordinates over the basis (e_1, \dots, e_n) are the components of A^j , that is, we can write

$$\begin{aligned} v_1 &= a_{11}e_1 + \cdots + a_{n1}e_n, \\ &\quad \dots \\ v_n &= a_{1n}e_1 + \cdots + a_{nn}e_n. \end{aligned}$$

Since by Lemma 5.5, each D is a multilinear alternating map, by applying Lemma 5.4, we get

$$D(A) = D(v_1, \dots, v_n) = \left(\sum_{\pi \in \mathfrak{S}_n} \epsilon(\pi) a_{\pi(1)1} \cdots a_{\pi(n)n} \right) D(e_1, \dots, e_n),$$

where the sum ranges over all permutations π on $\{1, \dots, n\}$. But $D(e_1, \dots, e_n) = D(I_n)$, and by Lemma 5.5, we have $D(I_n) = 1$. Thus,

$$D(A) = \sum_{\pi \in \mathfrak{S}_n} \epsilon(\pi) a_{\pi(1)1} \cdots a_{\pi(n)n},$$

where the sum ranges over all permutations π on $\{1, \dots, n\}$. □

From now on, we will favor the notation $\det(A)$ over $D(A)$ for the determinant of a square matrix.

Remark: There is a geometric interpretation of determinants which we find quite illuminating. Given n linearly independent vectors (u_1, \dots, u_n) in \mathbb{R}^n , the set

$$P_n = \{\lambda_1 u_1 + \cdots + \lambda_n u_n \mid 0 \leq \lambda_i \leq 1, 1 \leq i \leq n\}$$

is called a *parallelotope*. If $n = 2$, then P_2 is a *parallelogram* and if $n = 3$, then P_3 is a *parallelepiped*, a skew box having u_1, u_2, u_3 as three of its corner sides. Then, it turns out that $\det(u_1, \dots, u_n)$ is the *signed volume* of the parallelotope P_n (where volume means n -dimensional volume). The sign of this volume accounts for the orientation of P_n in \mathbb{R}^n .

We can now prove some properties of determinants.

Corollary 5.7. *For every matrix $A \in M_n(K)$, we have $\det(A) = \det(A^\top)$.*

Proof. By Theorem 5.6, we have

$$\det(A) = \sum_{\pi \in \mathfrak{S}_n} \epsilon(\pi) a_{\pi(1)1} \cdots a_{\pi(n)n},$$

where the sum ranges over all permutations π on $\{1, \dots, n\}$. Since a permutation is invertible, every product

$$a_{\pi(1)1} \cdots a_{\pi(n)n}$$

can be rewritten as

$$a_{1\pi^{-1}(1)} \cdots a_{n\pi^{-1}(n)},$$

and since $\epsilon(\pi^{-1}) = \epsilon(\pi)$ and the sum is taken over all permutations on $\{1, \dots, n\}$, we have

$$\sum_{\pi \in \mathfrak{S}_n} \epsilon(\pi) a_{\pi(1)1} \cdots a_{\pi(n)n} = \sum_{\sigma \in \mathfrak{S}_n} \epsilon(\sigma) a_{1\sigma(1)} \cdots a_{n\sigma(n)},$$

where π and σ range over all permutations. But it is immediately verified that

$$\det(A^\top) = \sum_{\sigma \in \mathfrak{S}_n} \epsilon(\sigma) a_{1\sigma(1)} \cdots a_{n\sigma(n)}.$$

□

A useful consequence of Corollary 5.7 is that the determinant of a matrix is also a multilinear alternating map of its rows. This fact, combined with the fact that the determinant of a matrix is a multilinear alternating map of its columns is often useful for finding short-cuts in computing determinants. We illustrate this point on the following example which shows up in polynomial interpolation.

Example 5.2. Consider the so-called *Vandermonde determinant*

$$V(x_1, \dots, x_n) = \begin{vmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \\ x_1^2 & x_2^2 & \cdots & x_n^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{n-1} & x_2^{n-1} & \cdots & x_n^{n-1} \end{vmatrix}.$$

We claim that

$$V(x_1, \dots, x_n) = \prod_{1 \leq i < j \leq n} (x_j - x_i),$$

with $V(x_1, \dots, x_n) = 1$, when $n = 1$. We prove it by induction on $n \geq 1$. The case $n = 1$ is obvious. Assume $n \geq 2$. We proceed as follows: multiply row $n - 1$ by x_1 and subtract it from row n (the last row), then multiply row $n - 2$ by x_1 and subtract it from row $n - 1$,

etc, multiply row $i - 1$ by x_1 and subtract it from row i , until we reach row 1. We obtain the following determinant:

$$V(x_1, \dots, x_n) = \begin{vmatrix} 1 & 1 & \dots & 1 \\ 0 & x_2 - x_1 & \dots & x_n - x_1 \\ 0 & x_2(x_2 - x_1) & \dots & x_n(x_n - x_1) \\ \vdots & \vdots & \ddots & \vdots \\ 0 & x_2^{n-2}(x_2 - x_1) & \dots & x_n^{n-2}(x_n - x_1) \end{vmatrix}$$

Now, expanding this determinant according to the first column and using multilinearity, we can factor $(x_i - x_1)$ from the column of index $i - 1$ of the matrix obtained by deleting the first row and the first column, and thus

$$V(x_1, \dots, x_n) = (x_2 - x_1)(x_3 - x_1) \cdots (x_n - x_1)V(x_2, \dots, x_n),$$

which establishes the induction step.

Lemma 5.4 can be reformulated nicely as follows.

Proposition 5.8. *Let $f: E \times \dots \times E \rightarrow F$ be an n -linear alternating map. Let (u_1, \dots, u_n) and (v_1, \dots, v_n) be two families of n vectors, such that*

$$\begin{aligned} v_1 &= a_{11}u_1 + \dots + a_{1n}u_n, \\ &\dots \\ v_n &= a_{n1}u_1 + \dots + a_{nn}u_n. \end{aligned}$$

Equivalently, letting

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix}$$

assume that we have

$$\begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix} = A \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix}.$$

Then,

$$f(v_1, \dots, v_n) = \det(A)f(u_1, \dots, u_n).$$

Proof. The only difference with Lemma 5.4 is that here, we are using A^\top instead of A . Thus, by Lemma 5.4 and Corollary 5.7, we get the desired result. \square

As a consequence, we get the very useful property that the determinant of a product of matrices is the product of the determinants of these matrices.

Proposition 5.9. *For any two $n \times n$ -matrices A and B , we have $\det(AB) = \det(A) \det(B)$.*

Proof. We use Proposition 5.8 as follows: let (e_1, \dots, e_n) be the standard basis of K^n , and let

$$\begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix} = AB \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}.$$

Then, we get

$$\det(w_1, \dots, w_n) = \det(AB) \det(e_1, \dots, e_n) = \det(AB),$$

since $\det(e_1, \dots, e_n) = 1$. Now, letting

$$\begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix} = B \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix},$$

we get

$$\det(v_1, \dots, v_n) = \det(B),$$

and since

$$\begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix} = A \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix},$$

we get

$$\det(w_1, \dots, w_n) = \det(A) \det(v_1, \dots, v_n) = \det(A) \det(B).$$

□

It should be noted that all the results of this section, up to now, also holds when K is a commutative ring, and not necessarily a field. We can now characterize when an $n \times n$ -matrix A is invertible in terms of its determinant $\det(A)$.

5.4 Inverse Matrices and Determinants

In the next two sections, K is a commutative ring and when needed, a field.

Definition 5.6. Let K be a commutative ring. Given a matrix $A \in M_n(K)$, let $\tilde{A} = (b_{ij})$ be the matrix defined such that

$$b_{ij} = (-1)^{i+j} \det(A_{ji}),$$

the cofactor of a_{ji} . The matrix \tilde{A} is called the *adjugate* of A , and each matrix A_{ji} is called a *minor* of the matrix A .



Note the reversal of the indices in

$$b_{ij} = (-1)^{i+j} \det(A_{ji}).$$

Thus, \tilde{A} is the transpose of the matrix of cofactors of elements of A .

We have the following proposition.

Proposition 5.10. *Let K be a commutative ring. For every matrix $A \in M_n(K)$, we have*

$$A\tilde{A} = \tilde{A}A = \det(A)I_n.$$

As a consequence, A is invertible iff $\det(A)$ is invertible, and if so, $A^{-1} = (\det(A))^{-1}\tilde{A}$.

Proof. If $\tilde{A} = (b_{ij})$ and $A\tilde{A} = (c_{ij})$, we know that the entry c_{ij} in row i and column j of $A\tilde{A}$ is

$$c_{ij} = a_{i1}b_{1j} + \cdots + a_{ik}b_{kj} + \cdots + a_{in}b_{nj},$$

which is equal to

$$a_{i1}(-1)^{j+1} \det(A_{j1}) + \cdots + a_{in}(-1)^{j+n} \det(A_{jn}).$$

If $j = i$, then we recognize the expression of the expansion of $\det(A)$ according to the i -th row:

$$c_{ii} = \det(A) = a_{i1}(-1)^{i+1} \det(A_{i1}) + \cdots + a_{in}(-1)^{i+n} \det(A_{in}).$$

If $j \neq i$, we can form the matrix A' by replacing the j -th row of A by the i -th row of A . Now, the matrix A_{jk} obtained by deleting row j and column k from A is equal to the matrix A'_{jk} obtained by deleting row j and column k from A' , since A and A' only differ by the j -th row. Thus,

$$\det(A_{jk}) = \det(A'_{jk}),$$

and we have

$$c_{ij} = a_{i1}(-1)^{j+1} \det(A'_{j1}) + \cdots + a_{in}(-1)^{j+n} \det(A'_{jn}).$$

However, this is the expansion of $\det(A')$ according to the j -th row, since the j -th row of A' is equal to the i -th row of A , and since A' has two identical rows i and j , because \det is an alternating map of the rows (see an earlier remark), we have $\det(A') = 0$. Thus, we have shown that $c_{ii} = \det(A)$, and $c_{ij} = 0$, when $j \neq i$, and so

$$A\tilde{A} = \det(A)I_n.$$

It is also obvious from the definition of \tilde{A} , that

$$\tilde{A}^\top = \widetilde{A^\top}.$$

Then, applying the first part of the argument to A^\top , we have

$$A^\top \widetilde{A^\top} = \det(A^\top) I_n,$$

and since, $\det(A^\top) = \det(A)$, $\tilde{A}^\top = \widetilde{A^\top}$, and $(\tilde{A}A)^\top = A^\top \tilde{A}^\top$, we get

$$\det(A) I_n = A^\top \widetilde{A^\top} = A^\top \tilde{A}^\top = (\tilde{A}A)^\top,$$

that is,

$$(\tilde{A}A)^\top = \det(A) I_n,$$

which yields

$$\tilde{A}A = \det(A) I_n,$$

since $I_n^\top = I_n$. This proves that

$$A\tilde{A} = \tilde{A}A = \det(A) I_n.$$

As a consequence, if $\det(A)$ is invertible, we have $A^{-1} = (\det(A))^{-1} \tilde{A}$. Conversely, if A is invertible, from $AA^{-1} = I_n$, by Proposition 5.9, we have $\det(A) \det(A^{-1}) = 1$, and $\det(A)$ is invertible. \square

When K is a field, an element $a \in K$ is invertible iff $a \neq 0$. In this case, the second part of the proposition can be stated as A is invertible iff $\det(A) \neq 0$. Note in passing that this method of computing the inverse of a matrix is usually not practical.

We now consider some applications of determinants to linear independence and to solving systems of linear equations. Although these results hold for matrices over an integral domain, their proofs require more sophisticated methods (it is necessary to use the fraction field of the integral domain, K). Therefore, we assume again that K is a field.

Let A be an $n \times n$ -matrix, x a column vectors of variables, and b another column vector, and let A^1, \dots, A^n denote the columns of A . Observe that the system of equation $Ax = b$,

$$\begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}$$

is equivalent to

$$x_1 A^1 + \cdots + x_j A^j + \cdots + x_n A^n = b,$$

since the equation corresponding to the i -th row is in both cases

$$a_{i1}x_1 + \cdots + a_{ij}x_j + \cdots + a_{in}x_n = b_i.$$

First, we characterize linear independence of the column vectors of a matrix A in terms of its determinant.

Proposition 5.11. *Given an $n \times n$ -matrix A over a field K , the columns A^1, \dots, A^n of A are linearly dependent iff $\det(A) = \det(A^1, \dots, A^n) = 0$. Equivalently, A has rank n iff $\det(A) \neq 0$.*

Proof. First, assume that the columns A^1, \dots, A^n of A are linearly dependent. Then, there are $x_1, \dots, x_n \in K$, such that

$$x_1 A^1 + \cdots + x_j A^j + \cdots + x_n A^n = 0,$$

where $x_j \neq 0$ for some j . If we compute

$$\det(A^1, \dots, x_1 A^1 + \cdots + x_j A^j + \cdots + x_n A^n, \dots, A^n) = \det(A^1, \dots, 0, \dots, A^n) = 0,$$

where 0 occurs in the j -th position, by multilinearity, all terms containing two identical columns A^k for $k \neq j$ vanish, and we get

$$x_j \det(A^1, \dots, A^n) = 0.$$

Since $x_j \neq 0$ and K is a field, we must have $\det(A^1, \dots, A^n) = 0$.

Conversely, we show that if the columns A^1, \dots, A^n of A are linearly independent, then $\det(A^1, \dots, A^n) \neq 0$. If the columns A^1, \dots, A^n of A are linearly independent, then they form a basis of K^n , and we can express the standard basis (e_1, \dots, e_n) of K^n in terms of A^1, \dots, A^n . Thus, we have

$$\begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix} = \begin{pmatrix} b_{11} & b_{12} & \cdots & b_{1n} \\ b_{21} & b_{22} & \cdots & b_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ b_{n1} & b_{n2} & \cdots & b_{nn} \end{pmatrix} \begin{pmatrix} A^1 \\ A^2 \\ \vdots \\ A^n \end{pmatrix},$$

for some matrix $B = (b_{ij})$, and by Proposition 5.8, we get

$$\det(e_1, \dots, e_n) = \det(B) \det(A^1, \dots, A^n),$$

and since $\det(e_1, \dots, e_n) = 1$, this implies that $\det(A^1, \dots, A^n) \neq 0$ (and $\det(B) \neq 0$). For the second assertion, recall that the rank of a matrix is equal to the maximum number of linearly independent columns, and the conclusion is clear. \square

If we combine Proposition 5.11 with Proposition 4.30, we obtain the following criterion for finding the rank of a matrix.

Proposition 5.12. *Given any $m \times n$ matrix A over a field K (typically $K = \mathbb{R}$ or $K = \mathbb{C}$), the rank of A is the maximum natural number r such that there is an $r \times r$ submatrix B of A obtained by selecting r rows and r columns of A , and such that $\det(B) \neq 0$.*

5.5 Systems of Linear Equations and Determinants

We now characterize when a system of linear equations of the form $Ax = b$ has a unique solution.

Proposition 5.13. *Given an $n \times n$ -matrix A over a field K , the following properties hold:*

- (1) *For every column vector b , there is a unique column vector x such that $Ax = b$ iff the only solution to $Ax = 0$ is the trivial vector $x = 0$, iff $\det(A) \neq 0$.*
- (2) *If $\det(A) \neq 0$, the unique solution of $Ax = b$ is given by the expressions*

$$x_j = \frac{\det(A^1, \dots, A^{j-1}, b, A^{j+1}, \dots, A^n)}{\det(A^1, \dots, A^{j-1}, A^j, A^{j+1}, \dots, A^n)},$$

known as Cramer's rules.

- (3) *The system of linear equations $Ax = 0$ has a nonzero solution iff $\det(A) = 0$.*

Proof. Assume that $Ax = b$ has a single solution x_0 , and assume that $Ay = 0$ with $y \neq 0$. Then,

$$A(x_0 + y) = Ax_0 + Ay = Ax_0 + 0 = b,$$

and $x_0 + y \neq x_0$ is another solution of $Ax = b$, contradicting the hypothesis that $Ax = b$ has a single solution x_0 . Thus, $Ax = 0$ only has the trivial solution. Now, assume that $Ax = 0$ only has the trivial solution. This means that the columns A^1, \dots, A^n of A are linearly independent, and by Proposition 5.11, we have $\det(A) \neq 0$. Finally, if $\det(A) \neq 0$, by Proposition 5.10, this means that A is invertible, and then, for every b , $Ax = b$ is equivalent to $x = A^{-1}b$, which shows that $Ax = b$ has a single solution.

- (2) Assume that $Ax = b$. If we compute

$$\det(A^1, \dots, x_1 A^1 + \dots + x_j A^j + \dots + x_n A^n, \dots, A^n) = \det(A^1, \dots, b, \dots, A^n),$$

where b occurs in the j -th position, by multilinearity, all terms containing two identical columns A^k for $k \neq j$ vanish, and we get

$$x_j \det(A^1, \dots, A^n) = \det(A^1, \dots, A^{j-1}, b, A^{j+1}, \dots, A^n),$$

for every j , $1 \leq j \leq n$. Since we assumed that $\det(A) = \det(A^1, \dots, A^n) \neq 0$, we get the desired expression.

- (3) Note that $Ax = 0$ has a nonzero solution iff A^1, \dots, A^n are linearly dependent (as observed in the proof of Proposition 5.11), which, by Proposition 5.11, is equivalent to $\det(A) = 0$. □

As pleasing as Cramer's rules are, it is usually impractical to solve systems of linear equations using the above expressions.

5.6 Determinant of a Linear Map

We close this chapter with the notion of determinant of a linear map $f: E \rightarrow E$.

Given a vector space E of finite dimension n , given a basis (u_1, \dots, u_n) of E , for every linear map $f: E \rightarrow E$, if $M(f)$ is the matrix of f w.r.t. the basis (u_1, \dots, u_n) , we can define $\det(f) = \det(M(f))$. If (v_1, \dots, v_n) is any other basis of E , and if P is the change of basis matrix, by Corollary 3.5, the matrix of f with respect to the basis (v_1, \dots, v_n) is $P^{-1}M(f)P$. Now, by proposition 5.9, we have

$$\det(P^{-1}M(f)P) = \det(P^{-1})\det(M(f))\det(P) = \det(P^{-1})\det(P)\det(M(f)) = \det(M(f)).$$

Thus, $\det(f)$ is indeed independent of the basis of E .

Definition 5.7. Given a vector space E of finite dimension, for any linear map $f: E \rightarrow E$, we define the *determinant* $\det(f)$ of f as the determinant $\det(M(f))$ of the matrix of f in any basis (since, from the discussion just before this definition, this determinant does not depend on the basis).

Then, we have the following proposition.

Proposition 5.14. *Given any vector space E of finite dimension n , a linear map $f: E \rightarrow E$ is invertible iff $\det(f) \neq 0$.*

Proof. The linear map $f: E \rightarrow E$ is invertible iff its matrix $M(f)$ in any basis is invertible (by Proposition 3.2), iff $\det(M(f)) \neq 0$, by Proposition 5.10. \square

Given a vector space of finite dimension n , it is easily seen that the set of bijective linear maps $f: E \rightarrow E$ such that $\det(f) = 1$ is a group under composition. This group is a subgroup of the general linear group $\mathbf{GL}(E)$. It is called the *special linear group (of E)*, and it is denoted by $\mathbf{SL}(E)$, or when $E = K^n$, by $\mathbf{SL}(n, K)$, or even by $\mathbf{SL}(n)$.

5.7 The Cayley–Hamilton Theorem

We conclude this chapter with an interesting and important application of Proposition 5.10, the *Cayley–Hamilton theorem*. The results of this section apply to matrices over any commutative ring K . First, we need the concept of the characteristic polynomial of a matrix.

Definition 5.8. If K is any commutative ring, for every $n \times n$ matrix $A \in M_n(K)$, the *characteristic polynomial* $P_A(X)$ of A is the determinant

$$P_A(X) = \det(XI - A).$$

The characteristic polynomial $P_A(X)$ is a polynomial in $K[X]$, the ring of polynomials in the indeterminate X with coefficients in the ring K . For example, when $n = 2$, if

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix},$$

then

$$P_A(X) = \begin{vmatrix} X - a & -b \\ -c & X - d \end{vmatrix} = X^2 - (a + d)X + ad - bc.$$

We can substitute the matrix A for the variable X in the polynomial $P_A(X)$, obtaining a matrix P_A . If we write

$$P_A(X) = X^n + c_1X^{n-1} + \cdots + c_n,$$

then

$$P_A = A^n + c_1A^{n-1} + \cdots + c_nI.$$

We have the following remarkable theorem.

Theorem 5.15. (*Cayley–Hamilton*) *If K is any commutative ring, for every $n \times n$ matrix $A \in M_n(K)$, if we let*

$$P_A(X) = X^n + c_1X^{n-1} + \cdots + c_n$$

be the characteristic polynomial of A , then

$$P_A = A^n + c_1A^{n-1} + \cdots + c_nI = 0.$$

Proof. We can view the matrix $B = XI - A$ as a matrix with coefficients in the polynomial ring $K[X]$, and then we can form the matrix \tilde{B} which is the transpose of the matrix of cofactors of elements of B . Each entry in \tilde{B} is an $(n-1) \times (n-1)$ determinant, and thus a polynomial of degree at most $n-1$, so we can write \tilde{B} as

$$\tilde{B} = X^{n-1}B_0 + X^{n-2}B_1 + \cdots + B_{n-1},$$

for some matrices B_0, \dots, B_{n-1} with coefficients in K . For example, when $n = 2$, we have

$$B = \begin{pmatrix} X - a & -b \\ -c & X - d \end{pmatrix}, \quad \tilde{B} = \begin{pmatrix} X - d & b \\ c & X - a \end{pmatrix} = X \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} + \begin{pmatrix} -d & b \\ c & -a \end{pmatrix}.$$

By Proposition 5.10, we have

$$B\tilde{B} = \det(B)I = P_A(X)I.$$

On the other hand, we have

$$B\tilde{B} = (XI - A)(X^{n-1}B_0 + X^{n-2}B_1 + \cdots + X^{n-j-1}B_j + \cdots + B_{n-1}),$$

and by multiplying out the right-hand side, we get

$$B\tilde{B} = X^n D_0 + X^{n-1} D_1 + \cdots + X^{n-j} D_j + \cdots + D_n,$$

with

$$\begin{aligned} D_0 &= B_0 \\ D_1 &= B_1 - AB_0 \\ &\vdots \\ D_j &= B_j - AB_{j-1} \\ &\vdots \\ D_{n-1} &= B_{n-1} - AB_{n-2} \\ D_n &= -AB_{n-1}. \end{aligned}$$

Since

$$P_A(X)I = (X^n + c_1 X^{n-1} + \cdots + c_n)I,$$

the equality

$$X^n D_0 + X^{n-1} D_1 + \cdots + D_n = (X^n + c_1 X^{n-1} + \cdots + c_n)I$$

is an equality between two matrices, so it requires that all corresponding entries are equal, and since these are polynomials, the coefficients of these polynomials must be identical, which is equivalent to the set of equations

$$\begin{aligned} I &= B_0 \\ c_1 I &= B_1 - AB_0 \\ &\vdots \\ c_j I &= B_j - AB_{j-1} \\ &\vdots \\ c_{n-1} I &= B_{n-1} - AB_{n-2} \\ c_n I &= -AB_{n-1}, \end{aligned}$$

for all j , with $1 \leq j \leq n-1$. If we multiply the first equation by A^n , the last by I , and generally the $(j+1)$ th by A^{n-j} , when we add up all these new equations, we see that the right-hand side adds up to 0, and we get our desired equation

$$A^n + c_1 A^{n-1} + \cdots + c_n I = 0,$$

as claimed. □

As a concrete example, when $n = 2$, the matrix

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

satisfies the equation

$$A^2 - (a + d)A + (ad - bc)I = 0.$$

Most readers will probably find the proof of Theorem 5.15 rather clever but very mysterious and unmotivated. The conceptual difficulty is that we really need to understand how polynomials in one variable “act” on vectors, in terms of the matrix A . This can be done and yields a more “natural” proof. Actually, the reasoning is simpler and more general if we free ourselves from matrices and instead consider a finite-dimensional vector space E and some given linear map $f: E \rightarrow E$. Given any polynomial $p(X) = a_0X^n + a_1X^{n-1} + \cdots + a_n$ with coefficients in the field K , we define the *linear map* $p(f): E \rightarrow E$ by

$$p(f) = a_0f^n + a_1f^{n-1} + \cdots + a_n\text{id},$$

where $f^k = f \circ \cdots \circ f$, the k -fold composition of f with itself. Note that

$$p(f)(u) = a_0f^n(u) + a_1f^{n-1}(u) + \cdots + a_nu,$$

for every vector $u \in E$. Then, we define a new kind of scalar multiplication $\cdot: K[X] \times E \rightarrow E$ by polynomials as follows: for every polynomial $p(X) \in K[X]$, for every $u \in E$,

$$p(X) \cdot u = p(f)(u).$$

It is easy to verify that this is a “good action,” which means that

$$\begin{aligned} p \cdot (u + v) &= p \cdot u + p \cdot v \\ (p + q) \cdot u &= p \cdot u + q \cdot u \\ (pq) \cdot u &= p \cdot (q \cdot u) \\ 1 \cdot u &= u, \end{aligned}$$

for all $p, q \in K[X]$ and all $u, v \in E$. With this new scalar multiplication, E is a $K[X]$ -module.

If $p = \lambda$ is just a scalar in K (a polynomial of degree 0), then

$$\lambda \cdot u = (\lambda\text{id})(u) = \lambda u,$$

which means that K acts on E by scalar multiplication as before. If $p(X) = X$ (the monomial X), then

$$X \cdot u = f(u).$$

Now, if we pick a basis (e_1, \dots, e_n) , if a polynomial $p(X) \in K[X]$ has the property that

$$p(X) \cdot e_i = 0, \quad i = 1, \dots, n,$$

then this means that $p(f)(e_i) = 0$ for $i = 1, \dots, n$, which means that the linear map $p(f)$ vanishes on E . We can also check, as we did in Section 5.2, that if A and B are two $n \times n$ matrices and if (u_1, \dots, u_n) are any n vectors, then

$$A \cdot \left(B \cdot \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix} \right) = (AB) \cdot \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix}.$$

This suggests the plan of attack for our second proof of the Cayley–Hamilton theorem. For simplicity, we prove the theorem for vector spaces over a field. The proof goes through for a free module over a commutative ring.

Theorem 5.16. (*Cayley–Hamilton*) *For every finite-dimensional vector space over a field K , for every linear map $f: E \rightarrow E$, for every basis (e_1, \dots, e_n) , if A is the matrix over f over the basis (e_1, \dots, e_n) and if*

$$P_A(X) = X^n + c_1 X^{n-1} + \dots + c_n$$

is the characteristic polynomial of A , then

$$P_A(f) = f^n + c_1 f^{n-1} + \dots + c_n \text{id} = 0.$$

Proof. Since the columns of A consist of the vector $f(e_j)$ expressed over the basis (e_1, \dots, e_n) , we have

$$f(e_j) = \sum_{i=1}^n a_{ij} e_i, \quad 1 \leq j \leq n.$$

Using our action of $K[X]$ on E , the above equations can be expressed as

$$X \cdot e_j = \sum_{i=1}^n a_{ij} \cdot e_i, \quad 1 \leq j \leq n,$$

which yields

$$\sum_{i=1}^{j-1} -a_{ij} \cdot e_i + (X - a_{jj}) \cdot e_j + \sum_{i=j+1}^n -a_{ij} \cdot e_i = 0, \quad 1 \leq j \leq n.$$

Observe that the transpose of the characteristic polynomial shows up, so the above system can be written as

$$\begin{pmatrix} X - a_{11} & -a_{21} & \cdots & -a_{n1} \\ -a_{12} & X - a_{22} & \cdots & -a_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ -a_{1n} & -a_{2n} & \cdots & X - a_{nn} \end{pmatrix} \cdot \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}.$$

If we let $B = XI - A^\top$, then as in the previous proof, if \tilde{B} is the transpose of the matrix of cofactors of B , we have

$$\tilde{B}B = \det(B)I = \det(XI - A^\top)I = \det(XI - A)I = P_AI.$$

But then, since

$$B \cdot \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix},$$

and since \tilde{B} is matrix whose entries are polynomials in $K[X]$, it makes sense to multiply on the left by \tilde{B} and we get

$$\tilde{B} \cdot B \cdot \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix} = (\tilde{B}B) \cdot \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix} = P_AI \cdot \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix} = \tilde{B} \cdot \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix};$$

that is,

$$P_A \cdot e_j = 0, \quad j = 1, \dots, n,$$

which proves that $P_A(f) = 0$, as claimed. \square

If K is a field, then the characteristic polynomial of a linear map $f: E \rightarrow E$ is independent of the basis (e_1, \dots, e_n) chosen in E . To prove this, observe that the matrix of f over another basis will be of the form $P^{-1}AP$, for some invertible matrix P , and then

$$\begin{aligned} \det(XI - P^{-1}AP) &= \det(XP^{-1}IP - P^{-1}AP) \\ &= \det(P^{-1}(XI - A)P) \\ &= \det(P^{-1}) \det(XI - A) \det(P) \\ &= \det(XI - A). \end{aligned}$$

Therefore, the characteristic polynomial of a linear map is intrinsic to f , and it is denoted by P_f .

The zeros (roots) of the characteristic polynomial of a linear map f are called the *eigenvalues* of f . They play an important role in theory and applications. We will come back to this topic later on.

5.8 Permanents

Recall that the explicit formula for the determinant of an $n \times n$ matrix is

$$\det(A) = \sum_{\pi \in \mathfrak{S}_n} \epsilon(\pi) a_{\pi(1)1} \cdots a_{\pi(n)n}.$$

If we drop the sign $\epsilon(\pi)$ of every permutation from the above formula, we obtain a quantity known as the *permanent*:

$$\text{per}(A) = \sum_{\pi \in \mathfrak{S}_n} a_{\pi(1)1} \cdots a_{\pi(n)n}.$$

Permanents and determinants were investigated as early as 1812 by Cauchy. It is clear from the above definition that the permanent is a multilinear and symmetric form. We also have

$$\text{per}(A) = \text{per}(A^T),$$

and the following unsigned version of the Laplace expansion formula:

$$\text{per}(A) = a_{i1}\text{per}(A_{i1}) + \cdots + a_{ij}\text{per}(A_{ij}) + \cdots + a_{in}\text{per}(A_{in}),$$

for $i = 1, \dots, n$. However, unlike determinants which have a clear geometric interpretation as signed volumes, permanents do not have any natural geometric interpretation. Furthermore, determinants can be evaluated efficiently, for example using the conversion to row reduced echelon form, but computing the permanent is hard.

Permanents turn out to have various combinatorial interpretations. One of these is in terms of perfect matchings of bipartite graphs which we now discuss.

Recall that a *bipartite* (undirected) graph $G = (V, E)$ is a graph whose set of nodes V can be partitioned into two nonempty disjoint subsets V_1 and V_2 , such that every edge $e \in E$ has one endpoint in V_1 and one endpoint in V_2 . An example of a bipartite graph with 14 nodes is shown in Figure 5.8; its nodes are partitioned into the two sets $\{x_1, x_2, x_3, x_4, x_5, x_6, x_7\}$ and $\{y_1, y_2, y_3, y_4, y_5, y_6, y_7\}$.

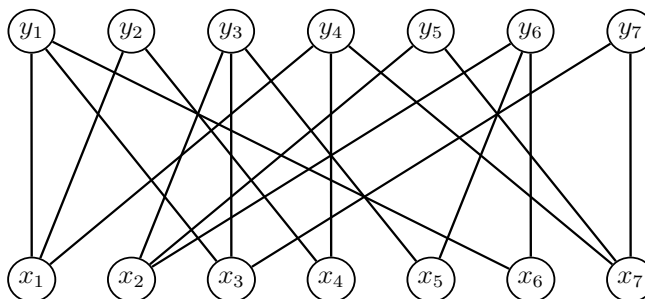
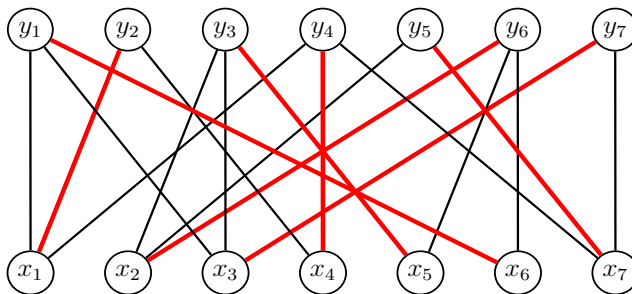


Figure 5.1: A bipartite graph G .

A *matching* in a graph $G = (V, E)$ (bipartite or not) is a set M of pairwise non-adjacent edges, which means that no two edges in M share a common vertex. A *perfect matching* is a matching such that every node in V is incident to some edge in the matching M (every node in V is an endpoint of some edge in M). Figure 5.8 shows a perfect matching (in red) in the bipartite graph G .

Figure 5.2: A perfect matching in the bipartite graph G .

Obviously, a perfect matching in a bipartite graph can exist only if its set of nodes has a partition in two blocks of equal size, say $\{x_1, \dots, x_m\}$ and $\{y_1, \dots, y_m\}$. Then, there is a bijection between perfect matchings and bijections $\pi: \{x_1, \dots, x_m\} \rightarrow \{y_1, \dots, y_m\}$ such that $\pi(x_i) = y_j$ iff there is an edge between x_i and y_j .

Now, every bipartite graph G with a partition of its nodes into two sets of equal size as above is represented by an $m \times m$ matrix $A = (a_{ij})$ such that $a_{ij} = 1$ iff there is an edge between x_i and y_j , and $a_{ij} = 0$ otherwise. Using the interpretation of perfect matchings as bijections $\pi: \{x_1, \dots, x_m\} \rightarrow \{y_1, \dots, y_m\}$, we see that *the permanent $\text{per}(A)$ of the $(0, 1)$ -matrix A representing the bipartite graph G counts the number of perfect matchings in G .*

In a famous paper published in 1979, Leslie Valiant proves that computing the permanent is a $\#P$ -complete problem. Such problems are suspected to be intractable. It is known that if a polynomial-time algorithm existed to solve a $\#P$ -complete problem, then we would have $P = NP$, which is believed to be very unlikely.

Another combinatorial interpretation of the permanent can be given in terms of systems of distinct representatives. Given a finite set S , let (A_1, \dots, A_n) be any sequence of nonempty subsets of S (not necessarily distinct). A *system of distinct representatives* (for short *SDR*) of the sets A_1, \dots, A_n is a sequence of n distinct elements (a_1, \dots, a_n) , with $a_i \in A_i$ for $i = 1, \dots, n$. The number of SDR's of a sequence of sets plays an important role in combinatorics. Now, if $S = \{1, 2, \dots, n\}$ and if we associate to any sequence (A_1, \dots, A_n) of nonempty subsets of S the matrix $A = (a_{ij})$ defined such that $a_{ij} = 1$ if $j \in A_i$ and $a_{ij} = 0$ otherwise, then *the permanent $\text{per}(A)$ counts the number of SDR's of the set A_1, \dots, A_n .*

This interpretation of permanents in terms of SDR's can be used to prove bounds for the permanents of various classes of matrices. Interested readers are referred to van Lint and Wilson [113] (Chapters 11 and 12). In particular, a proof of a theorem known as *Van der Waerden conjecture* is given in Chapter 12. This theorem states that for any $n \times n$ matrix A with nonnegative entries in which all row-sums and column-sums are 1 (doubly stochastic matrices), we have

$$\text{per}(A) \geq \frac{n!}{n^n},$$

with equality for the matrix in which all entries are equal to $1/n$.

5.9 Further Readings

Thorough expositions of the material covered in Chapters 2–4 and 5 can be found in Strang [105, 104], Lax [71], Lang [67], Artin [4], Mac Lane and Birkhoff [73], Hoffman and Kunze [62], Bourbaki [14, 15], Van Der Waerden [112], Serre [96], Horn and Johnson [57], and Bertin [12]. These notions of linear algebra are nicely put to use in classical geometry, see Berger [8, 9], Tisseron [109] and Dieudonné [28].

Chapter 6

Gaussian Elimination, LU -Factorization, Cholesky Factorization, Reduced Row Echelon Form

6.1 Motivating Example: Curve Interpolation

Curve interpolation is a problem that arises frequently in computer graphics and in robotics (path planning). There are many ways of tackling this problem and in this section we will describe a solution using *cubic splines*. Such splines consist of cubic Bézier curves. They are often used because they are cheap to implement and give more flexibility than quadratic Bézier curves.

A *cubic Bézier curve* $C(t)$ (in \mathbb{R}^2 or \mathbb{R}^3) is specified by a list of four *control points* (b_0, b_1, b_2, b_3) and is given parametrically by the equation

$$C(t) = (1-t)^3 b_0 + 3(1-t)^2 t b_1 + 3(1-t) t^2 b_2 + t^3 b_3.$$

Clearly, $C(0) = b_0$, $C(1) = b_3$, and for $t \in [0, 1]$, the point $C(t)$ belongs to the convex hull of the control points b_0, b_1, b_2, b_3 . The polynomials

$$(1-t)^3, \quad 3(1-t)^2 t, \quad 3(1-t) t^2, \quad t^3$$

are the *Bernstein polynomials* of degree 3.

Typically, we are only interested in the curve segment corresponding to the values of t in the interval $[0, 1]$. Still, the placement of the control points drastically affects the shape of the curve segment, which can even have a self-intersection; See Figures 6.1, 6.2, 6.3 illustrating various configurations.

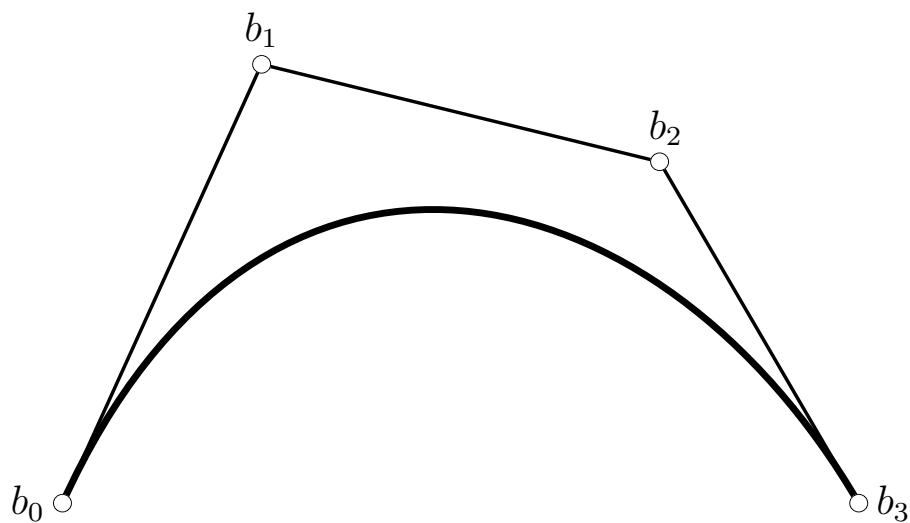


Figure 6.1: A “standard” Bézier curve

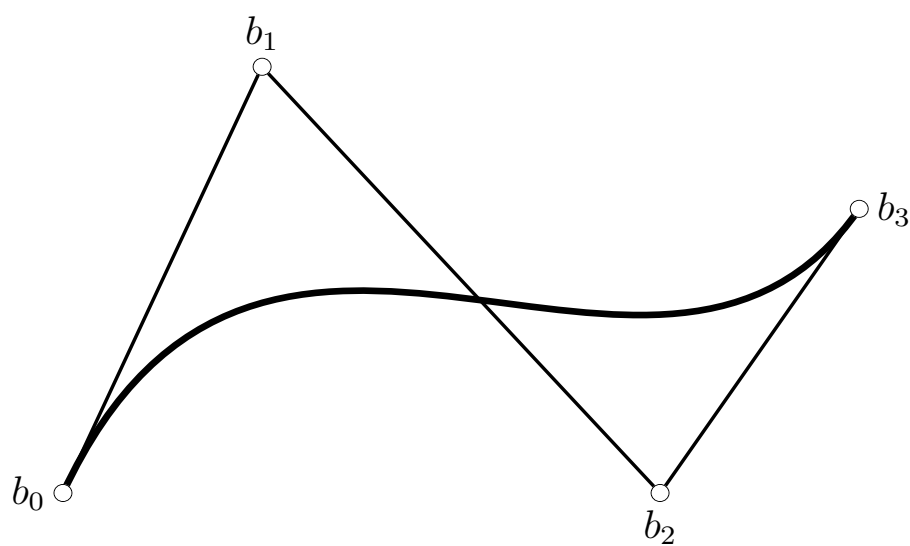


Figure 6.2: A Bézier curve with an inflexion point

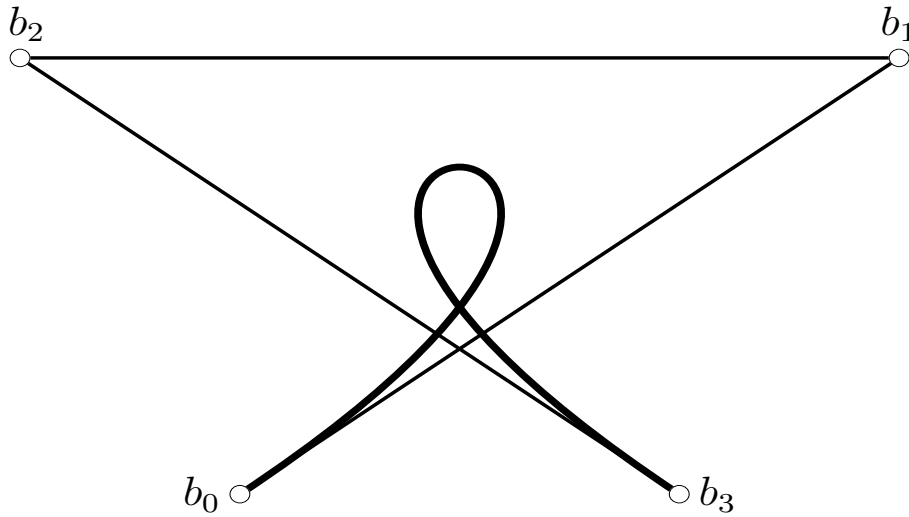


Figure 6.3: A self-intersecting Bézier curve

Interpolation problems require finding curves passing through some given data points and possibly satisfying some extra constraints.

A *Bézier spline curve* F is a curve which is made up of curve segments which are Bézier curves, say C_1, \dots, C_m ($m \geq 2$). We will assume that F defined on $[0, m]$, so that for $i = 1, \dots, m$,

$$F(t) = C_i(t - i + 1), \quad i - 1 \leq t \leq i.$$

Typically, some smoothness is required between any two junction points, that is, between any two points $C_i(1)$ and $C_{i+1}(0)$, for $i = 1, \dots, m - 1$. We require that $C_i(1) = C_{i+1}(0)$ (C^0 -continuity), and typically that the derivatives of C_i at 1 and of C_{i+1} at 0 agree up to second order derivatives. This is called C^2 -continuity, and it ensures that the tangents agree as well as the curvatures.

There are a number of interpolation problems, and we consider one of the most common problems which can be stated as follows:

Problem: Given $N + 1$ data points x_0, \dots, x_N , find a C^2 cubic spline curve F such that $F(i) = x_i$ for all i , $0 \leq i \leq N$ ($N \geq 2$).

A way to solve this problem is to find $N + 3$ auxiliary points d_{-1}, \dots, d_{N+1} , called *de Boor control points*, from which N Bézier curves can be found. Actually,

$$d_{-1} = x_0 \quad \text{and} \quad d_{N+1} = x_N$$

so we only need to find $N + 1$ points d_0, \dots, d_N .

It turns out that the C^2 -continuity constraints on the N Bézier curves yield only $N - 1$ equations, so d_0 and d_N can be chosen arbitrarily. In practice, d_0 and d_N are chosen according to various *end conditions*, such as prescribed velocities at x_0 and x_N . For the time being, we will assume that d_0 and d_N are given.

Figure 6.4 illustrates an interpolation problem involving $N + 1 = 7 + 1 = 8$ data points. The control points d_0 and d_7 were chosen arbitrarily.

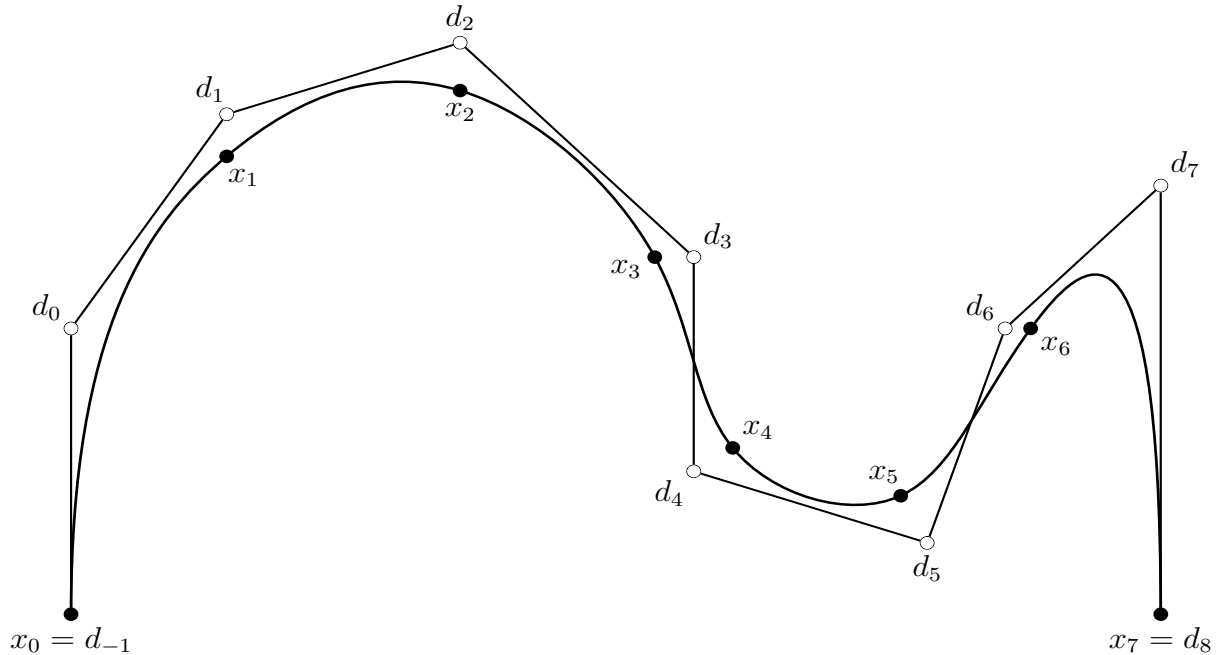


Figure 6.4: A C^2 cubic interpolation spline curve passing through the points $x_0, x_1, x_2, x_3, x_4, x_5, x_6, x_7$

It can be shown that d_1, \dots, d_{N-1} are given by the linear system

$$\begin{pmatrix} \frac{7}{2} & 1 & & & \\ 1 & 4 & 1 & & 0 \\ & \ddots & \ddots & \ddots & \\ 0 & & 1 & 4 & 1 \\ & & & 1 & \frac{7}{2} \end{pmatrix} \begin{pmatrix} d_1 \\ d_2 \\ \vdots \\ d_{N-2} \\ d_{N-1} \end{pmatrix} = \begin{pmatrix} 6x_1 - \frac{3}{2}d_0 \\ 6x_2 \\ \vdots \\ 6x_{N-2} \\ 6x_{N-1} - \frac{3}{2}d_N \end{pmatrix}.$$

We will show later that the above matrix is invertible because it is strictly diagonally dominant.

Once the above system is solved, the Bézier cubics C_1, \dots, C_N are determined as follows (we assume $N \geq 2$): For $2 \leq i \leq N-1$, the control points $(b_0^i, b_1^i, b_2^i, b_3^i)$ of C_i are given by

$$\begin{aligned} b_0^i &= x_{i-1} \\ b_1^i &= \frac{2}{3}d_{i-1} + \frac{1}{3}d_i \\ b_2^i &= \frac{1}{3}d_{i-1} + \frac{2}{3}d_i \\ b_3^i &= x_i. \end{aligned}$$

The control points $(b_0^1, b_1^1, b_2^1, b_3^1)$ of C_1 are given by

$$\begin{aligned} b_0^1 &= x_0 \\ b_1^1 &= d_0 \\ b_2^1 &= \frac{1}{2}d_0 + \frac{1}{2}d_1 \\ b_3^1 &= x_1, \end{aligned}$$

and the control points $(b_0^N, b_1^N, b_2^N, b_3^N)$ of C_N are given by

$$\begin{aligned} b_0^N &= x_{N-1} \\ b_1^N &= \frac{1}{2}d_{N-1} + \frac{1}{2}d_N \\ b_2^N &= d_N \\ b_3^N &= x_N. \end{aligned}$$

We will now describe various methods for solving linear systems. Since the matrix of the above system is tridiagonal, there are specialized methods which are more efficient than the general methods. We will discuss a few of these methods.

6.2 Gaussian Elimination and LU-Factorization

Let A be an $n \times n$ matrix, let $b \in \mathbb{R}^n$ be an n -dimensional vector and assume that A is invertible. Our goal is to solve the system $Ax = b$. Since A is assumed to be invertible, we know that this system has a unique solution $x = A^{-1}b$. Experience shows that two counter-intuitive facts are revealed:

- (1) One should avoid computing the inverse A^{-1} of A explicitly. This is because this would amount to solving the n linear systems $Au^{(j)} = e_j$ for $j = 1, \dots, n$, where $e_j = (0, \dots, 1, \dots, 0)$ is the j th canonical basis vector of \mathbb{R}^n (with a 1 in the j th slot). By doing so, we would replace the resolution of a single system by the resolution of n systems, and we would still have to multiply A^{-1} by b .

- (2) One does not solve (large) linear systems by computing determinants (using Cramer's formulae). This is because this method requires a number of additions (resp. multiplications) proportional to $(n+1)!$ (resp. $(n+2)!$).

The key idea on which most direct methods (as opposed to iterative methods, that look for an approximation of the solution) are based is that if A is an upper-triangular matrix, which means that $a_{ij} = 0$ for $1 \leq j < i \leq n$ (resp. lower-triangular, which means that $a_{ij} = 0$ for $1 \leq i < j \leq n$), then computing the solution x is trivial. Indeed, say A is an upper-triangular matrix

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n-2} & a_{1n-1} & a_{1n} \\ 0 & a_{22} & \cdots & a_{2n-2} & a_{2n-1} & a_{2n} \\ 0 & 0 & \ddots & \vdots & \vdots & \vdots \\ & & & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & a_{n-1n-1} & a_{n-1n} \\ 0 & 0 & \cdots & 0 & 0 & a_{nn} \end{pmatrix}.$$

Then, $\det(A) = a_{11}a_{22}\cdots a_{nn} \neq 0$, which implies that $a_{ii} \neq 0$ for $i = 1, \dots, n$, and we can solve the system $Ax = b$ from bottom-up by *back-substitution*. That is, first we compute x_n from the last equation, next plug this value of x_n into the next to the last equation and compute x_{n-1} from it, *etc.* This yields

$$\begin{aligned} x_n &= a_{nn}^{-1}b_n \\ x_{n-1} &= a_{n-1n-1}^{-1}(b_{n-1} - a_{n-1n}x_n) \\ &\vdots \\ x_1 &= a_{11}^{-1}(b_1 - a_{12}x_2 - \cdots - a_{1n}x_n). \end{aligned}$$

Note that the use of determinants can be avoided to prove that if A is invertible then $a_{ii} \neq 0$ for $i = 1, \dots, n$. Indeed, it can be shown directly (by induction) that an upper (or lower) triangular matrix is invertible iff all its diagonal entries are nonzero.

If A is lower-triangular, we solve the system from top-down by *forward-substitution*.

Thus, what we need is a method for transforming a matrix to an equivalent one in upper-triangular form. This can be done by *elimination*. Let us illustrate this method on the following example:

$$\begin{array}{rrcrcl} 2x & + & y & + & z & = & 5 \\ 4x & - & 6y & & & = & -2 \\ -2x & + & 7y & + & 2z & = & 9. \end{array}$$

We can eliminate the variable x from the second and the third equation as follows: Subtract twice the first equation from the second and add the first equation to the third. We get the

new system

$$\begin{array}{rclcl} 2x & + & y & + & z & = & 5 \\ & - & 8y & - & 2z & = & -12 \\ & & 8y & + & 3z & = & 14. \end{array}$$

This time, we can eliminate the variable y from the third equation by adding the second equation to the third:

$$\begin{array}{rclcl} 2x & + & y & + & z & = & 5 \\ & - & 8y & - & 2z & = & -12 \\ & & & & z & = & 2. \end{array}$$

This last system is upper-triangular. Using back-substitution, we find the solution: $z = 2$, $y = 1$, $x = 1$.

Observe that we have performed only row operations. The general method is to iteratively eliminate variables using simple row operations (namely, adding or subtracting a multiple of a row to another row of the matrix) while simultaneously applying these operations to the vector b , to obtain a system, $MAx = Mb$, where MA is upper-triangular. Such a method is called *Gaussian elimination*. However, one extra twist is needed for the method to work in all cases: It may be necessary to permute rows, as illustrated by the following example:

$$\begin{array}{rclcl} x & + & y & + & z & = & 1 \\ x & + & y & + & 3z & = & 1 \\ 2x & + & 5y & + & 8z & = & 1. \end{array}$$

In order to eliminate x from the second and third row, we subtract the first row from the second and we subtract twice the first row from the third:

$$\begin{array}{rclcrcl} x & + & y & + & z & = & 1 \\ & & & & 2z & = & 0 \\ & & 3y & + & 6z & = & -1. \end{array}$$

Now, the trouble is that y does not occur in the second row; so, we can't eliminate y from the third row by adding or subtracting a multiple of the second row to it. The remedy is simple: Permute the second and the third row! We get the system:

$$\begin{array}{rclcl} x & + & y & + & z & = & 1 \\ & & 3y & + & 6z & = & -1 \\ & & & & 2z & = & 0, \end{array}$$

which is already in triangular form. Another example where some permutations are needed is:

$$\begin{array}{rcl} & & z & = & 1 \\ -2x & + & 7y & + & 2z & = & 1 \\ 4x & - & 6y & & & = & -1. \end{array}$$

First, we permute the first and the second row, obtaining

$$\begin{array}{rclcl} -2x & + & 7y & + & 2z & = & 1 \\ & & & & z & = & 1 \\ 4x & - & 6y & & & = & -1, \end{array}$$

and then, we add twice the first row to the third, obtaining:

$$\begin{array}{rclcl} -2x & + & 7y & + & 2z & = & 1 \\ & & & & z & = & 1 \\ & & 8y & + & 4z & = & 1. \end{array}$$

Again, we permute the second and the third row, getting

$$\begin{array}{rclcl} -2x & + & 7y & + & 2z & = & 1 \\ & & 8y & + & 4z & = & 1 \\ & & & & z & = & 1, \end{array}$$

an upper-triangular system. Of course, in this example, z is already solved and we could have eliminated it first, but for the general method, we need to proceed in a systematic fashion.

We now describe the method of *Gaussian Elimination* applied to a linear system $Ax = b$, where A is assumed to be invertible. We use the variable k to keep track of the stages of elimination. Initially, $k = 1$.

- (1) The first step is to pick some nonzero entry a_{i_1} in the first column of A . Such an entry must exist, since A is invertible (otherwise, the first column of A would be the zero vector, and the columns of A would not be linearly independent. Equivalently, we would have $\det(A) = 0$). The actual choice of such an element has some impact on the numerical stability of the method, but this will be examined later. For the time being, we assume that some arbitrary choice is made. This chosen element is called the *pivot* of the elimination step and is denoted π_1 (so, in this first step, $\pi_1 = a_{i_1}$).
- (2) Next, we permute the row (i) corresponding to the pivot with the first row. Such a step is called *pivoting*. So, after this permutation, the first element of the first row is nonzero.
- (3) We now eliminate the variable x_1 from all rows except the first by adding suitable multiples of the first row to these rows. More precisely we add $-a_{i_1}/\pi_1$ times the first row to the i th row for $i = 2, \dots, n$. At the end of this step, all entries in the first column are zero except the first.
- (4) Increment k by 1. If $k = n$, stop. Otherwise, $k < n$, and then iteratively repeat steps (1), (2), (3) on the $(n - k + 1) \times (n - k + 1)$ subsystem obtained by deleting the first $k - 1$ rows and $k - 1$ columns from the current system.

If we let $A_1 = A$ and $A_k = (a_{ij}^k)$ be the matrix obtained after $k - 1$ elimination steps ($2 \leq k \leq n$), then the k th elimination step is applied to the matrix A_k of the form

$$A_k = \begin{pmatrix} a_{11}^k & a_{12}^k & \cdots & \cdots & \cdots & a_{1n}^k \\ & a_{22}^k & \cdots & \cdots & \cdots & a_{2n}^k \\ & & \ddots & \vdots & & \vdots \\ & & & a_{kk}^k & \cdots & a_{kn}^k \\ & & & \vdots & & \vdots \\ & & & a_{nk}^k & \cdots & a_{nn}^k \end{pmatrix}.$$

Actually, note that

$$a_{ij}^k = a_{ij}^i$$

for all i, j with $1 \leq i \leq k - 2$ and $i \leq j \leq n$, since the first $k - 1$ rows remain unchanged after the $(k - 1)$ th step.

We will prove later that $\det(A_k) = \pm \det(A)$. Consequently, A_k is invertible. The fact that A_k is invertible iff A is invertible can also be shown without determinants from the fact that there is some invertible matrix M_k such that $A_k = M_k A$, as we will see shortly.

Since A_k is invertible, some entry a_{ik}^k with $k \leq i \leq n$ is nonzero. Otherwise, the last $n - k + 1$ entries in the first k columns of A_k would be zero, and the first k columns of A_k would yield k vectors in \mathbb{R}^{k-1} . But then, the first k columns of A_k would be linearly dependent and A_k would not be invertible, a contradiction.

So, one of the entries a_{ik}^k with $k \leq i \leq n$ can be chosen as pivot, and we permute the k th row with the i th row, obtaining the matrix $\alpha^k = (\alpha_{jl}^k)$. The new pivot is $\pi_k = \alpha_{kk}^k$, and we zero the entries $i = k + 1, \dots, n$ in column k by adding $-\alpha_{ik}^k/\pi_k$ times row k to row i . At the end of this step, we have A_{k+1} . Observe that the first $k - 1$ rows of A_k are identical to the first $k - 1$ rows of A_{k+1} .

It is easy to figure out what kind of matrices perform the elementary row operations used during Gaussian elimination. The key point is that if $A = PB$, where A, B are $m \times n$ matrices and P is a square matrix of dimension m , if (as usual) we denote the rows of A and B by A_1, \dots, A_m and B_1, \dots, B_m , then the formula

$$a_{ij} = \sum_{k=1}^m p_{ik} b_{kj}$$

giving the (i, j) th entry in A shows that the i th row of A is a *linear combination* of the rows of B :

$$A_i = p_{i1}B_1 + \cdots + p_{im}B_m.$$

Therefore, *multiplication of a matrix on the left by a square matrix performs row operations*. Similarly, multiplication of a matrix on the right by a square matrix performs column operations

The permutation of the k th row with the i th row is achieved by multiplying A on the left by the *transposition matrix* $P(i, k)$, which is the matrix obtained from the identity matrix by permuting rows i and k , *i.e.*,

$$P(i, k) = \begin{pmatrix} 1 & & & & & & \\ & 1 & & & & & \\ & & 0 & & & 1 & \\ & & & 1 & & & \\ & & & & \ddots & & \\ & & & & & 1 & \\ & 1 & & & & 0 & \\ & & & & & & 1 & \\ & & & & & & & 1 \end{pmatrix}.$$

Observe that $\det(P(i, k)) = -1$. Furthermore, $P(i, k)$ is *symmetric* ($P(i, k)^\top = P(i, k)$), and

$$P(i, k)^{-1} = P(i, k).$$

During the permutation step (2), if row k and row i need to be permuted, the matrix A is multiplied on the left by the matrix P_k such that $P_k = P(i, k)$, else we set $P_k = I$.

Adding β times row j to row i is achieved by multiplying A on the left by the *elementary matrix*,

$$E_{i,j;\beta} = I + \beta e_{ij},$$

where

$$(e_{ij})_{kl} = \begin{cases} 1 & \text{if } k = i \text{ and } l = j \\ 0 & \text{if } k \neq i \text{ or } l \neq j, \end{cases}$$

i.e.,

$$E_{i,j;\beta} = \begin{pmatrix} 1 & & & & & & \\ & 1 & & & & & \\ & & 1 & & & & \\ & & & 1 & & & \\ & & & & \ddots & & \\ & & & & & 1 & \\ & \beta & & & & & 1 \\ & & & & & & & 1 \end{pmatrix} \quad \text{or} \quad E_{i,j;\beta} = \begin{pmatrix} 1 & & & & & & \\ & 1 & & & & & \\ & & 1 & & & \beta & \\ & & & 1 & & & \\ & & & & \ddots & & \\ & & & & & 1 & \\ & & & & & & 1 & \\ & & & & & & & 1 \end{pmatrix}.$$

On the left, $i > j$, and on the right, $i < j$. Observe that the inverse of $E_{i,j;\beta} = I + \beta e_{ij}$ is $E_{i,j;-\beta} = I - \beta e_{ij}$ and that $\det(E_{i,j;\beta}) = 1$. Therefore, during step 3 (the elimination step), the matrix A is multiplied on the left by a product E_k of matrices of the form $E_{i,k;\beta_{i,k}}$, with $i > k$.

Consequently, we see that

$$A_{k+1} = E_k P_k A_k,$$

and then

$$A_k = E_{k-1} P_{k-1} \cdots E_1 P_1 A.$$

This justifies the claim made earlier that $A_k = M_k A$ for some invertible matrix M_k ; we can pick

$$M_k = E_{k-1} P_{k-1} \cdots E_1 P_1,$$

a product of invertible matrices.

The fact that $\det(P(i, k)) = -1$ and that $\det(E_{i,j;\beta}) = 1$ implies immediately the fact claimed above: We always have

$$\det(A_k) = \pm \det(A).$$

Furthermore, since

$$A_k = E_{k-1} P_{k-1} \cdots E_1 P_1 A$$

and since Gaussian elimination stops for $k = n$, the matrix

$$A_n = E_{n-1} P_{n-1} \cdots E_2 P_2 E_1 P_1 A$$

is upper-triangular. Also note that if we let $M = E_{n-1} P_{n-1} \cdots E_2 P_2 E_1 P_1$, then $\det(M) = \pm 1$, and

$$\det(A) = \pm \det(A_n).$$

The matrices $P(i, k)$ and $E_{i,j;\beta}$ are called *elementary matrices*. We can summarize the above discussion in the following theorem:

Theorem 6.1. (*Gaussian Elimination*) *Let A be an $n \times n$ matrix (invertible or not). Then there is some invertible matrix M so that $U = MA$ is upper-triangular. The pivots are all nonzero iff A is invertible.*

Proof. We already proved the theorem when A is invertible, as well as the last assertion. Now, A is singular iff some pivot is zero, say at stage k of the elimination. If so, we must have $a_{i_k}^k = 0$ for $i = k, \dots, n$; but in this case, $A_{k+1} = A_k$ and we may pick $P_k = E_k = I$. \square

Remark: Obviously, the matrix M can be computed as

$$M = E_{n-1} P_{n-1} \cdots E_2 P_2 E_1 P_1,$$

but this expression is of no use. Indeed, what we need is M^{-1} ; when no permutations are needed, it turns out that M^{-1} can be obtained immediately from the matrices E_k 's, in fact, from their inverses, and no multiplications are necessary.

Remark: Instead of looking for an invertible matrix M so that MA is upper-triangular, we can look for an invertible matrix M so that MA is a diagonal matrix. Only a simple change to Gaussian elimination is needed. At every stage, k , after the pivot has been found and pivoting been performed, if necessary, in addition to adding suitable multiples of the k th row to the rows *below* row k in order to zero the entries in column k for $i = k + 1, \dots, n$, also add suitable multiples of the k th row to the rows *above* row k in order to zero the entries in column k for $i = 1, \dots, k - 1$. Such steps are also achieved by multiplying on the left by elementary matrices $E_{i,k;\beta_{i,k}}$, except that $i < k$, so that these matrices are not lower-triangular matrices. Nevertheless, at the end of the process, we find that $A_n = MA$, is a diagonal matrix.

This method is called the *Gauss-Jordan factorization*. Because it is more expansive than Gaussian elimination, this method is not used much in practice. However, Gauss-Jordan factorization can be used to compute the inverse of a matrix A . Indeed, we find the j th column of A^{-1} by solving the system $Ax^{(j)} = e_j$ (where e_j is the j th canonical basis vector of \mathbb{R}^n). By applying Gauss-Jordan, we are led to a system of the form $D_j x^{(j)} = M_j e_j$, where D_j is a diagonal matrix, and we can immediately compute $x^{(j)}$.

It remains to discuss the choice of the pivot, and also conditions that guarantee that no permutations are needed during the Gaussian elimination process. We begin by stating a necessary and sufficient condition for an invertible matrix to have an LU -factorization (*i.e.*, Gaussian elimination does not require pivoting).

We say that an invertible matrix A has an LU -factorization if it can be written as $A = LU$, where U is upper-triangular invertible and L is lower-triangular, with $L_{ii} = 1$ for $i = 1, \dots, n$.

A lower-triangular matrix with diagonal entries equal to 1 is called a *unit lower-triangular* matrix. Given an $n \times n$ matrix $A = (a_{ij})$, for any k with $1 \leq k \leq n$, let $A[1..k, 1..k]$ denote the submatrix of A whose entries are a_{ij} , where $1 \leq i, j \leq k$.

Proposition 6.2. *Let A be an invertible $n \times n$ -matrix. Then, A has an LU -factorization $A = LU$ iff every matrix $A[1..k, 1..k]$ is invertible for $k = 1, \dots, n$. Furthermore, when A has an LU -factorization, we have*

$$\det(A[1..k, 1..k]) = \pi_1 \pi_2 \cdots \pi_k, \quad k = 1, \dots, n,$$

where π_k is the pivot obtained after $k - 1$ elimination steps. Therefore, the k th pivot is given by

$$\pi_k = \begin{cases} a_{11} = \det(A[1..1, 1..1]) & \text{if } k = 1 \\ \frac{\det(A[1..k, 1..k])}{\det(A[1..k-1, 1..k-1])} & \text{if } k = 2, \dots, n. \end{cases}$$

Proof. First, assume that $A = LU$ is an LU -factorization of A . We can write

$$A = \begin{pmatrix} A[1..k, 1..k] & A_2 \\ A_3 & A_4 \end{pmatrix} = \begin{pmatrix} L_1 & 0 \\ L_3 & L_4 \end{pmatrix} \begin{pmatrix} U_1 & U_2 \\ 0 & U_4 \end{pmatrix} = \begin{pmatrix} L_1 U_1 & L_1 U_2 \\ L_3 U_1 & L_3 U_2 + L_4 U_4 \end{pmatrix},$$

where L_1, L_4 are unit lower-triangular and U_1, U_4 are upper-triangular. Thus,

$$A[1..k, 1..k] = L_1 U_1,$$

and since U is invertible, U_1 is also invertible (the determinant of U is the product of the diagonal entries in U , which is the product of the diagonal entries in U_1 and U_4). As L_1 is invertible (since its diagonal entries are equal to 1), we see that $A[1..k, 1..k]$ is invertible for $k = 1, \dots, n$.

Conversely, assume that $A[1..k, 1..k]$ is invertible for $k = 1, \dots, n$. We just need to show that Gaussian elimination does not need pivoting. We prove by induction on k that the k th step does not need pivoting.

This holds for $k = 1$, since $A[1..1, 1..1] = (a_{11})$, so $a_{11} \neq 0$. Assume that no pivoting was necessary for the first $k - 1$ steps ($2 \leq k \leq n - 1$). In this case, we have

$$E_{k-1} \cdots E_2 E_1 A = A_k,$$

where $L = E_{k-1} \cdots E_2 E_1$ is a unit lower-triangular matrix and $A_k[1..k, 1..k]$ is upper-triangular, so that $LA = A_k$ can be written as

$$\begin{pmatrix} L_1 & 0 \\ L_3 & L_4 \end{pmatrix} \begin{pmatrix} A[1..k, 1..k] & A_2 \\ A_3 & A_4 \end{pmatrix} = \begin{pmatrix} U_1 & B_2 \\ 0 & B_4 \end{pmatrix},$$

where L_1 is unit lower-triangular and U_1 is upper-triangular. But then,

$$L_1 A[1..k, 1..k] = U_1,$$

where L_1 is invertible (in fact, $\det(L_1) = 1$), and since by hypothesis $A[1..k, 1..k]$ is invertible, U_1 is also invertible, which implies that $(U_1)_{kk} \neq 0$, since U_1 is upper-triangular. Therefore, no pivoting is needed in step k , establishing the induction step. Since $\det(L_1) = 1$, we also have

$$\det(U_1) = \det(L_1 A[1..k, 1..k]) = \det(L_1) \det(A[1..k, 1..k]) = \det(A[1..k, 1..k]),$$

and since U_1 is upper-triangular and has the pivots π_1, \dots, π_k on its diagonal, we get

$$\det(A[1..k, 1..k]) = \pi_1 \pi_2 \cdots \pi_k, \quad k = 1, \dots, n,$$

as claimed. □

Remark: The use of determinants in the first part of the proof of Proposition 6.2 can be avoided if we use the fact that a triangular matrix is invertible iff all its diagonal entries are nonzero.

Corollary 6.3. (*LU-Factorization*) *Let A be an invertible $n \times n$ -matrix. If every matrix $A[1..k, 1..k]$ is invertible for $k = 1, \dots, n$, then Gaussian elimination requires no pivoting and yields an LU-factorization $A = LU$.*

Proof. We proved in Proposition 6.2 that in this case Gaussian elimination requires no pivoting. Then, since every elementary matrix $E_{i,k;\beta}$ is lower-triangular (since we always arrange that the pivot π_k occurs above the rows that it operates on), since $E_{i,k;\beta}^{-1} = E_{i,k;-\beta}$ and the E'_k s are products of $E_{i,k;\beta_{i,k}}$'s, from

$$E_{n-1} \cdots E_2 E_1 A = U,$$

where U is an upper-triangular matrix, we get

$$A = LU,$$

where $L = E_1^{-1} E_2^{-1} \cdots E_{n-1}^{-1}$ is a lower-triangular matrix. Furthermore, as the diagonal entries of each $E_{i,k;\beta}$ are 1, the diagonal entries of each E_k are also 1. \square

The reader should verify that the example below is indeed an LU -factorization.

$$\begin{pmatrix} 2 & 1 & 1 & 0 \\ 4 & 3 & 3 & 1 \\ 8 & 7 & 9 & 5 \\ 6 & 7 & 9 & 8 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 4 & 3 & 1 & 0 \\ 3 & 4 & 1 & 1 \end{pmatrix} \begin{pmatrix} 2 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 2 \end{pmatrix}.$$

One of the main reasons why the existence of an LU -factorization for a matrix A is interesting is that if we need to solve *several* linear systems $Ax = b$ corresponding to the same matrix A , we can do this cheaply by solving the two triangular systems

$$Lw = b, \quad \text{and} \quad Ux = w.$$

There is a certain asymmetry in the LU -decomposition $A = LU$ of an invertible matrix A . Indeed, the diagonal entries of L are all 1, but this is generally false for U . This asymmetry can be eliminated as follows: if

$$D = \text{diag}(u_{11}, u_{22}, \dots, u_{nn})$$

is the diagonal matrix consisting of the diagonal entries in U (the pivots), then we if let $U' = D^{-1}U$, we can write

$$A = LDU',$$

where L is lower-triangular, U' is upper-triangular, all diagonal entries of both L and U' are 1, and D is a diagonal matrix of pivots. Such a decomposition is called an LDU -factorization. We will see shortly that if A is symmetric, then $U' = L^\top$.

As we will see a bit later, symmetric positive definite matrices satisfy the condition of Proposition 6.2. Therefore, linear systems involving symmetric positive definite matrices can be solved by Gaussian elimination without pivoting. Actually, it is possible to do better: This is the Cholesky factorization.

The following easy proposition shows that, in principle, A can be premultiplied by some permutation matrix P , so that PA can be converted to upper-triangular form without using any pivoting. Permutations are discussed in some detail in Section 20.3, but for now we just need their definition. A *permutation matrix* is a square matrix that has a single 1 in every row and every column and zeros everywhere else. It is shown in Section 20.3 that every permutation matrix is a product of transposition matrices (the $P(i, k)$ s), and that P is invertible with inverse P^\top .

Proposition 6.4. *Let A be an invertible $n \times n$ -matrix. Then, there is some permutation matrix P so that $PA[1..k, 1..k]$ is invertible for $k = 1, \dots, n$.*

Proof. The case $n = 1$ is trivial, and so is the case $n = 2$ (we swap the rows if necessary). If $n \geq 3$, we proceed by induction. Since A is invertible, its columns are linearly independent; in particular, its first $n - 1$ columns are also linearly independent. Delete the last column of A . Since the remaining $n - 1$ columns are linearly independent, there are also $n - 1$ linearly independent rows in the corresponding $n \times (n - 1)$ matrix. Thus, there is a permutation of these n rows so that the $(n - 1) \times (n - 1)$ matrix consisting of the first $n - 1$ rows is invertible. But, then, there is a corresponding permutation matrix P_1 , so that the first $n - 1$ rows and columns of $P_1 A$ form an invertible matrix A' . Applying the induction hypothesis to the $(n - 1) \times (n - 1)$ matrix A' , we see that there some permutation matrix P_2 (leaving the n th row fixed), so that $P_2 P_1 A[1..k, 1..k]$ is invertible, for $k = 1, \dots, n - 1$. Since A is invertible in the first place and P_1 and P_2 are invertible, $P_1 P_2 A$ is also invertible, and we are done. \square

Remark: One can also prove Proposition 6.4 using a clever reordering of the Gaussian elimination steps suggested by Trefethen and Bau [110] (Lecture 21). Indeed, we know that if A is invertible, then there are permutation matrices P_i and products of elementary matrices E_i , so that

$$A_n = E_{n-1} P_{n-1} \cdots E_2 P_2 E_1 P_1 A,$$

where $U = A_n$ is upper-triangular. For example, when $n = 4$, we have $E_3 P_3 E_2 P_2 E_1 P_1 A = U$. We can define new matrices E'_1, E'_2, E'_3 which are still products of elementary matrices so that we have

$$E'_3 E'_2 E'_1 P_3 P_2 P_1 A = U.$$

Indeed, if we let $E'_3 = E_3$, $E'_2 = P_3 E_2 P_3^{-1}$, and $E'_1 = P_3 P_2 E_1 P_2^{-1} P_3^{-1}$, we easily verify that each E'_k is a product of elementary matrices and that

$$E'_3 E'_2 E'_1 P_3 P_2 P_1 = E_3 (P_3 E_2 P_3^{-1}) (P_3 P_2 E_1 P_2^{-1} P_3^{-1}) P_3 P_2 P_1 = E_3 P_3 E_2 P_2 E_1 P_1.$$

It can also be proved that E'_1, E'_2, E'_3 are lower triangular (see Theorem 6.5).

In general, we let

$$E'_k = P_{n-1} \cdots P_{k+1} E_k P_{k+1}^{-1} \cdots P_{n-1}^{-1},$$

and we have

$$E'_{n-1} \cdots E'_1 P_{n-1} \cdots P_1 A = U,$$

where each E'_j is a lower triangular matrix (see Theorem 6.5).

Using the above idea, we can prove the theorem below which also shows how to compute P, L and U using a simple adaptation of Gaussian elimination. We are not aware of a detailed proof of Theorem 6.5 in the standard texts. Although Golub and Van Loan [49] state a version of this theorem as their Theorem 3.1.4, they say that “The proof is a messy subscripting argument.” Meyer [80] also provides a sketch of proof (see the end of Section 3.10). In view of this situation, we offer a complete proof. It does involve a lot of subscripts and superscripts, but in our opinion, it contains some interesting techniques that go far beyond symbol manipulation.

Theorem 6.5. *For every invertible $n \times n$ -matrix A , the following hold:*

- (1) *There is some permutation matrix P , some upper-triangular matrix U , and some unit lower-triangular matrix L , so that $PA = LU$ (recall, $L_{ii} = 1$ for $i = 1, \dots, n$). Furthermore, if $P = I$, then L and U are unique and they are produced as a result of Gaussian elimination without pivoting.*
- (2) *If $E_{n-1} \dots E_1 A = U$ is the result of Gaussian elimination without pivoting, write as usual $A_k = E_{k-1} \dots E_1 A$ (with $A_k = (a_{ij}^k)$), and let $\ell_{ik} = a_{ik}^k / a_{kk}^k$, with $1 \leq k \leq n-1$ and $k+1 \leq i \leq n$. Then*

$$L = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ \ell_{21} & 1 & 0 & \cdots & 0 \\ \ell_{31} & \ell_{32} & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & 0 \\ \ell_{n1} & \ell_{n2} & \ell_{n3} & \cdots & 1 \end{pmatrix},$$

where the k th column of L is the k th column of E_k^{-1} , for $k = 1, \dots, n-1$.

- (3) *If $E_{n-1} P_{n-1} \cdots E_1 P_1 A = U$ is the result of Gaussian elimination with some pivoting, write $A_k = E_{k-1} P_{k-1} \cdots E_1 P_1 A$, and define E_j^k , with $1 \leq j \leq n-1$ and $j \leq k \leq n-1$, such that, for $j = 1, \dots, n-2$,*

$$\begin{aligned} E_j^j &= E_j \\ E_j^k &= P_k E_j^{k-1} P_k, \quad \text{for } k = j+1, \dots, n-1, \end{aligned}$$

and

$$E_{n-1}^{n-1} = E_{n-1}.$$

Then,

$$\begin{aligned} E_j^k &= P_k P_{k-1} \cdots P_{j+1} E_j P_{j+1} \cdots P_{k-1} P_k \\ U &= E_{n-1}^{n-1} \cdots E_1^{n-1} P_{n-1} \cdots P_1 A, \end{aligned}$$

and if we set

$$P = P_{n-1} \cdots P_1$$

$$L = (E_1^{n-1})^{-1} \cdots (E_{n-1}^{n-1})^{-1},$$

then

$$PA = LU.$$

Furthermore,

$$(E_j^k)^{-1} = I + \mathcal{E}_j^k, \quad 1 \leq j \leq n-1, j \leq k \leq n-1,$$

where \mathcal{E}_j^k is a lower triangular matrix of the form

$$\mathcal{E}_j^k = \begin{pmatrix} 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & \cdots & \ell_{j+1j}^k & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \ell_{nj}^k & 0 & \cdots & 0 \end{pmatrix},$$

we have

$$E_j^k = I - \mathcal{E}_j^k,$$

and

$$\mathcal{E}_j^k = P_k \mathcal{E}_j^{k-1}, \quad 1 \leq j \leq n-2, j+1 \leq k \leq n-1,$$

where $P_k = I$ or else $P_k = P(k, i)$ for some i such that $k+1 \leq i \leq n$; if $P_k \neq I$, this means that $(E_j^k)^{-1}$ is obtained from $(E_j^{k-1})^{-1}$ by permuting the entries on row i and k in column j . Because the matrices $(E_j^k)^{-1}$ are all lower triangular, the matrix L is also lower triangular.

In order to find L , define lower triangular matrices Λ_k of the form

$$\Lambda_k = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & \cdots & \cdots & 0 \\ \lambda_{21}^k & 0 & 0 & 0 & 0 & \vdots & \vdots & 0 \\ \lambda_{31}^k & \lambda_{32}^k & \ddots & 0 & 0 & \vdots & \vdots & 0 \\ \vdots & \vdots & \ddots & 0 & 0 & \vdots & \vdots & \vdots \\ \lambda_{k+11}^k & \lambda_{k+12}^k & \cdots & \lambda_{k+1k}^k & 0 & \cdots & \cdots & 0 \\ \lambda_{k+21}^k & \lambda_{k+22}^k & \cdots & \lambda_{k+2k}^k & 0 & \ddots & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \lambda_{n1}^k & \lambda_{n2}^k & \cdots & \lambda_{nk}^k & 0 & \cdots & \cdots & 0 \end{pmatrix}$$

to assemble the columns of L iteratively as follows: let

$$(-\ell_{k+1k}^k, \dots, -\ell_{nk}^k)$$

be the last $n - k$ elements of the k th column of E_k , and define Λ_k inductively by setting

$$\Lambda_1 = \begin{pmatrix} 0 & 0 & \cdots & 0 \\ \ell_{21}^1 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \ell_{n1}^1 & 0 & \cdots & 0 \end{pmatrix},$$

then for $k = 2, \dots, n - 1$, define

$$\Lambda'_k = P_k \Lambda_{k-1},$$

and

$$\Lambda_k = (I + \Lambda'_k) E_k^{-1} - I = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & \cdots & \cdots & 0 \\ \lambda'_{21}{}^{k-1} & 0 & 0 & 0 & 0 & \vdots & \vdots & 0 \\ \lambda'_{31}{}^{k-1} & \lambda'_{32}{}^{k-1} & \ddots & 0 & 0 & \vdots & \vdots & 0 \\ \vdots & \vdots & \ddots & 0 & 0 & \vdots & \vdots & \vdots \\ \lambda'_{k1}{}^{k-1} & \lambda'_{k2}{}^{k-1} & \cdots & \lambda'_{kk}{}^{k-1} & 0 & \cdots & \cdots & 0 \\ \lambda'_{k+11}{}^{k-1} & \lambda'_{k+12}{}^{k-1} & \cdots & \lambda'_{k+1,k-1}{}^{k-1} & \ell_{k+1,k}^k & \ddots & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \lambda'_{n1}{}^{k-1} & \lambda'_{n2}{}^{k-1} & \cdots & \lambda'_{nk}{}^{k-1} & \ell_{nk}^k & \cdots & \cdots & 0 \end{pmatrix},$$

with $P_k = I$ or $P_k = P(k, i)$ for some $i > k$. This means that in assembling L , row k and row i of Λ_{k-1} need to be permuted when a pivoting step permuting row k and row i of A_k is required. Then

$$\begin{aligned} I + \Lambda_k &= (E_1^k)^{-1} \cdots (E_k^k)^{-1} \\ \Lambda_k &= \mathcal{E}_1^k \cdots \mathcal{E}_k^k, \end{aligned}$$

for $k = 1, \dots, n - 1$, and therefore

$$L = I + \Lambda_{n-1}.$$

Proof. (1) The only part that has not been proved is the uniqueness part (when $P = I$). Assume that A is invertible and that $A = L_1 U_1 = L_2 U_2$, with L_1, L_2 unit lower-triangular and U_1, U_2 upper-triangular. Then, we have

$$L_2^{-1} L_1 = U_2 U_1^{-1}.$$

However, it is obvious that L_2^{-1} is lower-triangular and that U_1^{-1} is upper-triangular, and so $L_2^{-1} L_1$ is lower-triangular and $U_2 U_1^{-1}$ is upper-triangular. Since the diagonal entries of L_1 and L_2 are 1, the above equality is only possible if $U_2 U_1^{-1} = I$, that is, $U_1 = U_2$, and so $L_1 = L_2$.

(2) When $P = I$, we have $L = E_1^{-1}E_2^{-1}\cdots E_{n-1}^{-1}$, where E_k is the product of $n - k$ elementary matrices of the form $E_{i,k;-\ell_i}$, where $E_{i,k;-\ell_i}$ subtracts ℓ_i times row k from row i , with $\ell_{ik} = a_{ik}^k/a_{kk}^k$, $1 \leq k \leq n - 1$, and $k + 1 \leq i \leq n$. Then, it is immediately verified that

$$E_k = \begin{pmatrix} 1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 1 & 0 & \cdots & 0 \\ 0 & \cdots & -\ell_{k+1k} & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & -\ell_{nk} & 0 & \cdots & 1 \end{pmatrix},$$

and that

$$E_k^{-1} = \begin{pmatrix} 1 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 1 & 0 & \cdots & 0 \\ 0 & \cdots & \ell_{k+1k} & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \ell_{nk} & 0 & \cdots & 1 \end{pmatrix}.$$

If we define L_k by

$$L_k = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & \vdots & 0 \\ \ell_{21} & 1 & 0 & 0 & 0 & \vdots & 0 \\ \ell_{31} & \ell_{32} & \ddots & 0 & 0 & \vdots & 0 \\ \vdots & \vdots & \ddots & 1 & 0 & \vdots & 0 \\ \ell_{k+11} & \ell_{k+12} & \cdots & \ell_{k+1k} & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & 0 & \vdots & 0 \\ \ell_{n1} & \ell_{n2} & \cdots & \ell_{nk} & 0 & \cdots & 1 \end{pmatrix}$$

for $k = 1, \dots, n - 1$, we easily check that $L_1 = E_1^{-1}$, and that

$$L_k = L_{k-1}E_k^{-1}, \quad 2 \leq k \leq n - 1,$$

because multiplication on the right by E_k^{-1} adds ℓ_i times column i to column k (of the matrix L_{k-1}) with $i > k$, and column i of L_{k-1} has only the nonzero entry 1 as its i th element. Since

$$L_k = E_1^{-1} \cdots E_k^{-1}, \quad 1 \leq k \leq n - 1,$$

we conclude that $L = L_{n-1}$, proving our claim about the shape of L .

(3) First, we prove by induction on k that

$$A_{k+1} = E_k^k \cdots E_1^k P_k \cdots P_1 A, \quad k = 1, \dots, n - 2.$$

For $k = 1$, we have $A_2 = E_1 P_1 A = E_1^1 P_1 A$, since $E_1^1 = E_1$, so our assertion holds trivially.

Now, if $k \geq 2$,

$$A_{k+1} = E_k P_k A_k,$$

and by the induction hypothesis,

$$A_k = E_{k-1}^{k-1} \cdots E_2^{k-1} E_1^{k-1} P_{k-1} \cdots P_1 A.$$

Because P_k is either the identity or a transposition, $P_k^2 = I$, so by inserting occurrences of $P_k P_k$ as indicated below we can write

$$\begin{aligned} A_{k+1} &= E_k P_k A_k \\ &= E_k P_k E_{k-1}^{k-1} \cdots E_2^{k-1} E_1^{k-1} P_{k-1} \cdots P_1 A \\ &= E_k P_k E_{k-1}^{k-1} (P_k P_k) \cdots (P_k P_k) E_2^{k-1} (P_k P_k) E_1^{k-1} (P_k P_k) P_{k-1} \cdots P_1 A \\ &= E_k (P_k E_{k-1}^{k-1} P_k) \cdots (P_k E_2^{k-1} P_k) (P_k E_1^{k-1} P_k) P_k P_{k-1} \cdots P_1 A. \end{aligned}$$

Observe that P_k has been “moved” to the right of the elimination steps. However, by definition,

$$\begin{aligned} E_j^k &= P_k E_j^{k-1} P_k, \quad j = 1, \dots, k-1 \\ E_k^k &= E_k, \end{aligned}$$

so we get

$$A_{k+1} = E_k^k E_{k-1}^k \cdots E_2^k E_1^k P_k \cdots P_1 A,$$

establishing the induction hypothesis. For $k = n-2$, we get

$$U = A_{n-1} = E_{n-1}^{n-1} \cdots E_1^{n-1} P_{n-1} \cdots P_1 A,$$

as claimed, and the factorization $PA = LU$ with

$$\begin{aligned} P &= P_{n-1} \cdots P_1 \\ L &= (E_1^{n-1})^{-1} \cdots (E_{n-1}^{n-1})^{-1} \end{aligned}$$

is clear,

Since for $j = 1, \dots, n-2$, we have $E_j^j = E_j$,

$$E_j^k = P_k E_j^{k-1} P_k, \quad k = j+1, \dots, n-1,$$

since $E_{n-1}^{n-1} = E_{n-1}$ and $P_k^{-1} = P_k$, we get $(E_j^j)^{-1} = E_j^{-1}$ for $j = 1, \dots, n-1$, and for $j = 1, \dots, n-2$, we have

$$(E_j^k)^{-1} = P_k (E_j^{k-1})^{-1} P_k, \quad k = j+1, \dots, n-1.$$

Since

$$(E_j^{k-1})^{-1} = I + \mathcal{E}_j^{k-1}$$

and $P_k = P(k, i)$ is a transposition, $P_k^2 = I$, so we get

$$(E_j^k)^{-1} = P_k(E_j^{k-1})^{-1}P_k = P_k(I + \mathcal{E}_j^{k-1})P_k = P_k^2 + P_k \mathcal{E}_j^{k-1} P_k = I + P_k \mathcal{E}_j^{k-1} P_k.$$

Therefore, we have

$$(E_j^k)^{-1} = I + P_k \mathcal{E}_j^{k-1} P_k, \quad 1 \leq j \leq n-2, j+1 \leq k \leq n-1.$$

We prove for $j = 1, \dots, n-1$, that for $k = j, \dots, n-1$, each \mathcal{E}_j^k is a lower triangular matrix of the form

$$\mathcal{E}_j^k = \begin{pmatrix} 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & \cdots & \ell_{j+1j}^k & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \ell_{nj}^k & 0 & \cdots & 0 \end{pmatrix},$$

and that

$$\mathcal{E}_j^k = P_k \mathcal{E}_j^{k-1}, \quad 1 \leq j \leq n-2, j+1 \leq k \leq n-1,$$

with $P_k = I$ or $P_k = P(k, i)$ for some i such that $k+1 \leq i \leq n$.

For each j ($1 \leq j \leq n-1$) we proceed by induction on $k = j, \dots, n-1$. Since $(E_j^j)^{-1} = E_j^{-1}$ and since E_j^{-1} is of the above form, the base case holds.

For the induction step, we only need to consider the case where $P_k = P(k, i)$ is a transposition, since the case where $P_k = I$ is trivial. We have to figure out what $P_k \mathcal{E}_j^{k-1} P_k = P(k, i) \mathcal{E}_j^{k-1} P(k, i)$ is. However, since

$$\mathcal{E}_j^{k-1} = \begin{pmatrix} 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & \cdots & \ell_{j+1j}^{k-1} & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \ell_{nj}^{k-1} & 0 & \cdots & 0 \end{pmatrix},$$

and because $k+1 \leq i \leq n$ and $j \leq k-1$, multiplying \mathcal{E}_j^{k-1} on the right by $P(k, i)$ will permute *columns* i and k , which are columns of zeros, so

$$P(k, i) \mathcal{E}_j^{k-1} P(k, i) = P(k, i) \mathcal{E}_j^{k-1},$$

and thus,

$$(E_j^k)^{-1} = I + P(k, i) \mathcal{E}_j^{k-1},$$

which shows that

$$\mathcal{E}_j^k = P(k, i) \mathcal{E}_j^{k-1}.$$

We also know that multiplying $(\mathcal{E}_j^{k-1})^{-1}$ on the left by $P(k, i)$ will permute rows i and k , which shows that \mathcal{E}_j^k has the desired form, as claimed. Since all \mathcal{E}_j^k are strictly lower triangular, all $(E_j^k)^{-1} = I + \mathcal{E}_j^k$ are lower triangular, so the product

$$L = (E_1^{n-1})^{-1} \cdots (E_{n-1}^{n-1})^{-1}$$

is also lower triangular.

From the beginning of part (3), we know that

$$L = (E_1^{n-1})^{-1} \cdots (E_{n-1}^{n-1})^{-1}.$$

We prove by induction on k that

$$\begin{aligned} I + \Lambda_k &= (E_1^k)^{-1} \cdots (E_k^k)^{-1} \\ \Lambda_k &= \mathcal{E}_1^k \cdots \mathcal{E}_k^k, \end{aligned}$$

for $k = 1, \dots, n-1$.

If $k = 1$, we have $E_1^1 = E_1$ and

$$E_1 = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ -\ell_{21}^1 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ -\ell_{n1}^1 & 0 & \cdots & 1 \end{pmatrix}.$$

We get

$$(E_1^{-1})^{-1} = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ \ell_{21}^1 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \ell_{n1}^1 & 0 & \cdots & 1 \end{pmatrix} = I + \Lambda_1,$$

Since $(E_1^{-1})^{-1} = I + \mathcal{E}_1^1$, we also get $\Lambda_1 = \mathcal{E}_1^1$, and the base step holds.

Since $(E_j^k)^{-1} = I + \mathcal{E}_j^k$ with

$$\mathcal{E}_j^k = \begin{pmatrix} 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & \cdots & \ell_{j+1j}^k & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \ell_{nj}^k & 0 & \cdots & 0 \end{pmatrix},$$

as in part (2) for the computation involving the products of L_k 's, we get

$$(E_1^{k-1})^{-1} \cdots (E_{k-1}^{k-1})^{-1} = I + \mathcal{E}_1^{k-1} \cdots \mathcal{E}_{k-1}^{k-1}, \quad 2 \leq k \leq n. \quad (*)$$

Similarly, from the fact that $\mathcal{E}_j^{k-1} P(k, i) = \mathcal{E}_j^{k-1}$ if $i \geq k + 1$ and $j \leq k - 1$ and since

$$(E_j^k)^{-1} = I + P_k \mathcal{E}_j^{k-1}, \quad 1 \leq j \leq n - 2, j + 1 \leq k \leq n - 1,$$

we get

$$(E_1^k)^{-1} \cdots (E_{k-1}^k)^{-1} = I + P_k \mathcal{E}_1^{k-1} \cdots \mathcal{E}_{k-1}^{k-1}, \quad 2 \leq k \leq n - 1. \quad (**)$$

By the induction hypothesis,

$$I + \Lambda_{k-1} = (E_1^{k-1})^{-1} \cdots (E_{k-1}^{k-1})^{-1},$$

and from (*), we get

$$\Lambda_{k-1} = \mathcal{E}_1^{k-1} \cdots \mathcal{E}_{k-1}^{k-1}.$$

Using (**), we deduce that

$$(E_1^k)^{-1} \cdots (E_{k-1}^k)^{-1} = I + P_k \Lambda_{k-1}.$$

Since $E_k^k = E_k$, we obtain

$$(E_1^k)^{-1} \cdots (E_{k-1}^k)^{-1} (E_k^k)^{-1} = (I + P_k \Lambda_{k-1}) E_k^{-1}.$$

However, by definition

$$I + \Lambda_k = (I + P_k \Lambda_{k-1}) E_k^{-1},$$

which proves that

$$I + \Lambda_k = (E_1^k)^{-1} \cdots (E_{k-1}^k)^{-1} (E_k^k)^{-1}, \quad (\dagger)$$

and finishes the induction step for the proof of this formula.

If we apply equation (*) again with $k + 1$ in place of k , we have

$$(E_1^k)^{-1} \cdots (E_k^k)^{-1} = I + \mathcal{E}_1^k \cdots \mathcal{E}_k^k,$$

and together with (\dagger), we obtain,

$$\Lambda_k = \mathcal{E}_1^k \cdots \mathcal{E}_k^k,$$

also finishing the induction step for the proof of this formula. For $k = n - 1$ in (\dagger), we obtain the desired equation: $L = I + \Lambda_{n-1}$. \square

Part (3) of Theorem 6.5 shows the remarkable fact that in assembling the matrix L while performing Gaussian elimination with pivoting, the only change to the algorithm is to make the same transposition on the rows of L (really Λ_k , since the one's are not altered) that we make on the rows of A (really A_k) during a pivoting step involving row k and row i . We can also assemble P by starting with the identity matrix and applying to P the same row transpositions that we apply to A and Λ . Here is an example illustrating this method.

Consider the matrix

$$A = \begin{pmatrix} 1 & 2 & -3 & 4 \\ 4 & 8 & 12 & -8 \\ 2 & 3 & 2 & 1 \\ -3 & -1 & 1 & -4 \end{pmatrix}.$$

We set $P_0 = I_4$, and we can also set $\Lambda_0 = 0$. The first step is to permute row 1 and row 2, using the pivot 4. We also apply this permutation to P_0 :

$$A'_1 = \begin{pmatrix} 4 & 8 & 12 & -8 \\ 1 & 2 & -3 & 4 \\ 2 & 3 & 2 & 1 \\ -3 & -1 & 1 & -4 \end{pmatrix} \quad P_1 = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

Next, we subtract $1/4$ times row 1 from row 2, $1/2$ times row 1 from row 3, and add $3/4$ times row 1 to row 4, and start assembling Λ :

$$A_2 = \begin{pmatrix} 4 & 8 & 12 & -8 \\ 0 & 0 & -6 & 6 \\ 0 & -1 & -4 & 5 \\ 0 & 5 & 10 & -10 \end{pmatrix} \quad \Lambda_1 = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 1/4 & 0 & 0 & 0 \\ 1/2 & 0 & 0 & 0 \\ -3/4 & 0 & 0 & 0 \end{pmatrix} \quad P_1 = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

Next we permute row 2 and row 4, using the pivot 5. We also apply this permutation to Λ and P :

$$A'_3 = \begin{pmatrix} 4 & 8 & 12 & -8 \\ 0 & 5 & 10 & -10 \\ 0 & -1 & -4 & 5 \\ 0 & 0 & -6 & 6 \end{pmatrix} \quad \Lambda'_2 = \begin{pmatrix} 0 & 0 & 0 & 0 \\ -3/4 & 0 & 0 & 0 \\ 1/2 & 0 & 0 & 0 \\ 1/4 & 0 & 0 & 0 \end{pmatrix} \quad P_2 = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}.$$

Next we add $1/5$ times row 2 to row 3, and update Λ'_2 :

$$A_3 = \begin{pmatrix} 4 & 8 & 12 & -8 \\ 0 & 5 & 10 & -10 \\ 0 & 0 & -2 & 3 \\ 0 & 0 & -6 & 6 \end{pmatrix} \quad \Lambda_2 = \begin{pmatrix} 0 & 0 & 0 & 0 \\ -3/4 & 0 & 0 & 0 \\ 1/2 & -1/5 & 0 & 0 \\ 1/4 & 0 & 0 & 0 \end{pmatrix} \quad P_2 = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}.$$

Next we permute row 3 and row 4, using the pivot -6 . We also apply this permutation to Λ and P :

$$A'_4 = \begin{pmatrix} 4 & 8 & 12 & -8 \\ 0 & 5 & 10 & -10 \\ 0 & 0 & -6 & 6 \\ 0 & 0 & -2 & 3 \end{pmatrix} \quad \Lambda'_3 = \begin{pmatrix} 0 & 0 & 0 & 0 \\ -3/4 & 0 & 0 & 0 \\ 1/4 & 0 & 0 & 0 \\ 1/2 & -1/5 & 0 & 0 \end{pmatrix} \quad P_3 = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}.$$

Finally, we subtract $1/3$ times row 3 from row 4, and update Λ'_3 :

$$A_4 = \begin{pmatrix} 4 & 8 & 12 & -8 \\ 0 & 5 & 10 & -10 \\ 0 & 0 & -6 & 6 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad \Lambda_3 = \begin{pmatrix} 0 & 0 & 0 & 0 \\ -3/4 & 0 & 0 & 0 \\ 1/4 & 0 & 0 & 0 \\ 1/2 & -1/5 & 1/3 & 0 \end{pmatrix} \quad P_3 = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}.$$

Consequently, adding the identity to Λ_3 , we obtain

$$L = \begin{pmatrix} 1 & 0 & 0 & 0 \\ -3/4 & 1 & 0 & 0 \\ 1/4 & 0 & 1 & 0 \\ 1/2 & -1/5 & 1/3 & 1 \end{pmatrix}, \quad U = \begin{pmatrix} 4 & 8 & 12 & -8 \\ 0 & 5 & 10 & -10 \\ 0 & 0 & -6 & 6 \\ 0 & 0 & 0 & 1 \end{pmatrix}, \quad P = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}.$$

We check that

$$PA = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 2 & -3 & 4 \\ 4 & 8 & 12 & -8 \\ 2 & 3 & 2 & 1 \\ -3 & -1 & 1 & -4 \end{pmatrix} = \begin{pmatrix} 4 & 8 & 12 & -8 \\ -3 & -1 & 1 & -4 \\ 1 & 2 & -3 & 4 \\ 2 & 3 & 2 & 1 \end{pmatrix},$$

and that

$$LU = \begin{pmatrix} 1 & 0 & 0 & 0 \\ -3/4 & 1 & 0 & 0 \\ 1/4 & 0 & 1 & 0 \\ 1/2 & -1/5 & 1/3 & 1 \end{pmatrix} \begin{pmatrix} 4 & 8 & 12 & -8 \\ 0 & 5 & 10 & -10 \\ 0 & 0 & -6 & 6 \\ 0 & 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 4 & 8 & 12 & -8 \\ -3 & -1 & 1 & -4 \\ 1 & 2 & -3 & 4 \\ 2 & 3 & 2 & 1 \end{pmatrix} = PA.$$

Note that if one willing to overwrite the lower triangular part of the evolving matrix A , one can store the evolving Λ there, since these entries will eventually be zero anyway! There is also no need to save explicitly the permutation matrix P . One could instead record the permutation steps in an extra column (record the vector $(\pi(1), \dots, \pi(n))$ corresponding to the permutation π applied to the rows). We let the reader write such a bold and space-efficient version of LU -decomposition!

As a corollary of Theorem 6.5(1), we can show the following result.

Proposition 6.6. *If an invertible symmetric matrix A has an LU -decomposition, then A has a factorization of the form*

$$A = LDL^\top,$$

where L is a lower-triangular matrix whose diagonal entries are equal to 1, and where D consists of the pivots. Furthermore, such a decomposition is unique.

Proof. If A has an LU -factorization, then it has an LDU factorization

$$A = LDU,$$

where L is lower-triangular, U is upper-triangular, and the diagonal entries of both L and U are equal to 1. Since A is symmetric, we have

$$LDU = A = A^\top = U^\top DL^\top,$$

with U^\top lower-triangular and DL^\top upper-triangular. By the uniqueness of LU -factorization (part (1) of Theorem 6.5), we must have $L = U^\top$ (and $DU = DL^\top$), thus $U = L^\top$, as claimed. \square

Remark: It can be shown that Gaussian elimination + back-substitution requires $n^3/3 + O(n^2)$ additions, $n^3/3 + O(n^2)$ multiplications and $n^2/2 + O(n)$ divisions.

Let us now briefly comment on the choice of a pivot. Although theoretically, any pivot can be chosen, the possibility of roundoff errors implies that it is not a good idea to pick very small pivots. The following example illustrates this point. Consider the linear system

$$\begin{array}{rcrcrcrcrcl} 10^{-4}x & + & y & = & 1 \\ x & + & y & = & 2. \end{array}$$

Since 10^{-4} is nonzero, it can be taken as pivot, and we get

$$\begin{array}{rcrcrcrcrcrcl} 10^{-4}x & + & y & = & 1 \\ & & (1 - 10^4)y & = & 2 - 10^4. \end{array}$$

Thus, the exact solution is

$$x = \frac{10^4}{10^4 - 1}, \quad y = \frac{10^4 - 2}{10^4 - 1}.$$

However, if roundoff takes place on the fourth digit, then $10^4 - 1 = 9999$ and $10^4 - 2 = 9998$ will be rounded off both to 9990, and then the solution is $x = 0$ and $y = 1$, very far from the exact solution where $x \approx 1$ and $y \approx 1$. The problem is that we picked a very small pivot. If instead we permute the equations, the pivot is 1, and after elimination, we get the system

$$\begin{array}{rcrcrcrcrcrcl} x & + & y & = & 2 \\ & & (1 - 10^{-4})y & = & 1 - 2 \times 10^{-4}. \end{array}$$

This time, $1 - 10^{-4} = 0.9999$ and $1 - 2 \times 10^{-4} = 0.9998$ are rounded off to 0.999 and the solution is $x = 1, y = 1$, much closer to the exact solution.

To remedy this problem, one may use the strategy of *partial pivoting*. This consists of choosing during step k ($1 \leq k \leq n - 1$) one of the entries a_{ik}^k such that

$$|a_{ik}^k| = \max_{k \leq p \leq n} |a_{pk}^k|.$$

By maximizing the value of the pivot, we avoid dividing by undesirably small pivots.

Remark: A matrix, A , is called *strictly column diagonally dominant* iff

$$|a_{jj}| > \sum_{i=1, i \neq j}^n |a_{ij}|, \quad \text{for } j = 1, \dots, n$$

(resp. *strictly row diagonally dominant* iff

$$|a_{ii}| > \sum_{j=1, j \neq i}^n |a_{ij}|, \quad \text{for } i = 1, \dots, n.)$$

It has been known for a long time (before 1900, say by Hadamard) that if a matrix A is strictly column diagonally dominant (resp. strictly row diagonally dominant), then it is invertible. (This is a good exercise, try it!) It can also be shown that if A is strictly column diagonally dominant, then Gaussian elimination with partial pivoting does not actually require pivoting (See Problem 21.6 in Trefethen and Bau [110], or Question 2.19 in Demmel [27]).

Another strategy, called *complete pivoting*, consists in choosing some entry a_{ij}^k , where $k \leq i, j \leq n$, such that

$$|a_{ij}^k| = \max_{k \leq p, q \leq n} |a_{pq}^k|.$$

However, in this method, if the chosen pivot is not in column k , it is also necessary to permute columns. This is achieved by multiplying on the right by a permutation matrix. However, complete pivoting tends to be too expensive in practice, and partial pivoting is the method of choice.

A special case where the LU -factorization is particularly efficient is the case of tridiagonal matrices, which we now consider.

6.3 Gaussian Elimination of Tridiagonal Matrices

Consider the tridiagonal matrix

$$A = \begin{pmatrix} b_1 & c_1 & & & \\ a_2 & b_2 & c_2 & & \\ & a_3 & b_3 & c_3 & \\ & \ddots & \ddots & \ddots & \\ & & a_{n-2} & b_{n-2} & c_{n-2} \\ & & & a_{n-1} & b_{n-1} & c_{n-1} \\ & & & & a_n & b_n \end{pmatrix}.$$

Define the sequence

$$\delta_0 = 1, \quad \delta_1 = b_1, \quad \delta_k = b_k \delta_{k-1} - a_k c_{k-1} \delta_{k-2}, \quad 2 \leq k \leq n.$$

Proposition 6.7. *If A is the tridiagonal matrix above, then $\delta_k = \det(A[1..k, 1..k])$ for $k = 1, \dots, n$.*

Proof. By expanding $\det(A[1..k, 1..k])$ with respect to its last row, the proposition follows by induction on k . \square

Theorem 6.8. *If A is the tridiagonal matrix above and $\delta_k \neq 0$ for $k = 1, \dots, n$, then A has the following LU -factorization:*

$$A = \begin{pmatrix} 1 & & & & & \\ a_2 \frac{\delta_0}{\delta_1} & 1 & & & & \\ & a_3 \frac{\delta_1}{\delta_2} & 1 & & & \\ & & \ddots & \ddots & & \\ & & & a_{n-1} \frac{\delta_{n-3}}{\delta_{n-2}} & 1 & \\ & & & a_n \frac{\delta_{n-2}}{\delta_{n-1}} & & 1 \end{pmatrix} \begin{pmatrix} \frac{\delta_1}{\delta_0} & c_1 & & & & \\ & \frac{\delta_2}{\delta_1} & c_2 & & & \\ & & \frac{\delta_3}{\delta_2} & c_3 & & \\ & & & \ddots & \ddots & \\ & & & & \frac{\delta_{n-1}}{\delta_{n-2}} & c_{n-1} \\ & & & & & \frac{\delta_n}{\delta_{n-1}} \end{pmatrix}.$$

Proof. Since $\delta_k = \det(A[1..k, 1..k]) \neq 0$ for $k = 1, \dots, n$, by Theorem 6.5 (and Proposition 6.2), we know that A has a unique LU -factorization. Therefore, it suffices to check that the proposed factorization works. We easily check that

$$\begin{aligned} (LU)_{k,k+1} &= c_k, & 1 \leq k \leq n-1 \\ (LU)_{k,k-1} &= a_k, & 2 \leq k \leq n \\ (LU)_{kl} &= 0, & |k-l| \geq 2 \\ (LU)_{11} &= \frac{\delta_1}{\delta_0} = b_1 \\ (LU)_{kk} &= \frac{a_k c_{k-1} \delta_{k-2} + \delta_k}{\delta_{k-1}} = b_k, & 2 \leq k \leq n, \end{aligned}$$

since $\delta_k = b_k \delta_{k-1} - a_k c_{k-1} \delta_{k-2}$. \square

It follows that there is a simple method to solve a linear system $Ax = d$ where A is tridiagonal (and $\delta_k \neq 0$ for $k = 1, \dots, n$). For this, it is convenient to “squeeze” the diagonal matrix Δ defined such that $\Delta_{kk} = \delta_k / \delta_{k-1}$ into the factorization so that $A = (L\Delta)(\Delta^{-1}U)$, and if we let

$$z_1 = \frac{c_1}{b_1}, \quad z_k = c_k \frac{\delta_{k-1}}{\delta_k}, \quad 2 \leq k \leq n-1, \quad z_n = \frac{\delta_n}{\delta_{n-1}} = b_n - a_n z_{n-1},$$

$A = (L\Delta)(\Delta^{-1}U)$ is written as

$$A = \begin{pmatrix} \frac{c_1}{z_1} & & & & & \\ a_2 & \frac{c_2}{z_2} & & & & \\ & a_3 & \frac{c_3}{z_3} & & & \\ & & \ddots & \ddots & & \\ & & & a_{n-1} & \frac{c_{n-1}}{z_{n-1}} & \\ & & & & a_n & z_n \end{pmatrix} \begin{pmatrix} 1 & z_1 & & & & \\ & 1 & z_2 & & & \\ & & 1 & z_3 & & \\ & & & \ddots & \ddots & \\ & & & & 1 & z_{n-2} \\ & & & & & 1 & z_{n-1} \\ & & & & & & 1 \end{pmatrix}.$$

As a consequence, the system $Ax = d$ can be solved by constructing three sequences: First, the sequence

$$z_1 = \frac{c_1}{b_1}, \quad z_k = \frac{c_k}{b_k - a_k z_{k-1}}, \quad k = 2, \dots, n-1, \quad z_n = b_n - a_n z_{n-1},$$

corresponding to the recurrence $\delta_k = b_k \delta_{k-1} - a_k c_{k-1} \delta_{k-2}$ and obtained by dividing both sides of this equation by δ_{k-1} , next

$$w_1 = \frac{d_1}{b_1}, \quad w_k = \frac{d_k - a_k w_{k-1}}{b_k - a_k z_{k-1}}, \quad k = 2, \dots, n,$$

corresponding to solving the system $L\Delta w = d$, and finally

$$x_n = w_n, \quad x_k = w_k - z_k x_{k+1}, \quad k = n-1, n-2, \dots, 1,$$

corresponding to solving the system $\Delta^{-1}Ux = w$.

Remark: It can be verified that this requires $3(n-1)$ additions, $3(n-1)$ multiplications, and $2n$ divisions, a total of $8n-6$ operations, which is much less than the $O(2n^3/3)$ required by Gaussian elimination in general.

We now consider the special case of symmetric positive definite matrices (SPD matrices). Recall that an $n \times n$ symmetric matrix A is *positive definite* iff

$$x^\top A x > 0 \quad \text{for all } x \in \mathbb{R}^n \text{ with } x \neq 0.$$

Equivalently, A is symmetric positive definite iff all its eigenvalues are strictly positive. The following facts about a symmetric positive definite matrix A are easily established (some left as an exercise):

- (1) The matrix A is invertible. (Indeed, if $Ax = 0$, then $x^\top Ax = 0$, which implies $x = 0$.)
- (2) We have $a_{ii} > 0$ for $i = 1, \dots, n$. (Just observe that for $x = e_i$, the i th canonical basis vector of \mathbb{R}^n , we have $e_i^\top A e_i = a_{ii} > 0$.)
- (3) For every $n \times n$ invertible matrix Z , the matrix $Z^\top A Z$ is symmetric positive definite iff A is symmetric positive definite.

Next, we prove that a symmetric positive definite matrix has a special LU -factorization of the form $A = BB^\top$, where B is a lower-triangular matrix whose diagonal elements are strictly positive. This is the *Cholesky factorization*.

6.4 SPD Matrices and the Cholesky Decomposition

First, we note that a symmetric positive definite matrix satisfies the condition of Proposition 6.2.

Proposition 6.9. *If A is a symmetric positive definite matrix, then $A[1..k, 1..k]$ is symmetric positive definite, and thus invertible for $k = 1, \dots, n$.*

Proof. Since A is symmetric, each $A[1..k, 1..k]$ is also symmetric. If $w \in \mathbb{R}^k$, with $1 \leq k \leq n$, we let $x \in \mathbb{R}^n$ be the vector with $x_i = w_i$ for $i = 1, \dots, k$ and $x_i = 0$ for $i = k+1, \dots, n$. Now, since A is symmetric positive definite, we have $x^\top A x > 0$ for all $x \in \mathbb{R}^n$ with $x \neq 0$. This holds in particular for all vectors x obtained from nonzero vectors $w \in \mathbb{R}^k$ as defined earlier, and clearly

$$x^\top A x = w^\top A[1..k, 1..k] w,$$

which implies that $A[1..k, 1..k]$ is positive definite. Thus, $A[1..k, 1..k]$ is also invertible. \square

Proposition 6.9 can be strengthened as follows: *A symmetric matrix A is positive definite iff $\det(A[1..k, 1..k]) > 0$ for $k = 1, \dots, n$.*

The above fact is known as *Sylvester's criterion*. We will prove it after establishing the Cholesky factorization.

Let A be an $n \times n$ symmetric positive definite matrix and write

$$A = \begin{pmatrix} a_{11} & W^\top \\ W & C \end{pmatrix},$$

where C is an $(n-1) \times (n-1)$ symmetric matrix and W is an $(n-1) \times 1$ matrix. Since A is symmetric positive definite, $a_{11} > 0$, and we can compute $\alpha = \sqrt{a_{11}}$. The trick is that we can factor A uniquely as

$$A = \begin{pmatrix} a_{11} & W^\top \\ W & C \end{pmatrix} = \begin{pmatrix} \alpha & 0 \\ W/\alpha & I \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & C - WW^\top/a_{11} \end{pmatrix} \begin{pmatrix} \alpha & W^\top/\alpha \\ 0 & I \end{pmatrix},$$

i.e., as $A = B_1 A_1 B_1^\top$, where B_1 is lower-triangular with positive diagonal entries. Thus, B_1 is invertible, and by fact (3) above, A_1 is also symmetric positive definite.

Theorem 6.10. (*Cholesky Factorization*) Let A be a symmetric positive definite matrix. Then, there is some lower-triangular matrix B so that $A = BB^\top$. Furthermore, B can be chosen so that its diagonal elements are strictly positive, in which case B is unique.

Proof. We proceed by induction on the dimension n of A . For $n = 1$, we must have $a_{11} > 0$, and if we let $\alpha = \sqrt{a_{11}}$ and $B = (\alpha)$, the theorem holds trivially. If $n \geq 2$, as we explained above, again we must have $a_{11} > 0$, and we can write

$$A = \begin{pmatrix} a_{11} & W^\top \\ W & C \end{pmatrix} = \begin{pmatrix} \alpha & 0 \\ W/\alpha & I \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & C - WW^\top/a_{11} \end{pmatrix} \begin{pmatrix} \alpha & W^\top/\alpha \\ 0 & I \end{pmatrix} = B_1 A_1 B_1^\top,$$

where $\alpha = \sqrt{a_{11}}$, the matrix B_1 is invertible and

$$A_1 = \begin{pmatrix} 1 & 0 \\ 0 & C - WW^\top/a_{11} \end{pmatrix}$$

is symmetric positive definite. However, this implies that $C - WW^\top/a_{11}$ is also symmetric positive definite (consider $x^\top A_1 x$ for every $x \in \mathbb{R}^n$ with $x \neq 0$ and $x_1 = 0$). Thus, we can apply the induction hypothesis to $C - WW^\top/a_{11}$ (which is an $(n-1) \times (n-1)$ matrix), and we find a unique lower-triangular matrix L with positive diagonal entries so that

$$C - WW^\top/a_{11} = LL^\top.$$

But then, we get

$$\begin{aligned} A &= \begin{pmatrix} \alpha & 0 \\ W/\alpha & I \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & C - WW^\top/a_{11} \end{pmatrix} \begin{pmatrix} \alpha & W^\top/\alpha \\ 0 & I \end{pmatrix} \\ &= \begin{pmatrix} \alpha & 0 \\ W/\alpha & I \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & LL^\top \end{pmatrix} \begin{pmatrix} \alpha & W^\top/\alpha \\ 0 & I \end{pmatrix} \\ &= \begin{pmatrix} \alpha & 0 \\ W/\alpha & I \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & L \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & L^\top \end{pmatrix} \begin{pmatrix} \alpha & W^\top/\alpha \\ 0 & I \end{pmatrix} \\ &= \begin{pmatrix} \alpha & 0 \\ W/\alpha & L \end{pmatrix} \begin{pmatrix} \alpha & W^\top/\alpha \\ 0 & L^\top \end{pmatrix}. \end{aligned}$$

Therefore, if we let

$$B = \begin{pmatrix} \alpha & 0 \\ W/\alpha & L \end{pmatrix},$$

we have a unique lower-triangular matrix with positive diagonal entries and $A = BB^\top$.

The uniqueness of the Cholesky decomposition can also be established using the uniqueness of an LU decomposition. Indeed, if $A = B_1 B_1^\top = B_2 B_2^\top$ where B_1 and B_2 are lower triangular with positive diagonal entries, if we let Δ_1 (resp. Δ_2) be the diagonal matrix consisting of the diagonal entries of B_1 (resp. B_2) so that $(\Delta_k)_{ii} = (B_k)_{ii}$ for $k = 1, 2$, then we have two LU decompositions

$$A = (B_1 \Delta_1^{-1})(\Delta_1 B_1^\top) = (B_2 \Delta_2^{-1})(\Delta_2 B_2^\top)$$

with $B_1\Delta_1^{-1}, B_2\Delta_2^{-1}$ unit lower triangular, and $\Delta_1B_1^\top, \Delta_2B_2^\top$ upper triangular. By uniqueness of LU factorization (Theorem 6.5(1)), we have

$$B_1\Delta_1^{-1} = B_2\Delta_2^{-1}, \quad \Delta_1B_1^\top = \Delta_2B_2^\top,$$

and the second equation yields

$$B_1\Delta_1 = B_2\Delta_2. \quad (*)$$

The diagonal entries of $B_1\Delta_1$ are $(B_1)_{ii}^2$ and similarly the diagonal entries of $B_2\Delta_2$ are $(B_2)_{ii}^2$, so the above equation implies that

$$(B_1)_{ii}^2 = (B_2)_{ii}^2, \quad i = 1, \dots, n.$$

Since the diagonal entries of both B_1 and B_2 are assumed to be positive, we must have

$$(B_1)_{ii} = (B_2)_{ii}, \quad i = 1, \dots, n;$$

that is, $\Delta_1 = \Delta_2$, and since both are invertible, we conclude from $(*)$ that $B_1 = B_2$. \square

The proof of Theorem 6.10 immediately yields an algorithm to compute B from A by solving for a lower triangular matrix B such that $A = BB^\top$. For $j = 1, \dots, n$,

$$b_{jj} = \left(a_{jj} - \sum_{k=1}^{j-1} b_{jk}^2 \right)^{1/2},$$

and for $i = j+1, \dots, n$ (and $j = 1, \dots, n-1$)

$$b_{ij} = \left(a_{ij} - \sum_{k=1}^{j-1} b_{ik}b_{jk} \right) / b_{jj}.$$

The above formulae are used to compute the j th column of B from top-down, using the first $j-1$ columns of B previously computed, and the matrix A .

The Cholesky factorization can be used to solve linear systems $Ax = b$ where A is symmetric positive definite: Solve the two systems $Bw = b$ and $B^\top x = w$.

Remark: It can be shown that this method requires $n^3/6 + O(n^2)$ additions, $n^3/6 + O(n^2)$ multiplications, $n^2/2 + O(n)$ divisions, and $O(n)$ square root extractions. Thus, the Cholesky method requires half of the number of operations required by Gaussian elimination (since Gaussian elimination requires $n^3/3 + O(n^2)$ additions, $n^3/3 + O(n^2)$ multiplications, and $n^2/2 + O(n)$ divisions). It also requires half of the space (only B is needed, as opposed to both L and U). Furthermore, it can be shown that Cholesky's method is numerically stable (see Trefethen and Bau [110], Lecture 23).

Remark: If $A = BB^\top$, where B is any invertible matrix, then A is symmetric positive definite.

Proof. Obviously, BB^\top is symmetric, and since B is invertible, B^\top is invertible, and from

$$x^\top Ax = x^\top BB^\top x = (B^\top x)^\top B^\top x,$$

it is clear that $x^\top Ax > 0$ if $x \neq 0$. □

We now give three more criteria for a symmetric matrix to be positive definite.

Proposition 6.11. *Let A be any $n \times n$ symmetric matrix. The following conditions are equivalent:*

- (a) *A is positive definite.*
- (b) *All principal minors of A are positive; that is: $\det(A[1..k, 1..k]) > 0$ for $k = 1, \dots, n$ (Sylvester's criterion).*
- (c) *A has an LU -factorization and all pivots are positive.*
- (d) *A has an LDL^\top -factorization and all pivots in D are positive.*

Proof. By Proposition 6.9, if A is symmetric positive definite, then each matrix $A[1..k, 1..k]$ is symmetric positive definite for $k = 1, \dots, n$. By the Cholesky decomposition, $A[1..k, 1..k] = Q^\top Q$ for some invertible matrix Q , so $\det(A[1..k, 1..k]) = \det(Q)^2 > 0$. This shows that (a) implies (b).

If $\det(A[1..k, 1..k]) > 0$ for $k = 1, \dots, n$, then each $A[1..k, 1..k]$ is invertible. By Proposition 6.2, the matrix A has an LU -factorization, and since the pivots π_k are given by

$$\pi_k = \begin{cases} a_{11} = \det(A[1..1, 1..1]) & \text{if } k = 1 \\ \frac{\det(A[1..k, 1..k])}{\det(A[1..k-1, 1..k-1])} & \text{if } k = 2, \dots, n, \end{cases}$$

we see that $\pi_k > 0$ for $k = 1, \dots, n$. Thus (b) implies (c).

Assume A has an LU -factorization and that the pivots are all positive. Since A is symmetric, this implies that A has a factorization of the form

$$A = LDL^\top,$$

with L lower-triangular with 1's on its diagonal, and where D is a diagonal matrix with positive entries on the diagonal (the pivots). This shows that (c) implies (d).

Given a factorization $A = LDL^\top$ with all pivots in D positive, if we form the diagonal matrix

$$\sqrt{D} = \text{diag}(\sqrt{\pi_1}, \dots, \sqrt{\pi_n})$$

and if we let $B = L\sqrt{D}$, then we have

$$Q = BB^\top,$$

with B lower-triangular and invertible. By the remark before Proposition 6.11, A is positive definite. Hence, (d) implies (a). □

Criterion (c) yields a simple computational test to check whether a symmetric matrix is positive definite. There is one more criterion for a symmetric matrix to be positive definite: its eigenvalues must be positive. We will have to learn about the spectral theorem for symmetric matrices to establish this criterion.

For more on the stability analysis and efficient implementation methods of Gaussian elimination, LU -factoring and Cholesky factoring, see Demmel [27], Trefethen and Bau [110], Ciarlet [24], Golub and Van Loan [49], Meyer [80], Strang [104, 105], and Kincaid and Cheney [63].

6.5 Reduced Row Echelon Form

Gaussian elimination described in Section 6.2 can also be applied to rectangular matrices. This yields a method for determining whether a system $Ax = b$ is solvable, and a description of all the solutions when the system is solvable, for any rectangular $m \times n$ matrix A .

It turns out that the discussion is simpler if we rescale all pivots to be 1, and for this we need a third kind of elementary matrix. For any $\lambda \neq 0$, let $E_{i,\lambda}$ be the $n \times n$ diagonal matrix

$$E_{i,\lambda} = \begin{pmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & 1 & & & \\ & & & \lambda & & \\ & & & & 1 & \\ & & & & & \ddots \\ & & & & & & 1 \end{pmatrix},$$

with $(E_{i,\lambda})_{ii} = \lambda$ ($1 \leq i \leq n$). Note that $E_{i,\lambda}$ is also given by

$$E_{i,\lambda} = I + (\lambda - 1)e_{ii},$$

and that $E_{i,\lambda}$ is invertible with

$$E_{i,\lambda}^{-1} = E_{i,\lambda^{-1}}.$$

Now, after $k - 1$ elimination steps, if the bottom portion

$$(a_{kk}^k, a_{k+1k}^k, \dots, a_{mk}^k)$$

of the k th column of the current matrix A_k is nonzero so that a pivot π_k can be chosen, after a permutation of rows if necessary, we also divide row k by π_k to obtain the pivot 1, and not only do we zero all the entries $i = k + 1, \dots, m$ in column k , but also all the entries $i = 1, \dots, k - 1$, so that the only nonzero entry in column k is a 1 in row k . These row operations are achieved by multiplication on the left by elementary matrices.

If $a_{kk}^k = a_{k+1k}^k = \dots = a_{mk}^k = 0$, we move on to column $k + 1$.

The result is that after performing such elimination steps, we obtain a matrix that has a special shape known as a *reduced row echelon matrix*. Here is an example illustrating this process: Starting from the matrix

$$A_1 = \begin{pmatrix} 1 & 0 & 2 & 1 & 5 \\ 1 & 1 & 5 & 2 & 7 \\ 1 & 2 & 8 & 4 & 12 \end{pmatrix}$$

we perform the following steps

$$A_1 \longrightarrow A_2 = \begin{pmatrix} 1 & 0 & 2 & 1 & 5 \\ 0 & 1 & 3 & 1 & 2 \\ 0 & 2 & 6 & 3 & 7 \end{pmatrix},$$

by subtracting row 1 from row 2 and row 3;

$$A_2 \longrightarrow \begin{pmatrix} 1 & 0 & 2 & 1 & 5 \\ 0 & 2 & 6 & 3 & 7 \\ 0 & 1 & 3 & 1 & 2 \end{pmatrix} \longrightarrow \begin{pmatrix} 1 & 0 & 2 & 1 & 5 \\ 0 & 1 & 3 & 3/2 & 7/2 \\ 0 & 1 & 3 & 1 & 2 \end{pmatrix} \longrightarrow A_3 = \begin{pmatrix} 1 & 0 & 2 & 1 & 5 \\ 0 & 1 & 3 & 3/2 & 7/2 \\ 0 & 0 & 0 & -1/2 & -3/2 \end{pmatrix},$$

after choosing the pivot 2 and permuting row 2 and row 3, dividing row 2 by 2, and subtracting row 2 from row 3;

$$A_3 \longrightarrow \begin{pmatrix} 1 & 0 & 2 & 1 & 5 \\ 0 & 1 & 3 & 3/2 & 7/2 \\ 0 & 0 & 0 & 1 & 3 \end{pmatrix} \longrightarrow A_4 = \begin{pmatrix} 1 & 0 & 2 & 0 & 2 \\ 0 & 1 & 3 & 0 & -1 \\ 0 & 0 & 0 & 1 & 3 \end{pmatrix},$$

after dividing row 3 by $-1/2$, subtracting row 3 from row 1, and subtracting $(3/2) \times$ row 3 from row 2.

It is clear that columns 1, 2 and 4 are linearly independent, that column 3 is a linear combination of columns 1 and 2, and that column 5 is a linear combinations of columns 1, 2, 4.

In general, the sequence of steps leading to a reduced echelon matrix is not unique. For example, we could have chosen 1 instead of 2 as the second pivot in matrix A_2 . Nevertheless, the reduced row echelon matrix obtained from any given matrix is unique; that is, it does not depend on the the sequence of steps that are followed during the reduction process. This fact is not so easy to prove rigorously, but we will do it later.

If we want to solve a linear system of equations of the form $Ax = b$, we apply elementary row operations to both the matrix A and the right-hand side b . To do this conveniently, we form the *augmented matrix* (A, b) , which is the $m \times (n + 1)$ matrix obtained by adding b as an extra column to the matrix A . For example if

$$A = \begin{pmatrix} 1 & 0 & 2 & 1 \\ 1 & 1 & 5 & 2 \\ 1 & 2 & 8 & 4 \end{pmatrix} \quad \text{and} \quad b = \begin{pmatrix} 5 \\ 7 \\ 12 \end{pmatrix},$$

then the augmented matrix is

$$(A, b) = \begin{pmatrix} 1 & 0 & 2 & 1 & 5 \\ 1 & 1 & 5 & 2 & 7 \\ 1 & 2 & 8 & 4 & 12 \end{pmatrix}.$$

Now, for any matrix M , since

$$M(A, b) = (MA, Mb),$$

performing elementary row operations on (A, b) is equivalent to simultaneously performing operations on both A and b . For example, consider the system

$$\begin{array}{rrrrrcl} x_1 & & & + & 2x_3 & + & x_4 & = & 5 \\ x_1 & + & x_2 & + & 5x_3 & + & 2x_4 & = & 7 \\ x_1 & + & 2x_2 & + & 8x_3 & + & 4x_4 & = & 12. \end{array}$$

Its augmented matrix is the matrix

$$(A, b) = \begin{pmatrix} 1 & 0 & 2 & 1 & 5 \\ 1 & 1 & 5 & 2 & 7 \\ 1 & 2 & 8 & 4 & 12 \end{pmatrix}$$

considered above, so the reduction steps applied to this matrix yield the system

$$\begin{array}{rrrrcl} x_1 & & + & 2x_3 & & = & 2 \\ & x_2 & + & 3x_3 & & = & -1 \\ & & & & x_4 & = & 3. \end{array}$$

This reduced system has the same set of solutions as the original, and obviously x_3 can be chosen arbitrarily. Therefore, our system has infinitely many solutions given by

$$x_1 = 2 - 2x_3, \quad x_2 = -1 - 3x_3, \quad x_4 = 3,$$

where x_3 is arbitrary.

The following proposition shows that the set of solutions of a system $Ax = b$ is preserved by any sequence of row operations.

Proposition 6.12. *Given any $m \times n$ matrix A and any vector $b \in \mathbb{R}^m$, for any sequence of elementary row operations E_1, \dots, E_k , if $P = E_k \cdots E_1$ and $(A', b') = P(A, b)$, then the solutions of $Ax = b$ are the same as the solutions of $A'x = b'$.*

Proof. Since each elementary row operation E_i is invertible, so is P , and since $(A', b') = P(A, b)$, then $A' = PA$ and $b' = Pb$. If x is a solution of the original system $Ax = b$, then multiplying both sides by P we get $PAx = Pb$; that is, $A'x = b'$, so x is a solution of the new system. Conversely, assume that x is a solution of the new system, that is $A'x = b'$. Then, because $A' = PA$, $b' = Pb$, and P is invertible, we get

$$Ax = P^{-1}A'x = P^{-1}b' = b,$$

so x is a solution of the original system $Ax = b$. □

Another important fact is this:

Proposition 6.13. *Given a $m \times n$ matrix A , for any sequence of row operations E_1, \dots, E_k , if $P = E_k \cdots E_1$ and $B = PA$, then the subspaces spanned by the rows of A and the rows of B are identical. Therefore, A and B have the same row rank. Furthermore, the matrices A and B also have the same (column) rank.*

Proof. Since $B = PA$, from a previous observation, the rows of B are linear combinations of the rows of A , so the span of the rows of B is a subspace of the span of the rows of A . Since P is invertible, $A = P^{-1}B$, so by the same reasoning the span of the rows of A is a subspace of the span of the rows of B . Therefore, the subspaces spanned by the rows of A and the rows of B are identical, which implies that A and B have the same row rank.

Proposition 6.12 implies that the systems $Ax = 0$ and $Bx = 0$ have the same solutions. Since Ax is a linear combinations of the columns of A and Bx is a linear combinations of the columns of B , the maximum number of linearly independent columns in A is equal to the maximum number of linearly independent columns in B ; that is, A and B have the same rank. \square

Remark: The subspaces spanned by the columns of A and B can be different! However, their dimension must be the same.

Of course, we know from Proposition 4.29 that the row rank is equal to the column rank. We will see that the reduction to row echelon form provides another proof of this important fact. Let us now define precisely what is a reduced row echelon matrix.

Definition 6.1. A $m \times n$ matrix A is a *reduced row echelon matrix* iff the following conditions hold:

- (a) The first nonzero entry in every row is 1. This entry is called a *pivot*.
- (b) The first nonzero entry of row $i + 1$ is to the right of the first nonzero entry of row i .
- (c) The entries above a pivot are zero.

If a matrix satisfies the above conditions, we also say that it is in *reduced row echelon form*, for short *rref*.

Note that condition (b) implies that the entries below a pivot are also zero. For example, the matrix

$$A = \begin{pmatrix} 1 & 6 & 0 & 1 \\ 0 & 0 & 1 & 2 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

is a reduced row echelon matrix.

The following proposition shows that every matrix can be converted to a reduced row echelon form using row operations.

Proposition 6.14. *Given any $m \times n$ matrix A , there is a sequence of row operations E_1, \dots, E_k such that if $P = E_k \cdots E_1$, then $U = PA$ is a reduced row echelon matrix.*

Proof. We proceed by induction on m . If $m = 1$, then either all entries on this row are zero, so $A = 0$, or if a_j is the first nonzero entry in A , let $P = (a_j^{-1})$ (a 1×1 matrix); clearly, PA is a reduced row echelon matrix.

Let us now assume that $m \geq 2$. If $A = 0$ we are done, so let us assume that $A \neq 0$. Since $A \neq 0$, there is a leftmost column j which is nonzero, so pick any pivot $\pi = a_{ij}$ in the j th column, permute row i and row 1 if necessary, multiply the new first row by π^{-1} , and clear out the other entries in column j by subtracting suitable multiples of row 1. At the end of this process, we have a matrix A_1 that has the following shape:

$$A_1 = \begin{pmatrix} 0 & \cdots & 0 & 1 & * & \cdots & * \\ 0 & \cdots & 0 & 0 & * & \cdots & * \\ \vdots & & \vdots & \vdots & \vdots & & \vdots \\ 0 & \cdots & 0 & 0 & * & \cdots & * \end{pmatrix},$$

where $*$ stands for an arbitrary scalar, or more concisely

$$A_1 = \begin{pmatrix} 0 & 1 & B \\ 0 & 0 & D \end{pmatrix},$$

where D is a $(m-1) \times (n-j)$ matrix. If $j = n$, we are done. Otherwise, by the induction hypothesis applied to D , there is a sequence of row operations that converts D to a reduced row echelon matrix R' , and these row operations do not affect the first row of A_1 , which means that A_1 is reduced to a matrix of the form

$$R = \begin{pmatrix} 0 & 1 & B \\ 0 & 0 & R' \end{pmatrix}.$$

Because R' is a reduced row echelon matrix, the matrix R satisfies conditions (a) and (b) of the reduced row echelon form. Finally, the entries above all pivots in R' can be cleared out by subtracting suitable multiples of the rows of R' containing a pivot. The resulting matrix also satisfies condition (c), and the induction step is complete. \square

Remark: There is a **Matlab** function named **rref** that converts any matrix to its reduced row echelon form.

If A is any matrix and if R is a reduced row echelon form of A , the second part of Proposition 6.13 can be sharpened a little. Namely, *the rank of A is equal to the number of pivots in R .*

This is because the structure of a reduced row echelon matrix makes it clear that its rank is equal to the number of pivots.

Given a system of the form $Ax = b$, we can apply the reduction procedure to the augmented matrix (A, b) to obtain a reduced row echelon matrix (A', b') such that the system $A'x = b'$ has the same solutions as the original system $Ax = b$. The advantage of the reduced system $A'x = b'$ is that there is a simple test to check whether this system is solvable, and to find its solutions if it is solvable.

Indeed, if any row of the matrix A' is zero and if the corresponding entry in b' is nonzero, then it is a pivot and we have the “equation”

$$0 = 1,$$

which means that the system $A'x = b'$ has no solution. On the other hand, if there is no pivot in b' , then for every row i in which $b'_i \neq 0$, there is some column j in A' where the entry on row i is 1 (a pivot). Consequently, we can assign arbitrary values to the variable x_k if column k does not contain a pivot, and then solve for the pivot variables.

For example, if we consider the reduced row echelon matrix

$$(A', b') = \begin{pmatrix} 1 & 6 & 0 & 1 & 0 \\ 0 & 0 & 1 & 2 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix},$$

there is no solution to $A'x = b'$ because the third equation is $0 = 1$. On the other hand, the reduced system

$$(A', b') = \begin{pmatrix} 1 & 6 & 0 & 1 & 1 \\ 0 & 0 & 1 & 2 & 3 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

has solutions. We can pick the variables x_2, x_4 corresponding to nonpivot columns arbitrarily, and then solve for x_3 (using the second equation) and x_1 (using the first equation).

The above reasoning proved the following theorem:

Theorem 6.15. *Given any system $Ax = b$ where A is a $m \times n$ matrix, if the augmented matrix (A, b) is a reduced row echelon matrix, then the system $Ax = b$ has a solution iff there is no pivot in b . In that case, an arbitrary value can be assigned to the variable x_j if column j does not contain a pivot.*

Nonpivot variables are often called *free variables*.

Putting Proposition 6.14 and Theorem 6.15 together we obtain a criterion to decide whether a system $Ax = b$ has a solution: Convert the augmented system (A, b) to a row reduced echelon matrix (A', b') and check whether b' has no pivot.

Remark: When writing a program implementing row reduction, we may stop when the last column of the matrix A is reached. In this case, the test whether the system $Ax = b$ is

solvable is that the row-reduced matrix A' has no zero row of index $i > r$ such that $b'_i \neq 0$ (where r is the number of pivots, and b' is the row-reduced right-hand side).

If we have a *homogeneous system* $Ax = 0$, which means that $b = 0$, of course $x = 0$ is always a solution, but Theorem 6.15 implies that if the system $Ax = 0$ has more variables than equations, then it has some nonzero solution (we call it a *nontrivial solution*).

Proposition 6.16. *Given any homogeneous system $Ax = 0$ of m equations in n variables, if $m < n$, then there is a nonzero vector $x \in \mathbb{R}^n$ such that $Ax = 0$.*

Proof. Convert the matrix A to a reduced row echelon matrix A' . We know that $Ax = 0$ iff $A'x = 0$. If r is the number of pivots of A' , we must have $r \leq m$, so by Theorem 6.15 we may assign arbitrary values to $n - r > 0$ nonpivot variables and we get nontrivial solutions. \square

Theorem 6.15 can also be used to characterize when a square matrix is invertible. First, note the following simple but important fact:

If a square $n \times n$ matrix A is a row reduced echelon matrix, then either A is the identity or the bottom row of A is zero.

Proposition 6.17. *Let A be a square matrix of dimension n . The following conditions are equivalent:*

- (a) *The matrix A can be reduced to the identity by a sequence of elementary row operations.*
- (b) *The matrix A is a product of elementary matrices.*
- (c) *The matrix A is invertible.*
- (d) *The system of homogeneous equations $Ax = 0$ has only the trivial solution $x = 0$.*

Proof. First, we prove that (a) implies (b). If (a) can be reduced to the identity by a sequence of row operations E_1, \dots, E_p , this means that $E_p \cdots E_1 A = I$. Since each E_i is invertible, we get

$$A = E_1^{-1} \cdots E_p^{-1},$$

where each E_i^{-1} is also an elementary row operation, so (b) holds. Now if (b) holds, since elementary row operations are invertible, A is invertible, and (c) holds. If A is invertible, we already observed that the homogeneous system $Ax = 0$ has only the trivial solution $x = 0$, because from $Ax = 0$, we get $A^{-1}Ax = A^{-1}0$; that is, $x = 0$. It remains to prove that (d) implies (a), and for this we prove the contrapositive: if (a) does not hold, then (d) does not hold.

Using our basic observation about reducing square matrices, if A does not reduce to the identity, then A reduces to a row echelon matrix A' whose bottom row is zero. Say $A' = PA$, where P is a product of elementary row operations. Because the bottom row of A' is zero, the system $A'x = 0$ has at most $n - 1$ nontrivial equations, and by Proposition 6.16, this

system has a nontrivial solution x . But then, $Ax = P^{-1}A'x = 0$ with $x \neq 0$, contradicting the fact that the system $Ax = 0$ is assumed to have only the trivial solution. Therefore, (d) implies (a) and the proof is complete. \square

Proposition 6.17 yields a method for computing the inverse of an invertible matrix A : reduce A to the identity using elementary row operations, obtaining

$$E_p \cdots E_1 A = I.$$

Multiplying both sides by A^{-1} we get

$$A^{-1} = E_p \cdots E_1.$$

From a practical point of view, we can build up the product $E_p \cdots E_1$ by reducing to row echelon form the augmented $n \times 2n$ matrix (A, I_n) obtained by adding the n columns of the identity matrix to A . This is just another way of performing the Gauss–Jordan procedure.

Here is an example: let us find the inverse of the matrix

$$A = \begin{pmatrix} 5 & 4 \\ 6 & 5 \end{pmatrix}.$$

We form the 2×4 block matrix

$$(A, I) = \begin{pmatrix} 5 & 4 & 1 & 0 \\ 6 & 5 & 0 & 1 \end{pmatrix}$$

and apply elementary row operations to reduce A to the identity. For example:

$$(A, I) = \begin{pmatrix} 5 & 4 & 1 & 0 \\ 6 & 5 & 0 & 1 \end{pmatrix} \longrightarrow \begin{pmatrix} 5 & 4 & 1 & 0 \\ 1 & 1 & -1 & 1 \end{pmatrix}$$

by subtracting row 1 from row 2,

$$\begin{pmatrix} 5 & 4 & 1 & 0 \\ 1 & 1 & -1 & 1 \end{pmatrix} \longrightarrow \begin{pmatrix} 1 & 0 & 5 & -4 \\ 1 & 1 & -1 & 1 \end{pmatrix}$$

by subtracting $4 \times$ row 2 from row 1,

$$\begin{pmatrix} 1 & 0 & 5 & -4 \\ 1 & 1 & -1 & 1 \end{pmatrix} \longrightarrow \begin{pmatrix} 1 & 0 & 5 & -4 \\ 0 & 1 & -6 & 5 \end{pmatrix} = (I, A^{-1}),$$

by subtracting row 1 from row 2. Thus

$$A^{-1} = \begin{pmatrix} 5 & -4 \\ -6 & 5 \end{pmatrix}.$$

Proposition 6.17 can also be used to give an elementary proof of the fact that if a square matrix A has a left inverse B (resp. a right inverse B), so that $BA = I$ (resp. $AB = I$), then A is invertible and $A^{-1} = B$. This is an interesting exercise, try it!

For the sake of completeness, we prove that the reduced row echelon form of a matrix is unique. The neat proof given below is borrowed and adapted from W. Kahan.

Proposition 6.18. *Let A be any $m \times n$ matrix. If U and V are two reduced row echelon matrices obtained from A by applying two sequences of elementary row operations E_1, \dots, E_p and F_1, \dots, F_q , so that*

$$U = E_p \cdots E_1 A \quad \text{and} \quad V = F_q \cdots F_1 A,$$

then $U = V$ and $E_p \cdots E_1 = F_q \cdots F_1$. In other words, the reduced row echelon form of any matrix is unique.

Proof. Let

$$C = E_p \cdots E_1 F_1^{-1} \cdots F_q^{-1}$$

so that

$$U = CV \quad \text{and} \quad V = C^{-1}U.$$

We prove by induction on n that $U = V$ (and $C = I$).

Let ℓ_j denote the j th column of the identity matrix I_n , and let $u_j = U\ell_j$, $v_j = V\ell_j$, $c_j = C\ell_j$, and $a_j = A\ell_j$, be the j th column of U , V , C , and A respectively.

First, I claim that $u_j = 0$ iff $v_j = 0$, iff $a_j = 0$.

Indeed, if $v_j = 0$, then (because $U = CV$) $u_j = Cv_j = 0$, and if $u_j = 0$, then $v_j = C^{-1}u_j = 0$. Since $A = E_p \cdots E_1 U$, we also get $a_j = 0$ iff $u_j = 0$.

Therefore, we may simplify our task by striking out columns of zeros from U , V , and A , since they will have corresponding indices. We still use n to denote the number of columns of A . Observe that because U and V are reduced row echelon matrices with no zero columns, we must have $u_1 = v_1 = \ell_1$.

Claim. If U and V are reduced row echelon matrices without zero columns such that $U = CV$, for all $k \geq 1$, if $k \leq n$, then ℓ_k occurs in U iff ℓ_k occurs in V , and if ℓ_k does occur in U , then

1. ℓ_k occurs for the same index j_k in both U and V ;
2. the first j_k columns of U and V match;
3. the subsequent columns in U and V (of index $> j_k$) whose elements beyond the k th all vanish also match;
4. the first k columns of C match the first k columns of I_n .

We prove this claim by induction on k .

For the base case $k = 1$, we already know that $u_1 = v_1 = \ell_1$. We also have

$$c_1 = C\ell_1 = Cv_1 = u_1 = \ell_1.$$

If $v_j = \mu \ell_1$ for some $\mu \in \mathbb{R}$, then

$$u_j = U\ell_1 = CV\ell_1 = Cv_j = \mu C\ell_1 = \mu \ell_1 = v_j.$$

A similar argument using C^{-1} shows that if $u_j = \lambda \ell_1$, then $v_j = u_j$. Therefore, all the columns of U and V proportional to ℓ_1 match, which establishes the base case. Observe that if ℓ_2 appears in U , then it must appear in both U and V for the same index, and if not then $U = V$.

Next we now prove the induction step; this is only necessary if ℓ_{k+1} appears in both U , in which case, by (3) of the induction hypothesis, it appears in both U and V for the same index, say j_{k+1} . Thus $u_{j_{k+1}} = v_{j_{k+1}} = \ell_{k+1}$. It follows that

$$c_{k+1} = C\ell_{k+1} = Cv_{j_{k+1}} = u_{j_{k+1}} = \ell_{k+1},$$

so the first $k+1$ columns of C match the first $k+1$ columns of I_n .

Consider any subsequent column v_j (with $j > j_{k+1}$) whose elements beyond the $(k+1)$ th all vanish. Then, v_j is a linear combination of columns of V to the left of v_j , so

$$u_j = Cv_j = v_j.$$

because the first $k+1$ columns of C match the first column of I_n . Similarly, any subsequent column u_j (with $j > j_{k+1}$) whose elements beyond the $(k+1)$ th all vanish is equal to v_j . Therefore, all the subsequent columns in U and V (of index $> j_{k+1}$) whose elements beyond the $(k+1)$ th all vanish also match, which completes the induction hypothesis.

We can now prove that $U = V$ (recall that we may assume that U and V have no zero columns). We noted earlier that $u_1 = v_1 = \ell_1$, so there is a largest $k \leq n$ such that ℓ_k occurs in U . Then, the previous claim implies that all the columns of U and V match, which means that $U = V$. \square

The reduction to row echelon form also provides a method to describe the set of solutions of a linear system of the form $Ax = b$. First, we have the following simple result.

Proposition 6.19. *Let A be any $m \times n$ matrix and let $b \in \mathbb{R}^m$ be any vector. If the system $Ax = b$ has a solution, then the set Z of all solutions of this system is the set*

$$Z = x_0 + \text{Ker}(A) = \{x_0 + x \mid Ax = 0\},$$

where $x_0 \in \mathbb{R}^n$ is any solution of the system $Ax = b$, which means that $Ax_0 = b$ (x_0 is called a special solution), and where $\text{Ker}(A) = \{x \in \mathbb{R}^n \mid Ax = 0\}$, the set of solutions of the homogeneous system associated with $Ax = b$.

Proof. Assume that the system $Ax = b$ is solvable and let x_0 and x_1 be any two solutions so that $Ax_0 = b$ and $Ax_1 = b$. Subtracting the first equation from the second, we get

$$A(x_1 - x_0) = 0,$$

which means that $x_1 - x_0 \in \text{Ker}(A)$. Therefore, $Z \subseteq x_0 + \text{Ker}(A)$, where x_0 is a special solution of $Ax = b$. Conversely, if $Ax_0 = b$, then for any $z \in \text{Ker}(A)$, we have $Az = 0$, and so

$$A(x_0 + z) = Ax_0 + Az = b + 0 = b,$$

which shows that $x_0 + \text{Ker}(A) \subseteq Z$. Therefore, $Z = x_0 + \text{Ker}(A)$. \square

Given a linear system $Ax = b$, reduce the augmented matrix (A, b) to its row echelon form (A', b') . As we showed before, the system $Ax = b$ has a solution iff b' contains no pivot. Assume that this is the case. Then, if (A', b') has r pivots, which means that A' has r pivots since b' has no pivot, we know that the first r columns of I_n appear in A' .

We can permute the columns of A' and renumber the variables in x correspondingly so that the first r columns of I_n match the first r columns of A' , and then our reduced echelon matrix is of the form (R, b') with

$$R = \begin{pmatrix} I_r & F \\ 0_{m-r,r} & 0_{m-r,n-r} \end{pmatrix}$$

and

$$b' = \begin{pmatrix} d \\ 0_{m-r} \end{pmatrix},$$

where F is a $r \times (n - r)$ matrix and $d \in \mathbb{R}^r$. Note that R has $m - r$ zero rows.

Then, because

$$\begin{pmatrix} I_r & F \\ 0_{m-r,r} & 0_{m-r,n-r} \end{pmatrix} \begin{pmatrix} d \\ 0_{n-r} \end{pmatrix} = \begin{pmatrix} d \\ 0_{m-r} \end{pmatrix},$$

we see that

$$x_0 = \begin{pmatrix} d \\ 0_{n-r} \end{pmatrix}$$

is a special solution of $Rx = b'$, and thus to $Ax = b$. In other words, we get a special solution by assigning the first r components of b' to the pivot variables and setting the nonpivot variables (the *free variables*) to zero.

We can also find a basis of the kernel (nullspace) of A using F . If $x = (u, v)$ is in the kernel of A , with $u \in \mathbb{R}^r$ and $v \in \mathbb{R}^{n-r}$, then x is also in the kernel of R , which means that $Rx = 0$; that is,

$$\begin{pmatrix} I_r & F \\ 0_{m-r,r} & 0_{m-r,n-r} \end{pmatrix} \begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} u + Fv \\ 0_{m-r} \end{pmatrix} = \begin{pmatrix} 0_r \\ 0_{m-r} \end{pmatrix}.$$

Therefore, $u = -Fv$, and $\text{Ker}(A)$ consists of all vectors of the form

$$\begin{pmatrix} -Fv \\ v \end{pmatrix} = \begin{pmatrix} -F \\ I_{n-r} \end{pmatrix} v,$$

for any arbitrary $v \in \mathbb{R}^{n-r}$. It follows that the $n - r$ columns of the matrix

$$N = \begin{pmatrix} -F \\ I_{n-r} \end{pmatrix}$$

form a basis of the kernel of A . This is because N contains the identity matrix I_{n-r} as a submatrix, so the columns of N are linearly independent. In summary, if N^1, \dots, N^{n-r} are the columns of N , then the general solution of the equation $Ax = b$ is given by

$$x = \begin{pmatrix} d \\ 0_{n-r} \end{pmatrix} + x_{r+1}N^1 + \dots + x_nN^{n-r},$$

where x_{r+1}, \dots, x_n are the free variables; that is, the nonpivot variables.

In the general case where the columns corresponding to pivots are mixed with the columns corresponding to free variables, we find the special solution as follows. Let $i_1 < \dots < i_r$ be the indices of the columns corresponding to pivots. Then, assign b'_k to the pivot variable x_{i_k} for $k = 1, \dots, r$, and set all other variables to 0. To find a basis of the kernel, we form the $n - r$ vectors N^k obtained as follows. Let $j_1 < \dots < j_{n-r}$ be the indices of the columns corresponding to free variables. For every column j_k corresponding to a free variable ($1 \leq k \leq n - r$), form the vector N^k defined so that the entries $N^k_{i_1}, \dots, N^k_{i_r}$ are equal to the negatives of the first r entries in column j_k (flip the sign of these entries); let $N^k_{j_k} = 1$, and set all other entries to zero. The presence of the 1 in position j_k guarantees that N^1, \dots, N^{n-r} are linearly independent.

An illustration of the above method, consider the problem of finding a basis of the subspace V of $n \times n$ matrices $A \in M_n(\mathbb{R})$ satisfying the following properties:

1. The sum of the entries in every row has the same value (say c_1);
2. The sum of the entries in every column has the same value (say c_2).

It turns out that $c_1 = c_2$ and that the $2n - 2$ equations corresponding to the above conditions are linearly independent. We leave the proof of these facts as an interesting exercise. By the duality theorem, the dimension of the space V of matrices satisfying the above equations is $n^2 - (2n - 2)$. Let us consider the case $n = 4$. There are 6 equations, and the space V has dimension 10. The equations are

$$\begin{aligned} a_{11} + a_{12} + a_{13} + a_{14} - a_{21} - a_{22} - a_{23} - a_{24} &= 0 \\ a_{21} + a_{22} + a_{23} + a_{24} - a_{31} - a_{32} - a_{33} - a_{34} &= 0 \\ a_{31} + a_{32} + a_{33} + a_{34} - a_{41} - a_{42} - a_{43} - a_{44} &= 0 \\ a_{11} + a_{21} + a_{31} + a_{41} - a_{12} - a_{22} - a_{32} - a_{42} &= 0 \\ a_{12} + a_{22} + a_{32} + a_{42} - a_{13} - a_{23} - a_{33} - a_{43} &= 0 \\ a_{13} + a_{23} + a_{33} + a_{43} - a_{14} - a_{24} - a_{34} - a_{44} &= 0, \end{aligned}$$

and the corresponding matrix is

$$A = \begin{pmatrix} 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 \\ 1 & -1 & 0 & 0 & 1 & -1 & 0 & 0 & 1 & -1 & 0 & 0 & 1 & -1 & 0 & 0 \\ 0 & 1 & -1 & 0 & 0 & 1 & -1 & 0 & 0 & 1 & -1 & 0 & 0 & 1 & -1 & 0 \\ 0 & 0 & 1 & -1 & 0 & 0 & 1 & -1 & 0 & 0 & 1 & -1 & 0 & 0 & 1 & -1 \end{pmatrix}.$$

The result of performing the reduction to row echelon form yields the following matrix in rref:

$$U = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & -1 & -1 & -1 & 0 & -1 & -1 & -1 & 2 & 1 & 1 & 1 \\ 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & -1 & 0 & -1 & -1 \\ 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & -1 & -1 & 0 & -1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & -1 & -1 & -1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & -1 & -1 & -1 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 \end{pmatrix}$$

The list *pivlist* of indices of the pivot variables and the list *freelist* of indices of the free variables is given by

$$\textit{pivlist} = (1, 2, 3, 4, 5, 9),$$

$$\textit{freelist} = (6, 7, 8, 10, 11, 12, 13, 14, 15, 16).$$

After applying the algorithm to find a basis of the kernel of U , we find the following 16×10 matrix

$$BK = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & -2 & -1 & -1 & -1 \\ -1 & 0 & 0 & -1 & 0 & 0 & 1 & 0 & 1 & 1 \\ 0 & -1 & 0 & 0 & -1 & 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & -1 & 0 & 0 & -1 & 1 & 1 & 1 & 0 \\ -1 & -1 & -1 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & -1 & -1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}.$$

The reader should check that that in each column j of BK , the lowest 1 belongs to the row whose index is the j th element in *freelist*, and that in each column j of BK , the signs of

the entries whose indices belong to *pivlist* are the flipped signs of the 6 entries in the column U corresponding to the j th index in *freelist*. We can now read off from BK the 4×4 matrices that form a basis of V : every column of BK corresponds to a matrix whose rows have been concatenated. We get the following 10 matrices:

$$\begin{aligned}
 M_1 &= \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}, & M_2 &= \begin{pmatrix} 1 & 0 & -1 & 0 \\ -1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}, & M_3 &= \begin{pmatrix} 1 & 0 & 0 & -1 \\ -1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}, \\
 M_4 &= \begin{pmatrix} 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}, & M_5 &= \begin{pmatrix} 1 & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 \\ -1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}, & M_6 &= \begin{pmatrix} 1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 \\ -1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix}, \\
 M_7 &= \begin{pmatrix} -2 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{pmatrix}, & M_8 &= \begin{pmatrix} -1 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}, & M_9 &= \begin{pmatrix} -1 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}, \\
 M_{10} &= \begin{pmatrix} -1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.
 \end{aligned}$$

Recall that a *magic square* is a square matrix that satisfies the two conditions about the sum of the entries in each row and in each column to be the same number, and also the additional two constraints that the main descending and the main ascending diagonals add up to this common number. Furthermore, the entries are also required to be positive integers. For $n = 4$, the additional two equations are

$$\begin{aligned}
 a_{22} + a_{33} + a_{44} - a_{12} - a_{13} - a_{14} &= 0 \\
 a_{41} + a_{32} + a_{23} - a_{11} - a_{12} - a_{13} &= 0,
 \end{aligned}$$

and the 8 equations stating that a matrix is a magic square are linearly independent. Again, by running row elimination, we get a basis of the “generalized magic squares” whose entries are not restricted to be positive integers. We find a basis of 8 matrices. For $n = 3$, we find a basis of 3 matrices.

A magic square is said to be *normal* if its entries are precisely the integers $1, 2, \dots, n^2$. Then, since the sum of these entries is

$$1 + 2 + 3 + \dots + n^2 = \frac{n^2(n^2 + 1)}{2},$$

and since each row (and column) sums to the same number, this common value (the *magic sum*) is

$$\frac{n(n^2 + 1)}{2}.$$

It is easy to see that there are no normal magic squares for $n = 2$. For $n = 3$, the magic sum is 15, for $n = 4$, it is 34, and for $n = 5$, it is 65.

In the case $n = 3$, we have the additional condition that the rows and columns add up to 15, so we end up with a solution parametrized by two numbers x_1, x_2 ; namely,

$$\begin{pmatrix} x_1 + x_2 - 5 & 10 - x_2 & 10 - x_1 \\ 20 - 2x_1 - x_2 & 5 & 2x_1 + x_2 - 10 \\ x_1 & x_2 & 15 - x_1 - x_2 \end{pmatrix}.$$

Thus, in order to find a normal magic square, we have the additional inequality constraints

$$\begin{aligned} x_1 + x_2 &> 5 \\ x_1 &< 10 \\ x_2 &< 10 \\ 2x_1 + x_2 &< 20 \\ 2x_1 + x_2 &> 10 \\ x_1 &> 0 \\ x_2 &> 0 \\ x_1 + x_2 &< 15, \end{aligned}$$

and all 9 entries in the matrix must be distinct. After a tedious case analysis, we discover the remarkable fact that there is a unique normal magic square (up to rotations and reflections):

$$\begin{pmatrix} 2 & 7 & 6 \\ 9 & 5 & 1 \\ 4 & 3 & 8 \end{pmatrix}.$$

It turns out that there are 880 different normal magic squares for $n = 4$, and 275, 305, 224 normal magic squares for $n = 5$ (up to rotations and reflections). Even for $n = 4$, it takes a fair amount of work to enumerate them all! Finding the number of magic squares for $n > 5$ is an open problem!

Instead of performing elementary row operations on a matrix A , we can perform elementary columns operations, which means that we multiply A by elementary matrices on the right. As elementary row and column operations, $P(i, k)$, $E_{i,j;\beta}$, $E_{i,\lambda}$ perform the following actions:

1. As a row operation, $P(i, k)$ permutes row i and row k .

2. As a column operation, $P(i, k)$ permutes column i and column k .
3. The inverse of $P(i, k)$ is $P(i, k)$ itself.
4. As a row operation, $E_{i,j;\beta}$ adds β times row j to row i .
5. As a column operation, $E_{i,j;\beta}$ adds β times column i to column j (note the switch in the indices).
6. The inverse of $E_{i,j;\beta}$ is $E_{i,j;-\beta}$.
7. As a row operation, $E_{i,\lambda}$ multiplies row i by λ .
8. As a column operation, $E_{i,\lambda}$ multiplies column i by λ .
9. The inverse of $E_{i,\lambda}$ is $E_{i,\lambda^{-1}}$.

We can define the notion of a reduced column echelon matrix and show that every matrix can be reduced to a unique reduced column echelon form. Now, given any $m \times n$ matrix A , if we first convert A to its reduced row echelon form R , it is easy to see that we can apply elementary column operations that will reduce R to a matrix of the form

$$\begin{pmatrix} I_r & 0_{r,n-r} \\ 0_{m-r,r} & 0_{m-r,n-r} \end{pmatrix},$$

where r is the number of pivots (obtained during the row reduction). Therefore, for every $m \times n$ matrix A , there exist two sequences of elementary matrices E_1, \dots, E_p and F_1, \dots, F_q , such that

$$E_p \cdots E_1 A F_1 \cdots F_q = \begin{pmatrix} I_r & 0_{r,n-r} \\ 0_{m-r,r} & 0_{m-r,n-r} \end{pmatrix}.$$

The matrix on the right-hand side is called the *rank normal form* of A . Clearly, r is the rank of A . It is easy to see that the rank normal form also yields a proof of the fact that A and its transpose A^\top have the same rank.

6.6 Transvections and Dilatations

In this section, we characterize the linear isomorphisms of a vector space E that leave every vector in some hyperplane fixed. These maps turn out to be the linear maps that are represented in some suitable basis by elementary matrices of the form $E_{i,j;\beta}$ (transvections) or $E_{i,\lambda}$ (dilatations). Furthermore, the transvections generate the group $\mathbf{SL}(E)$, and the dilatations generate the group $\mathbf{GL}(E)$.

Let H be any hyperplane in E , and pick some (nonzero) vector $v \in E$ such that $v \notin H$, so that

$$E = H \oplus Kv.$$

Assume that $f: E \rightarrow E$ is a linear isomorphism such that $f(u) = u$ for all $u \in H$, and that f is not the identity. We have

$$f(v) = h + \alpha v, \quad \text{for some } h \in H \text{ and some } \alpha \in K,$$

with $\alpha \neq 0$, because otherwise we would have $f(v) = h = f(h)$ since $h \in H$, contradicting the injectivity of f ($v \neq h$ since $v \notin H$). For any $x \in E$, if we write

$$x = y + tv, \quad \text{for some } y \in H \text{ and some } t \in K,$$

then

$$f(x) = f(y) + f(tv) = y + tf(v) = y + th + t\alpha v,$$

and since $\alpha x = \alpha y + t\alpha v$, we get

$$\begin{aligned} f(x) - \alpha x &= (1 - \alpha)y + th \\ f(x) - x &= t(h + (\alpha - 1)v). \end{aligned}$$

Observe that if E is finite-dimensional, by picking a basis of E consisting of v and basis vectors of H , then the matrix of f is a lower triangular matrix whose diagonal entries are all 1 except the first entry which is equal to α . Therefore, $\det(f) = \alpha$.

Case 1. $\alpha \neq 1$.

We have $f(x) = \alpha x$ iff $(1 - \alpha)y + th = 0$ iff

$$y = \frac{t}{\alpha - 1}h.$$

Then, if we let $w = h + (\alpha - 1)v$, for $y = (t/(\alpha - 1))h$, we have

$$x = y + tv = \frac{t}{\alpha - 1}h + tv = \frac{t}{\alpha - 1}(h + (\alpha - 1)v) = \frac{t}{\alpha - 1}w,$$

which shows that $f(x) = \alpha x$ iff $x \in Kw$. Note that $w \notin H$, since $\alpha \neq 1$ and $v \notin H$. Therefore,

$$E = H \oplus Kw,$$

and f is the identity on H and a magnification by α on the line $D = Kw$.

Definition 6.2. Given a vector space E , for any hyperplane H in E , any nonzero vector $u \in E$ such that $u \notin H$, and any scalar $\alpha \neq 0, 1$, a linear map f such that $f(x) = x$ for all $x \in H$ and $f(x) = \alpha x$ for every $x \in D = Ku$ is called a *dilatation of hyperplane H , direction D , and scale factor α* .

If π_H and π_D are the projections of E onto H and D , then we have

$$f(x) = \pi_H(x) + \alpha\pi_D(x).$$

The inverse of f is given by

$$f^{-1}(x) = \pi_H(x) + \alpha^{-1}\pi_D(x).$$

When $\alpha = -1$, we have $f^2 = \text{id}$, and f is a symmetry about the hyperplane H in the direction D .

Case 2. $\alpha = 1$.

In this case,

$$f(x) - x = th,$$

that is, $f(x) - x \in Kh$ for all $x \in E$. Assume that the hyperplane H is given as the kernel of some linear form φ , and let $a = \varphi(v)$. We have $a \neq 0$, since $v \notin H$. For any $x \in E$, we have

$$\varphi(x - a^{-1}\varphi(x)v) = \varphi(x) - a^{-1}\varphi(x)\varphi(v) = \varphi(x) - \varphi(x) = 0,$$

which shows that $x - a^{-1}\varphi(x)v \in H$ for all $x \in E$. Since every vector in H is fixed by f , we get

$$\begin{aligned} x - a^{-1}\varphi(x)v &= f(x - a^{-1}\varphi(x)v) \\ &= f(x) - a^{-1}\varphi(x)f(v), \end{aligned}$$

so

$$f(x) = x + \varphi(x)(f(a^{-1}v) - a^{-1}v).$$

Since $f(z) - z \in Kh$ for all $z \in E$, we conclude that $u = f(a^{-1}v) - a^{-1}v = \beta h$ for some $\beta \in K$, so $\varphi(u) = 0$, and we have

$$f(x) = x + \varphi(x)u, \quad \varphi(u) = 0. \quad (*)$$

A linear map defined as above is denoted by $\tau_{\varphi,u}$.

Conversely for any linear map $f = \tau_{\varphi,u}$ given by equation (*), where φ is a nonzero linear form and u is some vector $u \in E$ such that $\varphi(u) = 0$, if $u = 0$ then f is the identity, so assume that $u \neq 0$. If so, we have $f(x) = x$ iff $\varphi(x) = 0$, that is, iff $x \in H$. We also claim that the inverse of f is obtained by changing u to $-u$. Actually, we check the slightly more general fact that

$$\tau_{\varphi,u} \circ \tau_{\varphi,v} = \tau_{\varphi,u+v}.$$

Indeed, using the fact that $\varphi(v) = 0$, we have

$$\begin{aligned} \tau_{\varphi,u}(\tau_{\varphi,v}(x)) &= \tau_{\varphi,v}(x) + \varphi(\tau_{\varphi,v}(v))u \\ &= \tau_{\varphi,v}(x) + (\varphi(x) + \varphi(v)\varphi(v))u \\ &= \tau_{\varphi,v}(x) + \varphi(x)u \\ &= x + \varphi(x)v + \varphi(x)u \\ &= x + \varphi(x)(u + v). \end{aligned}$$

For $v = -u$, we have $\tau_{\varphi, u+v} = \varphi_{\varphi, 0} = \text{id}$, so $\tau_{\varphi, u}^{-1} = \tau_{\varphi, -u}$, as claimed.

Therefore, we proved that every linear isomorphism of E that leaves every vector in some hyperplane H fixed and has the property that $f(x) - x \in H$ for all $x \in E$ is given by a map $\tau_{\varphi, u}$ as defined by equation (*), where φ is some nonzero linear form defining H and u is some vector in H . We have $\tau_{\varphi, u} = \text{id}$ iff $u = 0$.

Definition 6.3. Given any hyperplane H in E , for any nonzero linear form $\varphi \in E^*$ defining H (which means that $H = \text{Ker}(\varphi)$) and any nonzero vector $u \in H$, the linear map $\tau_{\varphi, u}$ given by

$$\tau_{\varphi, u}(x) = x + \varphi(x)u, \quad \varphi(u) = 0,$$

for all $x \in E$ is called a *transvection of hyperplane H and direction u* . The map $\tau_{\varphi, u}$ leaves every vector in H fixed, and $f(x) - x \in Ku$ for all $x \in E$.

The above arguments show the following result.

Proposition 6.20. Let $f: E \rightarrow E$ be a bijective linear map and assume that $f \neq \text{id}$ and that $f(x) = x$ for all $x \in H$, where H is some hyperplane in E . If there is some nonzero vector $u \in E$ such that $u \notin H$ and $f(u) - u \in H$, then f is a transvection of hyperplane H ; otherwise, f is a dilatation of hyperplane H .

Proof. Using the notation as above, for some $v \notin H$, we have $f(v) = h + \alpha v$ with $\alpha \neq 0$, and write $u = y + tv$ with $y \in H$ and $t \neq 0$ since $u \notin H$. If $f(u) - u \in H$, from

$$f(u) - u = t(h + (\alpha - 1)v),$$

we get $(\alpha - 1)v \in H$, and since $v \notin H$, we must have $\alpha = 1$, and we proved that f is a transvection. Otherwise, $\alpha \neq 0, 1$, and we proved that f is a dilatation. \square

If E is finite-dimensional, then $\alpha = \det(f)$, so we also have the following result.

Proposition 6.21. Let $f: E \rightarrow E$ be a bijective linear map of a finite-dimensional vector space E and assume that $f \neq \text{id}$ and that $f(x) = x$ for all $x \in H$, where H is some hyperplane in E . If $\det(f) = 1$, then f is a transvection of hyperplane H ; otherwise, f is a dilatation of hyperplane H .

Suppose that f is a dilatation of hyperplane H and direction u , and say $\det(f) = \alpha \neq 0, 1$. Pick a basis (u, e_2, \dots, e_n) of E where (e_2, \dots, e_n) is a basis of H . Then, the matrix of f is of the form

$$\begin{pmatrix} \alpha & 0 & \cdots & 0 \\ 0 & 1 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix},$$

which is an elementary matrix of the form $E_{1,\alpha}$. Conversely, it is clear that every elementary matrix of the form $E_{i,\alpha}$ with $\alpha \neq 0, 1$ is a dilatation.

Now, assume that f is a transvection of hyperplane H and direction $u \in H$. Pick some $v \notin H$, and pick some basis (u, e_3, \dots, e_n) of H , so that (v, u, e_3, \dots, e_n) is a basis of E . Since $f(v) - v \in Ku$, the matrix of f is of the form

$$\begin{pmatrix} 1 & 0 & \cdots & 0 \\ \alpha & 1 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix},$$

which is an elementary matrix of the form $E_{2,1;\alpha}$. Conversely, it is clear that every elementary matrix of the form $E_{i,j;\alpha}$ ($\alpha \neq 0$) is a transvection.

The following proposition is an interesting exercise that requires good mastery of the elementary row operations $E_{i,j;\beta}$.

Proposition 6.22. *Given any invertible $n \times n$ matrix A , there is a matrix S such that*

$$SA = \begin{pmatrix} I_{n-1} & 0 \\ 0 & \alpha \end{pmatrix} = E_{n,\alpha},$$

with $\alpha = \det(A)$, and where S is a product of elementary matrices of the form $E_{i,j;\beta}$; that is, S is a composition of transvections.

Surprisingly, every transvection is the composition of two dilatations!

Proposition 6.23. *If the field K is not of characteristic 2, then every transvection f of hyperplane H can be written as $f = d_2 \circ d_1$, where d_1, d_2 are dilatations of hyperplane H , where the direction of d_1 can be chosen arbitrarily.*

Proof. Pick some dilatation d_1 of hyperplane H and scale factor $\alpha \neq 0, 1$. Then, $d_2 = f \circ d_1^{-1}$ leaves every vector in H fixed, and $\det(d_2) = \alpha^{-1} \neq 1$. By Proposition 6.21, the linear map d_2 is a dilatation of hyperplane H , and we have $f = d_2 \circ d_1$, as claimed. \square

Observe that in Proposition 6.23, we can pick $\alpha = -1$; that is, every transvection of hyperplane H is the compositions of two symmetries about the hyperplane H , one of which can be picked arbitrarily.

Remark: Proposition 6.23 holds as long as $K \neq \{0, 1\}$.

The following important result is now obtained.

Theorem 6.24. *Let E be any finite-dimensional vector space over a field K of characteristic not equal to 2. Then, the group $\mathbf{SL}(E)$ is generated by the transvections, and the group $\mathbf{GL}(E)$ is generated by the dilatations.*

Proof. Consider any $f \in \mathbf{SL}(E)$, and let A be its matrix in any basis. By Proposition 6.22, there is a matrix S such that

$$SA = \begin{pmatrix} I_{n-1} & 0 \\ 0 & \alpha \end{pmatrix} = E_{n,\alpha},$$

with $\alpha = \det(A)$, and where S is a product of elementary matrices of the form $E_{i,j;\beta}$. Since $\det(A) = 1$, we have $\alpha = 1$, and the result is proved. Otherwise, $E_{n,\alpha}$ is a dilatation, S is a product of transvections, and by Proposition 6.23, every transvection is the composition of two dilatations, so the second result is also proved. \square

We conclude this section by proving that any two transvections are conjugate in $\mathbf{GL}(E)$. Let $\tau_{\varphi,u}$ ($u \neq 0$) be a transvection and let $g \in \mathbf{GL}(E)$ be any invertible linear map. We have

$$\begin{aligned} (g \circ \tau_{\varphi,u} \circ g^{-1})(x) &= g(g^{-1}(x) + \varphi(g^{-1}(x))u) \\ &= x + \varphi(g^{-1}(x))g(u). \end{aligned}$$

Let us find the hyperplane determined by the linear form $x \mapsto \varphi(g^{-1}(x))$. This is the set of vectors $x \in E$ such that $\varphi(g^{-1}(x)) = 0$, which holds iff $g^{-1}(x) \in H$ iff $x \in g(H)$. Therefore, $\text{Ker}(\varphi \circ g^{-1}) = g(H) = H'$, and we have $g(u) \in g(H) = H'$, so $g \circ \tau_{\varphi,u} \circ g^{-1}$ is the transvection of hyperplane $H' = g(H)$ and direction $u' = g(u)$ (with $u' \in H'$).

Conversely, let $\tau_{\psi,u'}$ be some transvection ($u' \neq 0$). Pick some vector v, v' such that $\varphi(v) = \psi(v') = 1$, so that

$$E = H \oplus Kv = H' \oplus v'.$$

There is a linear map $g \in \mathbf{GL}(E)$ such that $g(u) = u'$, $g(v) = v'$, and $g(H) = H'$. To define g , pick a basis $(v, u, e_2, \dots, e_{n-1})$ where (u, e_2, \dots, e_{n-1}) is a basis of H and pick a basis $(v', u', e'_2, \dots, e'_{n-1})$ where $(u', e'_2, \dots, e'_{n-1})$ is a basis of H' ; then g is defined so that $g(v) = v'$, $g(u) = u'$, and $g(e_i) = g(e'_i)$, for $i = 2, \dots, n-1$. If $n = 2$, then e_i and e'_i are missing. Then, we have

$$(g \circ \tau_{\varphi,u} \circ g^{-1})(x) = x + \varphi(g^{-1}(x))u'.$$

Now, $\varphi \circ g^{-1}$ also determines the hyperplane $H' = g(H)$, so we have $\varphi \circ g^{-1} = \lambda\psi$ for some nonzero λ in K . Since $v' = g(v)$, we get

$$\varphi(v) = \varphi \circ g^{-1}(v') = \lambda\psi(v'),$$

and since $\varphi(v) = \psi(v') = 1$, we must have $\lambda = 1$. It follows that

$$(g \circ \tau_{\varphi,u} \circ g^{-1})(x) = x + \psi(x)u' = \tau_{\psi,u'}(x).$$

In summary, we proved almost all parts the following result.

Proposition 6.25. *Let E be any finite-dimensional vector space. For every transvection $\tau_{\varphi,u}$ ($u \neq 0$) and every linear map $g \in \mathbf{GL}(E)$, the map $g \circ \tau_{\varphi,u} \circ g^{-1}$ is the transvection of hyperplane $g(H)$ and direction $g(u)$ (that is, $g \circ \tau_{\varphi,u} \circ g^{-1} = \tau_{\varphi \circ g^{-1}, g(u)}$). For every other transvection $\tau_{\psi,u'}$ ($u' \neq 0$), there is some $g \in \mathbf{GL}(E)$ such $\tau_{\psi,u'} = g \circ \tau_{\varphi,u} \circ g^{-1}$; in other words any two transvections ($\neq \text{id}$) are conjugate in $\mathbf{GL}(E)$. Moreover, if $n \geq 3$, then the linear isomorphism g as above can be chosen so that $g \in \mathbf{SL}(E)$.*

Proof. We just need to prove that if $n \geq 3$, then for any two transvections $\tau_{\varphi,u}$ and $\tau_{\psi,u'}$ ($u, u' \neq 0$), there is some $g \in \mathbf{SL}(E)$ such that $\tau_{\psi,u'} = g \circ \tau_{\varphi,u} \circ g^{-1}$. As before, we pick a basis $(v, u, e_2, \dots, e_{n-1})$ where (u, e_2, \dots, e_{n-1}) is a basis of H , we pick a basis $(v', u', e'_2, \dots, e'_{n-1})$ where $(u', e'_2, \dots, e'_{n-1})$ is a basis of H' , and we define g as the unique linear map such that $g(v) = v'$, $g(u) = u'$, and $g(e_i) = e'_i$, for $i = 1, \dots, n-1$. But, in this case, both H and $H' = g(H)$ have dimension at least 2, so in any basis of H' including u' , there is some basis vector e'_2 independent of u' , and we can rescale e'_2 in such a way that the matrix of g over the two bases has determinant $+1$. \square

6.7 Summary

The main concepts and results of this chapter are listed below:

- One does not solve (large) linear systems by computing determinants.
- *Upper-triangular* (*lower-triangular*) matrices.
- Solving by *back-substitution* (*forward-substitution*).
- *Gaussian elimination*.
- Permuting rows.
- The *pivot* of an elimination step; *pivoting*.
- *Transposition matrix*; *elementary matrix*.
- The *Gaussian elimination theorem* (Theorem 6.1).
- *Gauss-Jordan factorization*.
- *LU-factorization*; Necessary and sufficient condition for the existence of an *LU-factorization* (Proposition 6.2).
- *LDU-factorization*.
- “*PA = LU* theorem” (Theorem 6.5).
- *LDL^T-factorization* of a symmetric matrix.

- Avoiding small pivots: *partial pivoting*; *complete pivoting*.
- Gaussian elimination of tridiagonal matrices.
- *LU*-factorization of tridiagonal matrices.
- *Symmetric positive definite* matrices (SPD matrices).
- *Cholesky factorization* (Theorem 6.10).
- Criteria for a symmetric matrix to be positive definite; *Sylvester's criterion*.
- *Reduced row echelon form*.
- Reduction of a rectangular matrix to its row echelon form.
- Using the reduction to row echelon form to decide whether a system $Ax = b$ is solvable, and to find its solutions, using a *special* solution and a basis of the *homogeneous system* $Ax = 0$.
- *Magic squares*.
- *transvections and dilatations*.

Chapter 7

Vector Norms and Matrix Norms

7.1 Normed Vector Spaces

In order to define how close two vectors or two matrices are, and in order to define the convergence of sequences of vectors or matrices, we can use the notion of a norm. Recall that $\mathbb{R}_+ = \{x \in \mathbb{R} \mid x \geq 0\}$. Also recall that if $z = a + ib \in \mathbb{C}$ is a complex number, with $a, b \in \mathbb{R}$, then $\bar{z} = a - ib$ and $|z| = \sqrt{z\bar{z}} = \sqrt{a^2 + b^2}$ ($|z|$ is the *modulus* of z).

Definition 7.1. Let E be a vector space over a field K , where K is either the field \mathbb{R} of reals, or the field \mathbb{C} of complex numbers. A *norm* on E is a function $\|\cdot\|: E \rightarrow \mathbb{R}_+$, assigning a nonnegative real number $\|u\|$ to any vector $u \in E$, and satisfying the following conditions for all $x, y, z \in E$:

(N1) $\|x\| \geq 0$, and $\|x\| = 0$ iff $x = 0$. (positivity)

(N2) $\|\lambda x\| = |\lambda| \|x\|$. (homogeneity (or scaling))

(N3) $\|x + y\| \leq \|x\| + \|y\|$. (triangle inequality)

A vector space E together with a norm $\|\cdot\|$ is called a *normed vector space*.

By (N2), setting $\lambda = -1$, we obtain

$$\|-x\| = \|(-1)x\| = |-1| \|x\| = \|x\|;$$

that is, $\|-x\| = \|x\|$. From (N3), we have

$$\|x\| = \|x - y + y\| \leq \|x - y\| + \|y\|,$$

which implies that

$$\|x\| - \|y\| \leq \|x - y\|.$$

By exchanging x and y and using the fact that by (N2),

$$\|y - x\| = \|-(x - y)\| = \|x - y\|,$$

we also have

$$\|y\| - \|x\| \leq \|x - y\|.$$

Therefore,

$$|\|x\| - \|y\|| \leq \|x - y\|, \quad \text{for all } x, y \in E. \quad (*)$$

Observe that setting $\lambda = 0$ in (N2), we deduce that $\|0\| = 0$ without assuming (N1). Then, by setting $y = 0$ in (*), we obtain

$$|\|x\|| \leq \|x\|, \quad \text{for all } x \in E.$$

Therefore, the condition $\|x\| \geq 0$ in (N1) follows from (N2) and (N3), and (N1) can be replaced by the weaker condition

(N1') For all $x \in E$, if $\|x\| = 0$ then $x = 0$,

A function $\|\cdot\| : E \rightarrow \mathbb{R}$ satisfying axioms (N2) and (N3) is called a *seminorm*. From the above discussion, a seminorm also has the properties

$$\|x\| \geq 0 \text{ for all } x \in E, \text{ and } \|0\| = 0.$$

However, there may be nonzero vectors $x \in E$ such that $\|x\| = 0$. Let us give some examples of normed vector spaces.

Example 7.1.

1. Let $E = \mathbb{R}$, and $\|x\| = |x|$, the absolute value of x .
2. Let $E = \mathbb{C}$, and $\|z\| = |z|$, the modulus of z .
3. Let $E = \mathbb{R}^n$ (or $E = \mathbb{C}^n$). There are three standard norms. For every $(x_1, \dots, x_n) \in E$, we have the norm $\|x\|_1$, defined such that,

$$\|x\|_1 = |x_1| + \dots + |x_n|,$$

we have the *Euclidean norm* $\|x\|_2$, defined such that,

$$\|x\|_2 = (|x_1|^2 + \dots + |x_n|^2)^{\frac{1}{2}},$$

and the *sup-norm* $\|x\|_\infty$, defined such that,

$$\|x\|_\infty = \max\{|x_i| \mid 1 \leq i \leq n\}.$$

More generally, we define the ℓ_p -norm (for $p \geq 1$) by

$$\|x\|_p = (|x_1|^p + \dots + |x_n|^p)^{1/p}.$$

There are other norms besides the ℓ_p -norms; we urge the reader to find such norms.

Some work is required to show the triangle inequality for the ℓ_p -norm.

Proposition 7.1. *If E is a finite-dimensional vector space over \mathbb{R} or \mathbb{C} , for every real number $p \geq 1$, the ℓ_p -norm is indeed a norm.*

Proof. The cases $p = 1$ and $p = \infty$ are easy and left to the reader. If $p > 1$, then let $q > 1$ such that

$$\frac{1}{p} + \frac{1}{q} = 1.$$

We will make use of the following fact: for all $\alpha, \beta \in \mathbb{R}$, if $\alpha, \beta \geq 0$, then

$$\alpha\beta \leq \frac{\alpha^p}{p} + \frac{\beta^q}{q}. \quad (*)$$

To prove the above inequality, we use the fact that the exponential function $t \mapsto e^t$ satisfies the following convexity inequality:

$$e^{\theta x + (1-\theta)y} \leq \theta e^x + (1-\theta)e^y,$$

for all $x, y \in \mathbb{R}$ and all θ with $0 \leq \theta \leq 1$.

Since the case $\alpha\beta = 0$ is trivial, let us assume that $\alpha > 0$ and $\beta > 0$. If we replace θ by $1/p$, x by $p \log \alpha$ and y by $q \log \beta$, then we get

$$e^{\frac{1}{p}p \log \alpha + \frac{1}{q}q \log \beta} \leq \frac{1}{p}e^{p \log \alpha} + \frac{1}{q}e^{q \log \beta},$$

which simplifies to

$$\alpha\beta \leq \frac{\alpha^p}{p} + \frac{\beta^q}{q},$$

as claimed.

We will now prove that for any two vectors $u, v \in E$, we have

$$\sum_{i=1}^n |u_i v_i| \leq \|u\|_p \|v\|_q. \quad (**)$$

Since the above is trivial if $u = 0$ or $v = 0$, let us assume that $u \neq 0$ and $v \neq 0$. Then, the inequality $(*)$ with $\alpha = |u_i|/\|u\|_p$ and $\beta = |v_i|/\|v\|_q$ yields

$$\frac{|u_i v_i|}{\|u\|_p \|v\|_q} \leq \frac{|u_i|^p}{p \|u\|_p^p} + \frac{|v_i|^q}{q \|v\|_q^q},$$

for $i = 1, \dots, n$, and by summing up these inequalities, we get

$$\sum_{i=1}^n |u_i v_i| \leq \|u\|_p \|v\|_q,$$

as claimed. To finish the proof, we simply have to prove that property (N3) holds, since (N1) and (N2) are clear. Now, for $i = 1, \dots, n$, we can write

$$(|u_i| + |v_i|)^p = |u_i|(|u_i| + |v_i|)^{p-1} + |v_i|(|u_i| + |v_i|)^{p-1},$$

so that by summing up these equations we get

$$\sum_{i=1}^n (|u_i| + |v_i|)^p = \sum_{i=1}^n |u_i|(|u_i| + |v_i|)^{p-1} + \sum_{i=1}^n |v_i|(|u_i| + |v_i|)^{p-1},$$

and using the inequality (**), we get

$$\sum_{i=1}^n (|u_i| + |v_i|)^p \leq (\|u\|_p + \|v\|_p) \left(\sum_{i=1}^n (|u_i| + |v_i|)^{(p-1)q} \right)^{1/q}.$$

However, $1/p + 1/q = 1$ implies $pq = p + q$, that is, $(p-1)q = p$, so we have

$$\sum_{i=1}^n (|u_i| + |v_i|)^p \leq (\|u\|_p + \|v\|_p) \left(\sum_{i=1}^n (|u_i| + |v_i|)^p \right)^{1/q},$$

which yields

$$\left(\sum_{i=1}^n (|u_i| + |v_i|)^p \right)^{1/p} \leq \|u\|_p + \|v\|_p.$$

Since $|u_i + v_i| \leq |u_i| + |v_i|$, the above implies the triangle inequality $\|u + v\|_p \leq \|u\|_p + \|v\|_p$, as claimed. \square

For $p > 1$ and $1/p + 1/q = 1$, the inequality

$$\sum_{i=1}^n |u_i v_i| \leq \left(\sum_{i=1}^n |u_i|^p \right)^{1/p} \left(\sum_{i=1}^n |v_i|^q \right)^{1/q}$$

is known as *Hölder's inequality*. For $p = 2$, it is the *Cauchy-Schwarz inequality*.

Actually, if we define the *Hermitian inner product* $\langle -, - \rangle$ on \mathbb{C}^n by

$$\langle u, v \rangle = \sum_{i=1}^n u_i \bar{v}_i,$$

where $u = (u_1, \dots, u_n)$ and $v = (v_1, \dots, v_n)$, then

$$|\langle u, v \rangle| \leq \sum_{i=1}^n |u_i \bar{v}_i| = \sum_{i=1}^n |u_i v_i|,$$

so Hölder's inequality implies the inequality

$$|\langle u, v \rangle| \leq \|u\|_p \|v\|_q$$

also called *Hölder's inequality*, which, for $p = 2$ is the standard Cauchy–Schwarz inequality. The triangle inequality for the ℓ_p -norm,

$$\left(\sum_{i=1}^n (|u_i + v_i|)^p \right)^{1/p} \leq \left(\sum_{i=1}^n |u_i|^p \right)^{1/p} + \left(\sum_{i=1}^n |v_i|^p \right)^{1/p},$$

is known as *Minkowski's inequality*.

When we restrict the Hermitian inner product to real vectors, $u, v \in \mathbb{R}^n$, we get the *Euclidean inner product*

$$\langle u, v \rangle = \sum_{i=1}^n u_i v_i.$$

It is very useful to observe that if we represent (as usual) $u = (u_1, \dots, u_n)$ and $v = (v_1, \dots, v_n)$ (in \mathbb{R}^n) by column vectors, then their Euclidean inner product is given by

$$\langle u, v \rangle = u^\top v = v^\top u,$$

and when $u, v \in \mathbb{C}^n$, their Hermitian inner product is given by

$$\langle u, v \rangle = v^* u = \overline{u^* v}.$$

In particular, when $u = v$, in the complex case we get

$$\|u\|_2^2 = u^* u,$$

and in the real case, this becomes

$$\|u\|_2^2 = u^\top u.$$

As convenient as these notations are, we still recommend that you do not abuse them; the notation $\langle u, v \rangle$ is more intrinsic and still “works” when our vector space is infinite dimensional.

The following proposition is easy to show.

Proposition 7.2. *The following inequalities hold for all $x \in \mathbb{R}^n$ (or $x \in \mathbb{C}^n$):*

$$\begin{aligned} \|x\|_\infty &\leq \|x\|_1 \leq n\|x\|_\infty, \\ \|x\|_\infty &\leq \|x\|_2 \leq \sqrt{n}\|x\|_\infty, \\ \|x\|_2 &\leq \|x\|_1 \leq \sqrt{n}\|x\|_2. \end{aligned}$$

Proposition 7.2 is actually a special case of a very important result: in a finite-dimensional vector space, any two norms are equivalent.

Definition 7.2. Given any (real or complex) vector space E , two norms $\|\cdot\|_a$ and $\|\cdot\|_b$ are *equivalent* iff there exists some positive reals $C_1, C_2 > 0$, such that

$$\|u\|_a \leq C_1 \|u\|_b \quad \text{and} \quad \|u\|_b \leq C_2 \|u\|_a, \quad \text{for all } u \in E.$$

Given any norm $\|\cdot\|$ on a vector space of dimension n , for any basis (e_1, \dots, e_n) of E , observe that for any vector $x = x_1 e_1 + \dots + x_n e_n$, we have

$$\|x\| = \|x_1 e_1 + \dots + x_n e_n\| \leq |x_1| \|e_1\| + \dots + |x_n| \|e_n\| \leq C(|x_1| + \dots + |x_n|) = C \|x\|_1,$$

with $C = \max_{1 \leq i \leq n} \|e_i\|$ and

$$\|x\|_1 = \|x_1 e_1 + \dots + x_n e_n\| = |x_1| + \dots + |x_n|.$$

The above implies that

$$|\|u\| - \|v\|| \leq \|u - v\| \leq C \|u - v\|_1,$$

which means that the map $u \mapsto \|u\|$ is *continuous* with respect to the norm $\|\cdot\|_1$.

Let S_1^{n-1} be the unit sphere with respect to the norm $\|\cdot\|_1$, namely

$$S_1^{n-1} = \{x \in E \mid \|x\|_1 = 1\}.$$

Now, S_1^{n-1} is a closed and bounded subset of a finite-dimensional vector space, so by Heine–Borel (or equivalently, by Bolzano–Weierstrass), S_1^{n-1} is compact. On the other hand, it is a well known result of analysis that any continuous real-valued function on a nonempty compact set has a minimum and a maximum, and that they are achieved. Using these facts, we can prove the following important theorem:

Theorem 7.3. *If E is any real or complex vector space of finite dimension, then any two norms on E are equivalent.*

Proof. It is enough to prove that any norm $\|\cdot\|$ is equivalent to the 1-norm. We already proved that the function $x \mapsto \|x\|$ is continuous with respect to the norm $\|\cdot\|_1$ and we observed that the unit sphere S_1^{n-1} is compact. Now, we just recalled that because the function $f: x \mapsto \|x\|$ is continuous and because S_1^{n-1} is compact, the function f has a minimum m and a maximum M , and because $\|x\|$ is never zero on S_1^{n-1} , we must have $m > 0$. Consequently, we just proved that if $\|x\|_1 = 1$, then

$$0 < m \leq \|x\| \leq M,$$

so for any $x \in E$ with $x \neq 0$, we get

$$m \leq \|x\| / \|x\|_1 \leq M,$$

which implies

$$m \|x\|_1 \leq \|x\| \leq M \|x\|_1.$$

Since the above inequality holds trivially if $x = 0$, we just proved that $\|\cdot\|$ and $\|\cdot\|_1$ are equivalent, as claimed. \square

Next, we will consider norms on matrices.

7.2 Matrix Norms

For simplicity of exposition, we will consider the vector spaces $M_n(\mathbb{R})$ and $M_n(\mathbb{C})$ of square $n \times n$ matrices. Most results also hold for the spaces $M_{m,n}(\mathbb{R})$ and $M_{m,n}(\mathbb{C})$ of rectangular $m \times n$ matrices. Since $n \times n$ matrices can be multiplied, the idea behind matrix norms is that they should behave “well” with respect to matrix multiplication.

Definition 7.3. A *matrix norm* $\|\cdot\|$ on the space of square $n \times n$ matrices in $M_n(K)$, with $K = \mathbb{R}$ or $K = \mathbb{C}$, is a norm on the vector space $M_n(K)$, with the additional property called *submultiplicativity* that

$$\|AB\| \leq \|A\| \|B\|,$$

for all $A, B \in M_n(K)$. A norm on matrices satisfying the above property is often called a *submultiplicative* matrix norm.

Since $I^2 = I$, from $\|I\| = \|I^2\| \leq \|I\|^2$, we get $\|I\| \geq 1$, for every matrix norm.

Before giving examples of matrix norms, we need to review some basic definitions about matrices. Given any matrix $A = (a_{ij}) \in M_{m,n}(\mathbb{C})$, the *conjugate* \bar{A} of A is the matrix such that

$$\bar{A}_{ij} = \bar{a}_{ij}, \quad 1 \leq i \leq m, 1 \leq j \leq n.$$

The *transpose* of A is the $n \times m$ matrix A^\top such that

$$A_{ij}^\top = a_{ji}, \quad 1 \leq i \leq m, 1 \leq j \leq n.$$

The *adjoint* of A is the $n \times m$ matrix A^* such that

$$A^* = \overline{(A^\top)} = (\bar{A})^\top.$$

When A is a real matrix, $A^* = A^\top$. A matrix $A \in M_n(\mathbb{C})$ is *Hermitian* if

$$A^* = A.$$

If A is a real matrix ($A \in M_n(\mathbb{R})$), we say that A is *symmetric* if

$$A^\top = A.$$

A matrix $A \in M_n(\mathbb{C})$ is *normal* if

$$AA^* = A^*A,$$

and if A is a real matrix, it is *normal* if

$$AA^\top = A^\top A.$$

A matrix $U \in M_n(\mathbb{C})$ is *unitary* if

$$UU^* = U^*U = I.$$

A real matrix $Q \in M_n(\mathbb{R})$ is *orthogonal* if

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

Given any matrix $A = (a_{ij}) \in M_n(\mathbb{C})$, the *trace* $\text{tr}(A)$ of A is the sum of its diagonal elements

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

It is easy to show that the trace is a linear map, so that

$$\text{tr}(\lambda A) = \lambda \text{tr}(A)$$

and

$$\text{tr}(A + B) = \text{tr}(A) + \text{tr}(B).$$

Moreover, if A is an $m \times n$ matrix and B is an $n \times m$ matrix, it is not hard to show that

$$\text{tr}(AB) = \text{tr}(BA).$$

We also review eigenvalues and eigenvectors. We content ourselves with definition involving matrices. A more general treatment will be given later on (see Chapter 8).

Definition 7.4. Given any square matrix $A \in M_n(\mathbb{C})$, a complex number $\lambda \in \mathbb{C}$ is an *eigenvalue* of A if there is some *nonzero* vector $u \in \mathbb{C}^n$, such that

$$Au = \lambda u.$$

If λ is an eigenvalue of A , then the *nonzero* vectors $u \in \mathbb{C}^n$ such that $Au = \lambda u$ are called *eigenvectors of A associated with λ* ; together with the zero vector, these eigenvectors form a subspace of \mathbb{C}^n denoted by $E_\lambda(A)$, and called the *eigenspace associated with λ* .

Remark: Note that Definition 7.4 *requires an eigenvector to be nonzero*. A somewhat unfortunate consequence of this requirement is that the set of eigenvectors is *not* a subspace, since the zero vector is missing! On the positive side, whenever eigenvectors are involved, there is no need to say that they are nonzero. The fact that eigenvectors are nonzero is implicitly used in all the arguments involving them, so it seems safer (but perhaps not as elegant) to stipulate that eigenvectors should be nonzero.

If A is a square real matrix $A \in M_n(\mathbb{R})$, then we restrict Definition 7.4 to real eigenvalues $\lambda \in \mathbb{R}$ and real eigenvectors. However, it should be noted that although every complex matrix always has at least some complex eigenvalue, a real matrix may not have any real eigenvalues. For example, the matrix

$$A = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}$$

has the complex eigenvalues i and $-i$, but no real eigenvalues. Thus, typically, even for real matrices, we consider complex eigenvalues.

Observe that $\lambda \in \mathbb{C}$ is an eigenvalue of A
 iff $Au = \lambda u$ for some nonzero vector $u \in \mathbb{C}^n$
 iff $(\lambda I - A)u = 0$
 iff the matrix $\lambda I - A$ defines a linear map which has a nonzero kernel, that is,
 iff $\lambda I - A$ not invertible.

However, from Proposition 5.10, $\lambda I - A$ is not invertible iff

$$\det(\lambda I - A) = 0.$$

Now, $\det(\lambda I - A)$ is a polynomial of degree n in the indeterminate λ , in fact, of the form

$$\lambda^n - \operatorname{tr}(A)\lambda^{n-1} + \cdots + (-1)^n \det(A).$$

Thus, we see that the eigenvalues of A are the zeros (also called *roots*) of the above polynomial. Since every complex polynomial of degree n has exactly n roots, counted with their multiplicity, we have the following definition:

Definition 7.5. Given any square $n \times n$ matrix $A \in M_n(\mathbb{C})$, the polynomial

$$\det(\lambda I - A) = \lambda^n - \operatorname{tr}(A)\lambda^{n-1} + \cdots + (-1)^n \det(A)$$

is called the *characteristic polynomial* of A . The n (not necessarily distinct) roots $\lambda_1, \dots, \lambda_n$ of the characteristic polynomial are all the *eigenvalues* of A and constitute the *spectrum* of A . We let

$$\rho(A) = \max_{1 \leq i \leq n} |\lambda_i|$$

be the largest modulus of the eigenvalues of A , called the *spectral radius* of A .

Proposition 7.4. For any matrix norm $\|\cdot\|$ on $M_n(\mathbb{C})$ and for any square $n \times n$ matrix $A \in M_n(\mathbb{C})$, we have

$$\rho(A) \leq \|A\|.$$

Proof. Let λ be some eigenvalue of A for which $|\lambda|$ is maximum, that is, such that $|\lambda| = \rho(A)$. If $u (\neq 0)$ is any eigenvector associated with λ and if U is the $n \times n$ matrix whose columns are all u , then $Au = \lambda u$ implies

$$AU = \lambda U,$$

and since

$$|\lambda| \|U\| = \|\lambda U\| = \|AU\| \leq \|A\| \|U\|$$

and $U \neq 0$, we have $\|U\| \neq 0$, and get

$$\rho(A) = |\lambda| \leq \|A\|,$$

as claimed. □

Proposition 7.4 also holds for any real matrix norm $\| \cdot \|$ on $M_n(\mathbb{R})$ but the proof is more subtle and requires the notion of induced norm. We prove it after giving Definition 7.7.

Now, it turns out that if A is a real $n \times n$ symmetric matrix, then the eigenvalues of A are all real and there is some orthogonal matrix Q such that

$$A = Q \operatorname{diag}(\lambda_1, \dots, \lambda_n) Q^\top,$$

where $\operatorname{diag}(\lambda_1, \dots, \lambda_n)$ denotes the matrix whose only nonzero entries (if any) are its diagonal entries, which are the (real) eigenvalues of A . Similarly, if A is a complex $n \times n$ Hermitian matrix, then the eigenvalues of A are all real and there is some unitary matrix U such that

$$A = U \operatorname{diag}(\lambda_1, \dots, \lambda_n) U^*,$$

where $\operatorname{diag}(\lambda_1, \dots, \lambda_n)$ denotes the matrix whose only nonzero entries (if any) are its diagonal entries, which are the (real) eigenvalues of A .

We now return to matrix norms. We begin with the so-called *Frobenius norm*, which is just the norm $\| \cdot \|_2$ on \mathbb{C}^{n^2} , where the $n \times n$ matrix A is viewed as the vector obtained by concatenating together the rows (or the columns) of A . The reader should check that for any $n \times n$ complex matrix $A = (a_{ij})$,

$$\left(\sum_{i,j=1}^n |a_{ij}|^2 \right)^{1/2} = \sqrt{\operatorname{tr}(A^* A)} = \sqrt{\operatorname{tr}(A A^*)}.$$

Definition 7.6. The *Frobenius norm* $\| \cdot \|_F$ is defined so that for every square $n \times n$ matrix $A \in M_n(\mathbb{C})$,

$$\|A\|_F = \left(\sum_{i,j=1}^n |a_{ij}|^2 \right)^{1/2} = \sqrt{\operatorname{tr}(A A^*)} = \sqrt{\operatorname{tr}(A^* A)}.$$

The following proposition shows that the Frobenius norm is a matrix norm satisfying other nice properties.

Proposition 7.5. *The Frobenius norm $\| \cdot \|_F$ on $M_n(\mathbb{C})$ satisfies the following properties:*

- (1) *It is a matrix norm; that is, $\|AB\|_F \leq \|A\|_F \|B\|_F$, for all $A, B \in M_n(\mathbb{C})$.*
- (2) *It is unitarily invariant, which means that for all unitary matrices U, V , we have*

$$\|A\|_F = \|UA\|_F = \|AV\|_F = \|UAV\|_F.$$

- (3) *$\sqrt{\rho(A^* A)} \leq \|A\|_F \leq \sqrt{n} \sqrt{\rho(A^* A)}$, for all $A \in M_n(\mathbb{C})$.*

Proof. (1) The only property that requires a proof is the fact $\|AB\|_F \leq \|A\|_F \|B\|_F$. This follows from the Cauchy–Schwarz inequality:

$$\begin{aligned}\|AB\|_F^2 &= \sum_{i,j=1}^n \left| \sum_{k=1}^n a_{ik} b_{kj} \right|^2 \\ &\leq \sum_{i,j=1}^n \left(\sum_{h=1}^n |a_{ih}|^2 \right) \left(\sum_{k=1}^n |b_{kj}|^2 \right) \\ &= \left(\sum_{i,h=1}^n |a_{ih}|^2 \right) \left(\sum_{k,j=1}^n |b_{kj}|^2 \right) = \|A\|_F^2 \|B\|_F^2.\end{aligned}$$

(2) We have

$$\|A\|_F^2 = \operatorname{tr}(A^*A) = \operatorname{tr}(VV^*A^*A) = \operatorname{tr}(V^*A^*AV) = \|AV\|_F^2,$$

and

$$\|A\|_F^2 = \operatorname{tr}(A^*A) = \operatorname{tr}(A^*U^*UA) = \|UA\|_F^2.$$

The identity

$$\|A\|_F = \|UAV\|_F$$

follows from the previous two.

(3) It is well known that the trace of a matrix is equal to the sum of its eigenvalues. Furthermore, A^*A is symmetric positive semidefinite (which means that its eigenvalues are nonnegative), so $\rho(A^*A)$ is the largest eigenvalue of A^*A and

$$\rho(A^*A) \leq \operatorname{tr}(A^*A) \leq n\rho(A^*A),$$

which yields (3) by taking square roots. □

Remark: The Frobenius norm is also known as the *Hilbert-Schmidt norm* or the *Schur norm*. So many famous names associated with such a simple thing!

We now give another method for obtaining matrix norms using subordinate norms. First, we need a proposition that shows that in a finite-dimensional space, the linear map induced by a matrix is bounded, and thus continuous.

Proposition 7.6. *For every norm $\|\cdot\|$ on \mathbb{C}^n (or \mathbb{R}^n), for every matrix $A \in M_n(\mathbb{C})$ (or $A \in M_n(\mathbb{R})$), there is a real constant $C_A \geq 0$, such that*

$$\|Au\| \leq C_A \|u\|,$$

for every vector $u \in \mathbb{C}^n$ (or $u \in \mathbb{R}^n$ if A is real).

Proof. For every basis (e_1, \dots, e_n) of \mathbb{C}^n (or \mathbb{R}^n), for every vector $u = u_1 e_1 + \dots + u_n e_n$, we have

$$\begin{aligned}\|Au\| &= \|u_1 A(e_1) + \dots + u_n A(e_n)\| \\ &\leq |u_1| \|A(e_1)\| + \dots + |u_n| \|A(e_n)\| \\ &\leq C_1(|u_1| + \dots + |u_n|) = C_1 \|u\|_1,\end{aligned}$$

where $C_1 = \max_{1 \leq i \leq n} \|A(e_i)\|$. By Theorem 7.3, the norms $\|\cdot\|$ and $\|\cdot\|_1$ are equivalent, so there is some constant $C_2 > 0$ so that $\|u\|_1 \leq C_2 \|u\|$ for all u , which implies that

$$\|Au\| \leq C_A \|u\|,$$

where $C_A = C_1 C_2$. □

Proposition 7.6 says that every linear map on a finite-dimensional space is *bounded*. This implies that every linear map on a finite-dimensional space is continuous. Actually, it is not hard to show that a linear map on a normed vector space E is bounded iff it is continuous, regardless of the dimension of E .

Proposition 7.6 implies that for every matrix $A \in M_n(\mathbb{C})$ (or $A \in M_n(\mathbb{R})$),

$$\sup_{\substack{x \in \mathbb{C}^n \\ x \neq 0}} \frac{\|Ax\|}{\|x\|} \leq C_A.$$

Now, since $\|\lambda u\| = |\lambda| \|u\|$, for every nonzero vector x , we have

$$\frac{\|Ax\|}{\|x\|} = \frac{\|x\| \|A(x/\|x\|)\|}{\|x\| \|(x/\|x\|)\|} = \frac{\|A(x/\|x\|)\|}{\|(x/\|x\|)\|},$$

which implies that

$$\sup_{\substack{x \in \mathbb{C}^n \\ x \neq 0}} \frac{\|Ax\|}{\|x\|} = \sup_{\substack{x \in \mathbb{C}^n \\ \|x\|=1}} \|Ax\|.$$

Similarly

$$\sup_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{\|Ax\|}{\|x\|} = \sup_{\substack{x \in \mathbb{R}^n \\ \|x\|=1}} \|Ax\|.$$

The above considerations justify the following definition.

Definition 7.7. If $\|\cdot\|$ is any norm on \mathbb{C}^n , we define the function $\|\cdot\|$ on $M_n(\mathbb{C})$ by

$$\|A\| = \sup_{\substack{x \in \mathbb{C}^n \\ x \neq 0}} \frac{\|Ax\|}{\|x\|} = \sup_{\substack{x \in \mathbb{C}^n \\ \|x\|=1}} \|Ax\|.$$

The function $A \mapsto \|A\|$ is called the *subordinate matrix norm* or *operator norm* induced by the norm $\|\cdot\|$.

It is easy to check that the function $A \mapsto \|A\|$ is indeed a norm, and by definition, it satisfies the property

$$\|Ax\| \leq \|A\| \|x\|, \quad \text{for all } x \in \mathbb{C}^n.$$

A norm $\|\cdot\|$ on $M_n(\mathbb{C})$ satisfying the above property is said to be *subordinate* to the vector norm $\|\cdot\|$ on \mathbb{C}^n . As a consequence of the above inequality, we have

$$\|ABx\| \leq \|A\| \|Bx\| \leq \|A\| \|B\| \|x\|,$$

for all $x \in \mathbb{C}^n$, which implies that

$$\|AB\| \leq \|A\| \|B\| \quad \text{for all } A, B \in M_n(\mathbb{C}),$$

showing that $A \mapsto \|A\|$ is a matrix norm (it is submultiplicative).

Observe that the operator norm is also defined by

$$\|A\| = \inf\{\lambda \in \mathbb{R} \mid \|Ax\| \leq \lambda \|x\|, \text{ for all } x \in \mathbb{C}^n\}.$$

Since the function $x \mapsto \|Ax\|$ is continuous (because $|\|Ay\| - \|Ax\|| \leq \|Ay - Ax\| \leq C_A \|x - y\|$) and the unit sphere $S^{n-1} = \{x \in \mathbb{C}^n \mid \|x\| = 1\}$ is compact, there is some $x \in \mathbb{C}^n$ such that $\|x\| = 1$ and

$$\|Ax\| = \|A\|.$$

Equivalently, there is some $x \in \mathbb{C}^n$ such that $x \neq 0$ and

$$\|Ax\| = \|A\| \|x\|.$$

The definition of an operator norm also implies that

$$\|I\| = 1.$$

The above shows that the Frobenius norm is not a subordinate matrix norm (why?). The notion of subordinate norm can be slightly generalized.

Definition 7.8. If $K = \mathbb{R}$ or $K = \mathbb{C}$, for any norm $\|\cdot\|$ on $M_{m,n}(K)$, and for any two norms $\|\cdot\|_a$ on K^n and $\|\cdot\|_b$ on K^m , we say that the norm $\|\cdot\|$ is *subordinate* to the norms $\|\cdot\|_a$ and $\|\cdot\|_b$ if

$$\|Ax\|_b \leq \|A\| \|x\|_a \quad \text{for all } A \in M_{m,n}(K) \text{ and all } x \in K^n.$$

Remark: For any norm $\|\cdot\|$ on \mathbb{C}^n , we can define the function $\|\cdot\|_{\mathbb{R}}$ on $M_n(\mathbb{R})$ by

$$\|A\|_{\mathbb{R}} = \sup_{\substack{x \in \mathbb{R}^n \\ x \neq 0}} \frac{\|Ax\|}{\|x\|} = \sup_{\substack{x \in \mathbb{R}^n \\ \|x\|=1}} \|Ax\|.$$

The function $A \mapsto \|A\|_{\mathbb{R}}$ is a matrix norm on $M_n(\mathbb{R})$, and

$$\|A\|_{\mathbb{R}} \leq \|A\|,$$

for all real matrices $A \in M_n(\mathbb{R})$. However, it is possible to construct vector norms $\|\cdot\|$ on \mathbb{C}^n and *real* matrices A such that

$$\|A\|_{\mathbb{R}} < \|A\|.$$

In order to avoid this kind of difficulties, we define subordinate matrix norms over $M_n(\mathbb{C})$. Luckily, it turns out that $\|A\|_{\mathbb{R}} = \|A\|$ for the vector norms, $\|\cdot\|_1$, $\|\cdot\|_2$, and $\|\cdot\|_{\infty}$.

We now prove Proposition 7.4 for real matrix norms.

Proposition 7.7. *For any matrix norm $\|\cdot\|$ on $M_n(\mathbb{R})$ and for any square $n \times n$ matrix $A \in M_n(\mathbb{R})$, we have*

$$\rho(A) \leq \|A\|.$$

Proof. We follow the proof in Denis Serre's book [96]. If A is a real matrix, the problem is that the eigenvectors associated with the eigenvalue of maximum modulus may be complex. We use a trick based on the fact that for every matrix A (real or complex),

$$\rho(A^k) = (\rho(A))^k,$$

which is left as an exercise (use Proposition 8.5 which shows that if $(\lambda_1, \dots, \lambda_n)$ are the (not necessarily distinct) eigenvalues of A , then $(\lambda_1^k, \dots, \lambda_n^k)$ are the eigenvalues of A^k , for $k \geq 1$).

Pick any complex norm $\|\cdot\|_c$ on \mathbb{C}^n and let $\|\cdot\|_c$ denote the corresponding induced norm on matrices. The restriction of $\|\cdot\|_c$ to real matrices is a real norm that we also denote by $\|\cdot\|_c$. Now, by Theorem 7.3, since $M_n(\mathbb{R})$ has finite dimension n^2 , there is some constant $C > 0$ so that

$$\|A\|_c \leq C \|A\|, \quad \text{for all } A \in M_n(\mathbb{R}).$$

Furthermore, for every $k \geq 1$ and for every real $n \times n$ matrix A , by Proposition 7.4, $\rho(A^k) \leq \|A^k\|_c$, and because $\|\cdot\|$ is a matrix norm, $\|A^k\| \leq \|A\|^k$, so we have

$$(\rho(A))^k = \rho(A^k) \leq \|A^k\|_c \leq C \|A^k\| \leq C \|A\|^k,$$

for all $k \geq 1$. It follows that

$$\rho(A) \leq C^{1/k} \|A\|, \quad \text{for all } k \geq 1.$$

However because $C > 0$, we have $\lim_{k \rightarrow \infty} C^{1/k} = 1$ (we have $\lim_{k \rightarrow \infty} \frac{1}{k} \log(C) = 0$). Therefore, we conclude that

$$\rho(A) \leq \|A\|,$$

as desired. □

We now determine explicitly what are the subordinate matrix norms associated with the vector norms $\|\cdot\|_1$, $\|\cdot\|_2$, and $\|\cdot\|_{\infty}$.

Proposition 7.8. *For every square matrix $A = (a_{ij}) \in M_n(\mathbb{C})$, we have*

$$\begin{aligned}\|A\|_1 &= \sup_{\substack{x \in \mathbb{C}^n \\ \|x\|_1=1}} \|Ax\|_1 = \max_j \sum_{i=1}^n |a_{ij}| \\ \|A\|_\infty &= \sup_{\substack{x \in \mathbb{C}^n \\ \|x\|_\infty=1}} \|Ax\|_\infty = \max_i \sum_{j=1}^n |a_{ij}| \\ \|A\|_2 &= \sup_{\substack{x \in \mathbb{C}^n \\ \|x\|_2=1}} \|Ax\|_2 = \sqrt{\rho(A^*A)} = \sqrt{\rho(AA^*)}.\end{aligned}$$

Furthermore, $\|A^*\|_2 = \|A\|_2$, the norm $\|\cdot\|_2$ is unitarily invariant, which means that

$$\|A\|_2 = \|UAV\|_2$$

for all unitary matrices U, V , and if A is a normal matrix, then $\|A\|_2 = \rho(A)$.

Proof. For every vector u , we have

$$\|Au\|_1 = \sum_i \left| \sum_j a_{ij} u_j \right| \leq \sum_j |u_j| \sum_i |a_{ij}| \leq \left(\max_j \sum_i |a_{ij}| \right) \|u\|_1,$$

which implies that

$$\|A\|_1 \leq \max_j \sum_{i=1}^n |a_{ij}|.$$

It remains to show that equality can be achieved. For this let j_0 be some index such that

$$\max_j \sum_i |a_{ij}| = \sum_i |a_{ij_0}|,$$

and let $u_i = 0$ for all $i \neq j_0$ and $u_{j_0} = 1$.

In a similar way, we have

$$\|Au\|_\infty = \max_i \left| \sum_j a_{ij} u_j \right| \leq \left(\max_i \sum_j |a_{ij}| \right) \|u\|_\infty,$$

which implies that

$$\|A\|_\infty \leq \max_i \sum_{j=1}^n |a_{ij}|.$$

To achieve equality, let i_0 be some index such that

$$\max_i \sum_j |a_{ij}| = \sum_j |a_{i_0j}|.$$

The reader should check that the vector given by

$$u_j = \begin{cases} \frac{\bar{a}_{i_0 j}}{|a_{i_0 j}|} & \text{if } a_{i_0 j} \neq 0 \\ 1 & \text{if } a_{i_0 j} = 0 \end{cases}$$

works.

We have

$$\|A\|_2^2 = \sup_{\substack{x \in \mathbb{C}^n \\ x^* x = 1}} \|Ax\|_2^2 = \sup_{\substack{x \in \mathbb{C}^n \\ x^* x = 1}} x^* A^* A x.$$

Since the matrix A^*A is symmetric, it has real eigenvalues and it can be diagonalized with respect to an orthogonal matrix. These facts can be used to prove that the function $x \mapsto x^* A^* A x$ has a maximum on the sphere $x^* x = 1$ equal to the largest eigenvalue of A^*A , namely, $\rho(A^*A)$. We postpone the proof until we discuss optimizing quadratic functions. Therefore,

$$\|A\|_2 = \sqrt{\rho(A^*A)}.$$

Let us now prove that $\rho(A^*A) = \rho(AA^*)$. First, assume that $\rho(A^*A) > 0$. In this case, there is some eigenvector $u (\neq 0)$ such that

$$A^* A u = \rho(A^*A) u,$$

and since $\rho(A^*A) > 0$, we must have $Au \neq 0$. Since $Au \neq 0$,

$$A A^* (A u) = \rho(A^*A) A u$$

which means that $\rho(A^*A)$ is an eigenvalue of AA^* , and thus

$$\rho(A^*A) \leq \rho(AA^*).$$

Because $(A^*)^* = A$, by replacing A by A^* , we get

$$\rho(AA^*) \leq \rho(A^*A),$$

and so $\rho(A^*A) = \rho(AA^*)$.

If $\rho(A^*A) = 0$, then we must have $\rho(AA^*) = 0$, since otherwise by the previous reasoning we would have $\rho(A^*A) = \rho(AA^*) > 0$. Hence, in all case

$$\|A\|_2^2 = \rho(A^*A) = \rho(AA^*) = \|A^*\|_2^2.$$

For any unitary matrices U and V , it is an easy exercise to prove that $V^* A^* A V$ and $A^* A$ have the same eigenvalues, so

$$\|A\|_2^2 = \rho(A^*A) = \rho(V^* A^* A V) = \|A V\|_2^2,$$

and also

$$\|A\|_2^2 = \rho(A^*A) = \rho(A^*U^*UA) = \|UA\|_2^2.$$

Finally, if A is a normal matrix ($AA^* = A^*A$), it can be shown that there is some unitary matrix U so that

$$A = UDU^*,$$

where $D = \text{diag}(\lambda_1, \dots, \lambda_n)$ is a diagonal matrix consisting of the eigenvalues of A , and thus

$$A^*A = (UDU^*)^*UDU^* = UD^*U^*UDU^* = UD^*DU^*.$$

However, $D^*D = \text{diag}(|\lambda_1|^2, \dots, |\lambda_n|^2)$, which proves that

$$\rho(A^*A) = \rho(D^*D) = \max_i |\lambda_i|^2 = (\rho(A))^2,$$

so that $\|A\|_2 = \rho(A)$. □

The norm $\|A\|_2$ is often called the *spectral norm*. Observe that property (3) of proposition 7.5 says that

$$\|A\|_2 \leq \|A\|_F \leq \sqrt{n} \|A\|_2,$$

which shows that the Frobenius norm is an upper bound on the spectral norm. The Frobenius norm is much easier to compute than the spectral norm.

The reader will check that the above proof still holds if the matrix A is real, confirming the fact that $\|A\|_{\mathbb{R}} = \|A\|$ for the vector norms $\|\cdot\|_1$, $\|\cdot\|_2$, and $\|\cdot\|_{\infty}$. It is also easy to verify that the proof goes through for *rectangular* matrices, with the same formulae. Similarly, the Frobenius norm is also a norm on rectangular matrices. For these norms, whenever AB makes sense, we have

$$\|AB\| \leq \|A\| \|B\|.$$

Remark: Let $(E, \|\cdot\|)$ and $(F, \|\cdot\|)$ be two normed vector spaces (for simplicity of notation, we use the same symbol $\|\cdot\|$ for the norms on E and F ; this should not cause any confusion). Recall that a function $f: E \rightarrow F$ is *continuous* if for every $a \in E$, for every $\epsilon > 0$, there is some $\eta > 0$ such that for all $x \in E$,

$$\text{if } \|x - a\| \leq \eta \text{ then } \|f(x) - f(a)\| \leq \epsilon.$$

It is not hard to show that a *linear map* $f: E \rightarrow F$ is continuous iff there is some constant $C \geq 0$ such that

$$\|f(x)\| \leq C \|x\| \text{ for all } x \in E.$$

If so, we say that f is *bounded* (or a *linear bounded operator*). We let $\mathcal{L}(E; F)$ denote the set of all continuous (equivalently, bounded) linear maps from E to F . Then, we can define the *operator norm* (or *subordinate norm*) $\|\cdot\|$ on $\mathcal{L}(E; F)$ as follows: for every $f \in \mathcal{L}(E; F)$,

$$\|f\| = \sup_{\substack{x \in E \\ x \neq 0}} \frac{\|f(x)\|}{\|x\|} = \sup_{\substack{x \in E \\ \|x\|=1}} \|f(x)\|,$$

or equivalently by

$$\|f\| = \inf\{\lambda \in \mathbb{R} \mid \|f(x)\| \leq \lambda \|x\|, \text{ for all } x \in E\}.$$

It is not hard to show that the map $f \mapsto \|f\|$ is a norm on $\mathcal{L}(E; F)$ satisfying the property

$$\|f(x)\| \leq \|f\| \|x\|$$

for all $x \in E$, and that if $f \in \mathcal{L}(E; F)$ and $g \in \mathcal{L}(F; G)$, then

$$\|g \circ f\| \leq \|g\| \|f\|.$$

Operator norms play an important role in functional analysis, especially when the spaces E and F are *complete*.

The following proposition will be needed when we deal with the condition number of a matrix.

Proposition 7.9. *Let $\|\cdot\|$ be any matrix norm and let B be a matrix such that $\|B\| < 1$.*

(1) If $\|\cdot\|$ is a subordinate matrix norm, then the matrix $I + B$ is invertible and

$$\|(I + B)^{-1}\| \leq \frac{1}{1 - \|B\|}.$$

(2) If a matrix of the form $I + B$ is singular, then $\|B\| \geq 1$ for every matrix norm (not necessarily subordinate).

Proof. (1) Observe that $(I + B)u = 0$ implies $Bu = -u$, so

$$\|u\| = \|Bu\|.$$

Recall that

$$\|Bu\| \leq \|B\| \|u\|$$

for every subordinate norm. Since $\|B\| < 1$, if $u \neq 0$, then

$$\|Bu\| < \|u\|,$$

which contradicts $\|u\| = \|Bu\|$. Therefore, we must have $u = 0$, which proves that $I + B$ is injective, and thus bijective, i.e., invertible. Then, we have

$$(I + B)^{-1} + B(I + B)^{-1} = (I + B)(I + B)^{-1} = I,$$

so we get

$$(I + B)^{-1} = I - B(I + B)^{-1},$$

which yields

$$\|(I + B)^{-1}\| \leq 1 + \|B\| \|(I + B)^{-1}\|,$$

and finally,

$$\|(I + B)^{-1}\| \leq \frac{1}{1 - \|B\|}.$$

(2) If $I + B$ is singular, then -1 is an eigenvalue of B , and by Proposition 7.4, we get $\rho(B) \leq \|B\|$, which implies $1 \leq \rho(B) \leq \|B\|$. \square

The following result is needed to deal with the convergence of sequences of powers of matrices.

Proposition 7.10. *For every matrix $A \in M_n(\mathbb{C})$ and for every $\epsilon > 0$, there is some subordinate matrix norm $\|\cdot\|$ such that*

$$\|A\| \leq \rho(A) + \epsilon.$$

Proof. By Theorem 8.4, there exists some invertible matrix U and some upper triangular matrix T such that

$$A = UTU^{-1},$$

and say that

$$T = \begin{pmatrix} \lambda_1 & t_{12} & t_{13} & \cdots & t_{1n} \\ 0 & \lambda_2 & t_{23} & \cdots & t_{2n} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & \lambda_{n-1} & t_{n-1n} \\ 0 & 0 & \cdots & 0 & \lambda_n \end{pmatrix},$$

where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of A . For every $\delta \neq 0$, define the diagonal matrix

$$D_\delta = \text{diag}(1, \delta, \delta^2, \dots, \delta^{n-1}),$$

and consider the matrix

$$(UD_\delta)^{-1}A(UD_\delta) = D_\delta^{-1}TD_\delta = \begin{pmatrix} \lambda_1 & \delta t_{12} & \delta^2 t_{13} & \cdots & \delta^{n-1} t_{1n} \\ 0 & \lambda_2 & \delta t_{23} & \cdots & \delta^{n-2} t_{2n} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & \lambda_{n-1} & \delta t_{n-1n} \\ 0 & 0 & \cdots & 0 & \lambda_n \end{pmatrix}.$$

Now, define the function $\|\cdot\|: M_n(\mathbb{C}) \rightarrow \mathbb{R}$ by

$$\|B\| = \|(UD_\delta)^{-1}B(UD_\delta)\|_\infty,$$

for every $B \in M_n(\mathbb{C})$. Then it is easy to verify that the above function is the matrix norm subordinate to the vector norm

$$v \mapsto \|(UD_\delta)^{-1}v\|_\infty.$$

Furthermore, for every $\epsilon > 0$, we can pick δ so that

$$\sum_{j=i+1}^n |\delta^{j-i} t_{ij}| \leq \epsilon, \quad 1 \leq i \leq n-1,$$

and by definition of the norm $\|\cdot\|_\infty$, we get

$$\|A\| \leq \rho(A) + \epsilon,$$

which shows that the norm that we have constructed satisfies the required properties. \square

Note that equality is generally not possible; consider the matrix

$$A = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix},$$

for which $\rho(A) = 0 < \|A\|$, since $A \neq 0$.

7.3 Condition Numbers of Matrices

Unfortunately, there exist linear systems $Ax = b$ whose solutions are not stable under small perturbations of either b or A . For example, consider the system

$$\begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}.$$

The reader should check that it has the solution $x = (1, 1, 1, 1)$. If we perturb slightly the right-hand side, obtaining the new system

$$\begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix} \begin{pmatrix} x_1 + \Delta x_1 \\ x_2 + \Delta x_2 \\ x_3 + \Delta x_3 \\ x_4 + \Delta x_4 \end{pmatrix} = \begin{pmatrix} 32.1 \\ 22.9 \\ 33.1 \\ 30.9 \end{pmatrix},$$

the new solutions turns out to be $x = (9.2, -12.6, 4.5, -1.1)$. In other words, a relative error of the order $1/200$ in the data (here, b) produces a relative error of the order $10/1$ in the solution, which represents an amplification of the relative error of the order 2000 .

Now, let us perturb the matrix slightly, obtaining the new system

$$\begin{pmatrix} 10 & 7 & 8.1 & 7.2 \\ 7.08 & 5.04 & 6 & 5 \\ 8 & 5.98 & 9.98 & 9 \\ 6.99 & 4.99 & 9 & 9.98 \end{pmatrix} \begin{pmatrix} x_1 + \Delta x_1 \\ x_2 + \Delta x_2 \\ x_3 + \Delta x_3 \\ x_4 + \Delta x_4 \end{pmatrix} = \begin{pmatrix} 32 \\ 23 \\ 33 \\ 31 \end{pmatrix}.$$

This time, the solution is $x = (-81, 137, -34, 22)$. Again, a small change in the data alters the result rather drastically. Yet, the original system is symmetric, has determinant 1, and has integer entries. The problem is that the matrix of the system is badly conditioned, a concept that we will now explain.

Given an invertible matrix A , first, assume that we perturb b to $b + \delta b$, and let us analyze the change between the two exact solutions x and $x + \delta x$ of the two systems

$$\begin{aligned} Ax &= b \\ A(x + \delta x) &= b + \delta b. \end{aligned}$$

We also assume that we have some norm $\| \cdot \|$ and we use the subordinate matrix norm on matrices. From

$$\begin{aligned} Ax &= b \\ Ax + A\delta x &= b + \delta b, \end{aligned}$$

we get

$$\delta x = A^{-1}\delta b,$$

and we conclude that

$$\begin{aligned} \|\delta x\| &\leq \|A^{-1}\| \|\delta b\| \\ \|b\| &\leq \|A\| \|x\|. \end{aligned}$$

Consequently, the relative error in the result $\|\delta x\| / \|x\|$ is bounded in terms of the relative error $\|\delta b\| / \|b\|$ in the data as follows:

$$\frac{\|\delta x\|}{\|x\|} \leq (\|A\| \|A^{-1}\|) \frac{\|\delta b\|}{\|b\|}.$$

Now let us assume that A is perturbed to $A + \delta A$, and let us analyze the change between the exact solutions of the two systems

$$\begin{aligned} Ax &= b \\ (A + \Delta A)(x + \Delta x) &= b. \end{aligned}$$

The second equation yields $Ax + A\Delta x + \Delta A(x + \Delta x) = b$, and by subtracting the first equation we get

$$\Delta x = -A^{-1}\Delta A(x + \Delta x).$$

It follows that

$$\|\Delta x\| \leq \|A^{-1}\| \|\Delta A\| \|x + \Delta x\|,$$

which can be rewritten as

$$\frac{\|\Delta x\|}{\|x + \Delta x\|} \leq (\|A\| \|A^{-1}\|) \frac{\|\Delta A\|}{\|A\|}.$$

Observe that the above reasoning is valid even if the matrix $A + \Delta A$ is singular, as long as $x + \Delta x$ is a solution of the second system. Furthermore, if $\|\Delta A\|$ is small enough, it is not unreasonable to expect that the ratio $\|\Delta x\| / \|x + \Delta x\|$ is close to $\|\Delta x\| / \|x\|$. This will be made more precise later.

In summary, for each of the two perturbations, we see that the relative error in the result is bounded by the relative error in the data, *multiplied the number* $\|A\| \|A^{-1}\|$. In fact, this factor turns out to be optimal and this suggests the following definition:

Definition 7.9. For any subordinate matrix norm $\|\cdot\|$, for any invertible matrix A , the number

$$\text{cond}(A) = \|A\| \|A^{-1}\|$$

is called the *condition number* of A relative to $\|\cdot\|$.

The condition number $\text{cond}(A)$ measures the sensitivity of the linear system $Ax = b$ to variations in the data b and A ; a feature referred to as the *condition* of the system. Thus, when we says that a linear system is *ill-conditioned*, we mean that the condition number of its matrix is large. We can sharpen the preceding analysis as follows:

Proposition 7.11. Let A be an invertible matrix and let x and $x + \delta x$ be the solutions of the linear systems

$$\begin{aligned} Ax &= b \\ A(x + \delta x) &= b + \delta b. \end{aligned}$$

If $b \neq 0$, then the inequality

$$\frac{\|\delta x\|}{\|x\|} \leq \text{cond}(A) \frac{\|\delta b\|}{\|b\|}$$

holds and is the best possible. This means that for a given matrix A , there exist some vectors $b \neq 0$ and $\delta b \neq 0$ for which equality holds.

Proof. We already proved the inequality. Now, because $\|\cdot\|$ is a subordinate matrix norm, there exist some vectors $x \neq 0$ and $\delta b \neq 0$ for which

$$\|A^{-1}\delta b\| = \|A^{-1}\| \|\delta b\| \quad \text{and} \quad \|Ax\| = \|A\| \|x\|.$$

□

Proposition 7.12. *Let A be an invertible matrix and let x and $x + \Delta x$ be the solutions of the two systems*

$$\begin{aligned} Ax &= b \\ (A + \Delta A)(x + \Delta x) &= b. \end{aligned}$$

If $b \neq 0$, then the inequality

$$\frac{\|\Delta x\|}{\|x + \Delta x\|} \leq \text{cond}(A) \frac{\|\Delta A\|}{\|A\|}$$

holds and is the best possible. This means that given a matrix A , there exist a vector $b \neq 0$ and a matrix $\Delta A \neq 0$ for which equality holds. Furthermore, if $\|\Delta A\|$ is small enough (for instance, if $\|\Delta A\| < 1/\|A^{-1}\|$), we have

$$\frac{\|\Delta x\|}{\|x\|} \leq \text{cond}(A) \frac{\|\Delta A\|}{\|A\|} (1 + O(\|\Delta A\|));$$

in fact, we have

$$\frac{\|\Delta x\|}{\|x\|} \leq \text{cond}(A) \frac{\|\Delta A\|}{\|A\|} \left(\frac{1}{1 - \|A^{-1}\| \|\Delta A\|} \right).$$

Proof. The first inequality has already been proved. To show that equality can be achieved, let w be any vector such that $w \neq 0$ and

$$\|A^{-1}w\| = \|A^{-1}\| \|w\|,$$

and let $\beta \neq 0$ be any real number. Now, the vectors

$$\begin{aligned} \Delta x &= -\beta A^{-1}w \\ x + \Delta x &= w \\ b &= (A + \beta I)w \end{aligned}$$

and the matrix

$$\Delta A = \beta I$$

satisfy the equations

$$\begin{aligned} Ax &= b \\ (A + \Delta A)(x + \Delta x) &= b \\ \|\Delta x\| &= |\beta| \|A^{-1}w\| = \|\Delta A\| \|A^{-1}\| \|x + \Delta x\|. \end{aligned}$$

Finally, we can pick β so that $-\beta$ is not equal to any of the eigenvalues of A , so that $A + \Delta A = A + \beta I$ is invertible and b is nonzero.

If $\|\Delta A\| < 1/\|A^{-1}\|$, then

$$\|A^{-1}\Delta A\| \leq \|A^{-1}\| \|\Delta A\| < 1,$$

so by Proposition 7.9, the matrix $I + A^{-1}\Delta A$ is invertible and

$$\|(I + A^{-1}\Delta A)^{-1}\| \leq \frac{1}{1 - \|A^{-1}\Delta A\|} \leq \frac{1}{1 - \|A^{-1}\| \|\Delta A\|}.$$

Recall that we proved earlier that

$$\Delta x = -A^{-1}\Delta A(x + \Delta x),$$

and by adding x to both sides and moving the right-hand side to the left-hand side yields

$$(I + A^{-1}\Delta A)(x + \Delta x) = x,$$

and thus

$$x + \Delta x = (I + A^{-1}\Delta A)^{-1}x,$$

which yields

$$\begin{aligned} \Delta x &= ((I + A^{-1}\Delta A)^{-1} - I)x = (I + A^{-1}\Delta A)^{-1}(I - (I + A^{-1}\Delta A))x \\ &= -(I + A^{-1}\Delta A)^{-1}A^{-1}(\Delta A)x. \end{aligned}$$

From this and

$$\|(I + A^{-1}\Delta A)^{-1}\| \leq \frac{1}{1 - \|A^{-1}\| \|\Delta A\|},$$

we get

$$\|\Delta x\| \leq \frac{\|A^{-1}\| \|\Delta A\|}{1 - \|A^{-1}\| \|\Delta A\|} \|x\|,$$

which can be written as

$$\frac{\|\Delta x\|}{\|x\|} \leq \text{cond}(A) \frac{\|\Delta A\|}{\|A\|} \left(\frac{1}{1 - \|A^{-1}\| \|\Delta A\|} \right),$$

which is the kind of inequality that we were seeking. □

Remark: If A and b are perturbed simultaneously, so that we get the “perturbed” system

$$(A + \Delta A)(x + \delta x) = b + \delta b,$$

it can be shown that if $\|\Delta A\| < 1/\|A^{-1}\|$ (and $b \neq 0$), then

$$\frac{\|\Delta x\|}{\|x\|} \leq \frac{\text{cond}(A)}{1 - \|A^{-1}\| \|\Delta A\|} \left(\frac{\|\Delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right);$$

see Demmel [27], Section 2.2 and Horn and Johnson [57], Section 5.8.

We now list some properties of condition numbers and figure out what $\text{cond}(A)$ is in the case of the spectral norm (the matrix norm induced by $\|\cdot\|_2$). First, we need to introduce a very important factorization of matrices, the *singular value decomposition*, for short, *SVD*.

It can be shown that given any $n \times n$ matrix $A \in M_n(\mathbb{C})$, there exist two unitary matrices U and V , and a *real* diagonal matrix $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$, with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$, such that

$$A = V\Sigma U^*.$$

The nonnegative numbers $\sigma_1, \dots, \sigma_n$ are called the *singular values* of A .

If A is a real matrix, the matrices U and V are orthogonal matrices. The factorization $A = V\Sigma U^*$ implies that

$$A^*A = U\Sigma^2U^* \quad \text{and} \quad AA^* = V\Sigma^2V^*,$$

which shows that $\sigma_1^2, \dots, \sigma_n^2$ are the eigenvalues of *both* A^*A and AA^* , that the columns of U are corresponding eivenvectors for A^*A , and that the columns of V are corresponding eivenvectors for AA^* . In the case of a normal matrix if $\lambda_1, \dots, \lambda_n$ are the (complex) eigenvalues of A , then

$$\sigma_i = |\lambda_i|, \quad 1 \leq i \leq n.$$

Proposition 7.13. *For every invertible matrix $A \in M_n(\mathbb{C})$, the following properties hold:*

(1)

$$\begin{aligned} \text{cond}(A) &\geq 1, \\ \text{cond}(A) &= \text{cond}(A^{-1}) \\ \text{cond}(\alpha A) &= \text{cond}(A) \quad \text{for all } \alpha \in \mathbb{C} - \{0\}. \end{aligned}$$

(2) *If $\text{cond}_2(A)$ denotes the condition number of A with respect to the spectral norm, then*

$$\text{cond}_2(A) = \frac{\sigma_1}{\sigma_n},$$

where $\sigma_1 \geq \dots \geq \sigma_n$ are the singular values of A .

(3) *If the matrix A is normal, then*

$$\text{cond}_2(A) = \frac{|\lambda_1|}{|\lambda_n|},$$

where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of A sorted so that $|\lambda_1| \geq \dots \geq |\lambda_n|$.

(4) *If A is a unitary or an orthogonal matrix, then*

$$\text{cond}_2(A) = 1.$$

(5) The condition number $\text{cond}_2(A)$ is invariant under unitary transformations, which means that

$$\text{cond}_2(A) = \text{cond}_2(UA) = \text{cond}_2(AV),$$

for all unitary matrices U and V .

Proof. The properties in (1) are immediate consequences of the properties of subordinate matrix norms. In particular, $AA^{-1} = I$ implies

$$1 = \|I\| \leq \|A\| \|A^{-1}\| = \text{cond}(A).$$

(2) We showed earlier that $\|A\|_2^2 = \rho(A^*A)$, which is the square of the modulus of the largest eigenvalue of A^*A . Since we just saw that the eigenvalues of A^*A are $\sigma_1^2 \geq \cdots \geq \sigma_n^2$, where $\sigma_1, \dots, \sigma_n$ are the singular values of A , we have

$$\|A\|_2 = \sigma_1.$$

Now, if A is invertible, then $\sigma_1 \geq \cdots \geq \sigma_n > 0$, and it is easy to show that the eigenvalues of $(A^*A)^{-1}$ are $\sigma_n^{-2} \geq \cdots \geq \sigma_1^{-2}$, which shows that

$$\|A^{-1}\|_2 = \sigma_n^{-1},$$

and thus

$$\text{cond}_2(A) = \frac{\sigma_1}{\sigma_n}.$$

(3) This follows from the fact that $\|A\|_2 = \rho(A)$ for a normal matrix.

(4) If A is a unitary matrix, then $A^*A = AA^* = I$, so $\rho(A^*A) = 1$, and $\|A\|_2 = \sqrt{\rho(A^*A)} = 1$. We also have $\|A^{-1}\|_2 = \|A^*\|_2 = \sqrt{\rho(AA^*)} = 1$, and thus $\text{cond}(A) = 1$.

(5) This follows immediately from the unitary invariance of the spectral norm. \square

Proposition 7.13 (4) shows that unitary and orthogonal transformations are very well-conditioned, and part (5) shows that unitary transformations preserve the condition number.

In order to compute $\text{cond}_2(A)$, we need to compute the top and bottom singular values of A , which may be hard. The inequality

$$\|A\|_2 \leq \|A\|_F \leq \sqrt{n} \|A\|_2,$$

may be useful in getting an approximation of $\text{cond}_2(A) = \|A\|_2 \|A^{-1}\|_2$, if A^{-1} can be determined.

Remark: There is an interesting geometric characterization of $\text{cond}_2(A)$. If $\theta(A)$ denotes the least angle between the vectors Au and Av as u and v range over all pairs of orthonormal vectors, then it can be shown that

$$\text{cond}_2(A) = \cot(\theta(A)/2).$$

Thus, if A is nearly singular, then there will be some orthonormal pair u, v such that Au and Av are nearly parallel; the angle $\theta(A)$ will be small and $\cot(\theta(A)/2)$ will be large. For more details, see Horn and Johnson [57] (Section 5.8 and Section 7.4).

It should be noted that in general (if A is not a normal matrix) a matrix could have a very large condition number even if all its eigenvalues are identical! For example, if we consider the $n \times n$ matrix

$$A = \begin{pmatrix} 1 & 2 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & 2 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & 2 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 1 & 2 & 0 \\ 0 & 0 & \dots & 0 & 0 & 1 & 2 \\ 0 & 0 & \dots & 0 & 0 & 0 & 1 \end{pmatrix},$$

it turns out that $\text{cond}_2(A) \geq 2^{n-1}$.

A classical example of matrix with a very large condition number is the *Hilbert matrix* $H^{(n)}$, the $n \times n$ matrix with

$$H_{ij}^{(n)} = \left(\frac{1}{i+j-1} \right).$$

For example, when $n = 5$,

$$H^{(5)} = \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \frac{1}{6} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \frac{1}{6} & \frac{1}{7} \\ \frac{1}{4} & \frac{1}{5} & \frac{1}{6} & \frac{1}{7} & \frac{1}{8} \\ \frac{1}{5} & \frac{1}{6} & \frac{1}{7} & \frac{1}{8} & \frac{1}{9} \end{pmatrix}.$$

It can be shown that

$$\text{cond}_2(H^{(5)}) \approx 4.77 \times 10^5.$$

Hilbert introduced these matrices in 1894 while studying a problem in approximation theory. The Hilbert matrix $H^{(n)}$ is symmetric positive definite. A closed-form formula can be given for its determinant (it is a special form of the so-called *Cauchy determinant*). The inverse of $H^{(n)}$ can also be computed explicitly! It can be shown that

$$\text{cond}_2(H^{(n)}) = O((1 + \sqrt{2})^{4n} / \sqrt{n}).$$

Going back to our matrix

$$A = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix},$$

which is a symmetric, positive, definite, matrix, it can be shown that its eigenvalues, which in this case are also its singular values because A is SPD, are

$$\lambda_1 \approx 30.2887 > \lambda_2 \approx 3.858 > \lambda_3 \approx 0.8431 > \lambda_4 \approx 0.01015,$$

so that

$$\text{cond}_2(A) = \frac{\lambda_1}{\lambda_4} \approx 2984.$$

The reader should check that for the perturbation of the right-hand side b used earlier, the relative errors $\|\delta x\|/\|x\|$ and $\|\delta x\|/\|x\|$ satisfy the inequality

$$\frac{\|\delta x\|}{\|x\|} \leq \text{cond}(A) \frac{\|\delta b\|}{\|b\|}$$

and comes close to equality.

7.4 An Application of Norms: Solving Inconsistent Linear Systems

The problem of solving an inconsistent linear system $Ax = b$ often arises in practice. This is a system where b does not belong to the column space of A , usually with more equations than variables. Thus, such a system has no solution. Yet, we would still like to “solve” such a system, at least approximately.

Such systems often arise when trying to fit some data. For example, we may have a set of 3D data points

$$\{p_1, \dots, p_n\},$$

and we have reason to believe that these points are nearly coplanar. We would like to find a plane that best fits our data points. Recall that the equation of a plane is

$$\alpha x + \beta y + \gamma z + \delta = 0,$$

with $(\alpha, \beta, \gamma) \neq (0, 0, 0)$. Thus, every plane is either not parallel to the x -axis ($\alpha \neq 0$) or not parallel to the y -axis ($\beta \neq 0$) or not parallel to the z -axis ($\gamma \neq 0$).

Say we have reasons to believe that the plane we are looking for is not parallel to the z -axis. If we are wrong, in the least squares solution, one of the coefficients, α, β , will be very large. If $\gamma \neq 0$, then we may assume that our plane is given by an equation of the form

$$z = ax + by + d,$$

and we would like this equation to be satisfied for all the p_i 's, which leads to a system of n equations in 3 unknowns a, b, d , with $p_i = (x_i, y_i, z_i)$;

$$\begin{aligned} ax_1 + by_1 + d &= z_1 \\ &\vdots \\ ax_n + by_n + d &= z_n. \end{aligned}$$

However, if n is larger than 3, such a system generally has *no solution*. Since the above system can't be solved exactly, we can try to find a solution (a, b, d) that *minimizes the least-squares error*

$$\sum_{i=1}^n (ax_i + by_i + d - z_i)^2.$$

This is what Legendre and Gauss figured out in the early 1800's!

In general, given a linear system

$$Ax = b,$$

we solve the *least squares problem*: minimize $\|Ax - b\|_2^2$.

Fortunately, every $n \times m$ -matrix A can be written as

$$A = VDU^\top$$

where U and V are orthogonal and D is a rectangular diagonal matrix with non-negative entries (*singular value decomposition, or SVD*); see Chapter 16.

The SVD can be used to solve an inconsistent system. It is shown in Chapter 17 that there is a vector x of smallest norm minimizing $\|Ax - b\|_2$. It is given by the (Penrose) *pseudo-inverse* of A (itself given by the SVD).

It has been observed that solving in the least-squares sense may give too much weight to “outliers,” that is, points clearly outside the best-fit plane. In this case, it is preferable to minimize (the ℓ_1 -norm)

$$\sum_{i=1}^n |ax_i + by_i + d - z_i|.$$

This does not appear to be a linear problem, but we can use a trick to convert this minimization problem into a linear program (which means a problem involving linear constraints).

Note that $|x| = \max\{x, -x\}$. So, by introducing new variables e_1, \dots, e_n , our minimization problem is equivalent to the linear program (LP):

$$\begin{array}{ll} \text{minimize} & e_1 + \cdots + e_n \\ \text{subject to} & ax_i + by_i + d - z_i \leq e_i \\ & -(ax_i + by_i + d - z_i) \leq e_i \\ & 1 \leq i \leq n. \end{array}$$

Observe that the constraints are equivalent to

$$e_i \geq |ax_i + by_i + d - z_i|, \quad 1 \leq i \leq n.$$

For an optimal solution, we must have equality, since otherwise we could decrease some e_i and get an even better solution. Of course, we are no longer dealing with “pure” linear algebra, since our constraints are inequalities.

We prefer not getting into linear programming right now, but the above example provides a good reason to learn more about linear programming!

7.5 Summary

The main concepts and results of this chapter are listed below:

- *Norms* and *normed vector spaces*.
- The *triangle inequality*.
- The *Euclidean norm*; the ℓ_p -norms.
- *Hölder’s inequality*; the *Cauchy–Schwarz inequality*; *Minkowski’s inequality*.
- *Hermitian inner product* and *Euclidean inner product*.
- *Equivalent norms*.
- *All norms on a finite-dimensional vector space are equivalent* (Theorem 7.3).
- *Matrix norms*.
- *Hermitian, symmetric* and *normal matrices*. *Orthogonal* and *unitary matrices*.
- The *trace* of a matrix.
- *Eigenvalues* and *eigenvectors* of a matrix.
- The *characteristic polynomial* of a matrix.
- The *spectral radius* $\rho(A)$ of a matrix A .
- The *Frobenius norm*.
- The Frobenius norm is a *unitarily invariant* matrix norm.
- *Bounded linear maps*.
- *Subordinate matrix norms*.
- Characterization of the subordinate matrix norms for the vector norms $\|\cdot\|_1$, $\|\cdot\|_2$, and $\|\cdot\|_\infty$.

- The *spectral norm*.
- For every matrix $A \in M_n(\mathbb{C})$ and for every $\epsilon > 0$, there is some subordinate matrix norm $\| \cdot \|$ such that $\|A\| \leq \rho(A) + \epsilon$.
- *Condition numbers* of matrices.
- Perturbation analysis of linear systems.
- The *singular value decomposition* (SVD).
- Properties of conditions numbers. Characterization of $\text{cond}_2(A)$ in terms of the largest and smallest singular values of A .
- The *Hilbert matrix*: a very badly conditioned matrix.
- Solving inconsistent linear systems by the method of *least-squares*; *linear programming*.

Chapter 8

Eigenvectors and Eigenvalues

8.1 Eigenvectors and Eigenvalues of a Linear Map

Given a finite-dimensional vector space E , let $f: E \rightarrow E$ be any linear map. If, by luck, there is a basis (e_1, \dots, e_n) of E with respect to which f is represented by a *diagonal matrix*

$$D = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \lambda_n \end{pmatrix},$$

then the action of f on E is very simple; in every “direction” e_i , we have

$$f(e_i) = \lambda_i e_i.$$

We can think of f as a transformation that stretches or shrinks space along the direction e_1, \dots, e_n (at least if E is a real vector space). In terms of matrices, the above property translates into the fact that there is an invertible matrix P and a diagonal matrix D such that a matrix A can be factored as

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

When this happens, we say that f (or A) is *diagonalizable*, the λ_i s are called the *eigenvalues* of f , and the e_i s are *eigenvectors* of f . For example, we will see that every symmetric matrix can be diagonalized. Unfortunately, not every matrix can be diagonalized. For example, the matrix

$$A_1 = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$$

can't be diagonalized. Sometimes, a matrix fails to be diagonalizable because its eigenvalues do not belong to the field of coefficients, such as

$$A_2 = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix},$$

whose eigenvalues are $\pm i$. This is not a serious problem because A_2 can be diagonalized over the complex numbers. However, A_1 is a “fatal” case! Indeed, its eigenvalues are both 1 and the problem is that A_1 does not have enough eigenvectors to span E .

The next best thing is that there is a basis with respect to which f is represented by an *upper triangular* matrix. In this case we say that f can be *triangularized*, or that f is *triangularizable*. As we will see in Section 8.2, if all the eigenvalues of f belong to the field of coefficients K , then f can be triangularized. In particular, this is the case if $K = \mathbb{C}$.

Now, an alternative to triangularization is to consider the representation of f with respect to *two* bases (e_1, \dots, e_n) and (f_1, \dots, f_n) , rather than a single basis. In this case, if $K = \mathbb{R}$ or $K = \mathbb{C}$, it turns out that we can even pick these bases to be *orthonormal*, and we get a diagonal matrix Σ with *nonnegative entries*, such that

$$f(e_i) = \sigma_i f_i, \quad 1 \leq i \leq n.$$

The nonzero σ_i s are the *singular values* of f , and the corresponding representation is the *singular value decomposition*, or *SVD*. The SVD plays a very important role in applications, and will be considered in detail later.

In this section, we focus on the possibility of diagonalizing a linear map, and we introduce the relevant concepts to do so. Given a vector space E over a field K , let I denote the identity map on E .

Definition 8.1. Given any vector space E and any linear map $f: E \rightarrow E$, a scalar $\lambda \in K$ is called an *eigenvalue*, or *proper value*, or *characteristic value* of f if there is some nonzero vector $u \in E$ such that

$$f(u) = \lambda u.$$

Equivalently, λ is an eigenvalue of f if $\text{Ker}(\lambda I - f)$ is nontrivial (i.e., $\text{Ker}(\lambda I - f) \neq \{0\}$). A vector $u \in E$ is called an *eigenvector*, or *proper vector*, or *characteristic vector* of f if $u \neq 0$ and if there is some $\lambda \in K$ such that

$$f(u) = \lambda u;$$

the scalar λ is then an eigenvalue, and we say that u is an *eigenvector associated with* λ . Given any eigenvalue $\lambda \in K$, the nontrivial subspace $\text{Ker}(\lambda I - f)$ consists of all the eigenvectors associated with λ together with the zero vector; this subspace is denoted by $E_\lambda(f)$, or $E(\lambda, f)$, or even by E_λ , and is called the *eigenspace associated with* λ , or *proper subspace associated with* λ .

Note that distinct eigenvectors may correspond to the same eigenvalue, but distinct eigenvalues correspond to disjoint sets of eigenvectors.

Remark: As we emphasized in the remark following Definition 7.4, we *require an eigenvector to be nonzero*. This requirement seems to have more benefits than inconvenients, even though

it may be considered somewhat inelegant because the set of all eigenvectors associated with an eigenvalue is not a subspace since the zero vector is excluded.

Let us now assume that E is of finite dimension n . The next proposition shows that the eigenvalues of a linear map $f: E \rightarrow E$ are the roots of a polynomial associated with f .

Proposition 8.1. *Let E be any vector space of finite dimension n and let f be any linear map $f: E \rightarrow E$. The eigenvalues of f are the roots (in K) of the polynomial*

$$\det(\lambda I - f).$$

Proof. A scalar $\lambda \in K$ is an eigenvalue of f iff there is some nonzero vector $u \neq 0$ in E such that

$$f(u) = \lambda u$$

iff

$$(\lambda I - f)(u) = 0$$

iff $(\lambda I - f)$ is not invertible iff, by Proposition 5.14,

$$\det(\lambda I - f) = 0.$$

□

In view of the importance of the polynomial $\det(\lambda I - f)$, we have the following definition.

Definition 8.2. Given any vector space E of dimension n , for any linear map $f: E \rightarrow E$, the polynomial $P_f(X) = \chi_f(X) = \det(XI - f)$ is called the *characteristic polynomial* of f . For any square matrix A , the polynomial $P_A(X) = \chi_A(X) = \det(XI - A)$ is called the *characteristic polynomial* of A .

Note that we already encountered the characteristic polynomial in Section 5.7; see Definition 5.8.

Given any basis (e_1, \dots, e_n) , if $A = M(f)$ is the matrix of f w.r.t. (e_1, \dots, e_n) , we can compute the characteristic polynomial $\chi_f(X) = \det(XI - f)$ of f by expanding the following determinant:

$$\det(XI - A) = \begin{vmatrix} X - a_{11} & -a_{12} & \dots & -a_{1n} \\ -a_{21} & X - a_{22} & \dots & -a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ -a_{n1} & -a_{n2} & \dots & X - a_{nn} \end{vmatrix}.$$

If we expand this determinant, we find that

$$\chi_A(X) = \det(XI - A) = X^n - (a_{11} + \dots + a_{nn})X^{n-1} + \dots + (-1)^n \det(A).$$

The sum $\text{tr}(A) = a_{11} + \dots + a_{nn}$ of the diagonal elements of A is called the *trace* of A . Since we proved in Section 5.7 that the characteristic polynomial only depends on the linear map f , the above shows that $\text{tr}(A)$ has the same value for all matrices A representing f . Thus,

the *trace of a linear map* is well-defined; we have $\text{tr}(f) = \text{tr}(A)$ for any matrix A representing f .

Remark: The characteristic polynomial of a linear map is sometimes defined as $\det(f - XI)$. Since

$$\det(f - XI) = (-1)^n \det(XI - f),$$

this makes essentially no difference but the version $\det(XI - f)$ has the small advantage that the coefficient of X^n is $+1$.

If we write

$$\chi_A(X) = \det(XI - A) = X^n - \tau_1(A)X^{n-1} + \cdots + (-1)^k \tau_k(A)X^{n-k} + \cdots + (-1)^n \tau_n(A),$$

then we just proved that

$$\tau_1(A) = \text{tr}(A) \quad \text{and} \quad \tau_n(A) = \det(A).$$

It is also possible to express $\tau_k(A)$ in terms of determinants of certain submatrices of A . For any nonempty subset, $I \subseteq \{1, \dots, n\}$, say $I = \{i_1 < \dots < i_k\}$, let $A_{I,I}$ be the $k \times k$ submatrix of A whose j th column consists of the elements $a_{i_h i_j}$, where $h = 1, \dots, k$. Equivalently, $A_{I,I}$ is the matrix obtained from A by first selecting the columns whose indices belong to I , and then the rows whose indices also belong to I . Then, it can be shown that

$$\tau_k(A) = \sum_{\substack{I \subseteq \{1, \dots, n\} \\ |I|=k}} \det(A_{I,I}).$$

If all the roots, $\lambda_1, \dots, \lambda_n$, of the polynomial $\det(XI - A)$ belong to the field K , then we can write

$$\chi_A(X) = \det(XI - A) = (X - \lambda_1) \cdots (X - \lambda_n),$$

where some of the λ_i s may appear more than once. Consequently,

$$\chi_A(X) = \det(XI - A) = X^n - \sigma_1(\lambda)X^{n-1} + \cdots + (-1)^k \sigma_k(\lambda)X^{n-k} + \cdots + (-1)^n \sigma_n(\lambda),$$

where

$$\sigma_k(\lambda) = \sum_{\substack{I \subseteq \{1, \dots, n\} \\ |I|=k}} \prod_{i \in I} \lambda_i,$$

the k th symmetric function of the λ_i 's. From this, it is clear that

$$\sigma_k(\lambda) = \tau_k(A)$$

and, in particular, the product of the eigenvalues of f is equal to $\det(A) = \det(f)$, and the sum of the eigenvalues of f is equal to the trace $\text{tr}(A) = \text{tr}(f)$, of f ; for the record,

$$\begin{aligned} \text{tr}(f) &= \lambda_1 + \cdots + \lambda_n \\ \det(f) &= \lambda_1 \cdots \lambda_n, \end{aligned}$$

where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of f (and A), where some of the λ_i s may appear more than once. In particular, f is not invertible iff it admits 0 as an eigenvalue.

Remark: Depending on the field K , the characteristic polynomial $\chi_A(X) = \det(XI - A)$ may or may not have roots in K . This motivates considering *algebraically closed fields*, which are fields K such that every polynomial with coefficients in K has all its roots in K . For example, over $K = \mathbb{R}$, not every polynomial has real roots. If we consider the matrix

$$A = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix},$$

then the characteristic polynomial $\det(XI - A)$ has no real roots unless $\theta = k\pi$. However, over the field \mathbb{C} of complex numbers, every polynomial has roots. For example, the matrix above has the roots $\cos \theta \pm i \sin \theta = e^{\pm i\theta}$.

It is possible to show that every linear map f over a complex vector space E must have some (complex) eigenvalue without having recourse to determinants (and the characteristic polynomial). Let $n = \dim(E)$, pick any nonzero vector $u \in E$, and consider the sequence

$$u, f(u), f^2(u), \dots, f^n(u).$$

Since the above sequence has $n + 1$ vectors and E has dimension n , these vectors must be linearly dependent, so there are some complex numbers c_0, \dots, c_m , not all zero, such that

$$c_0 f^m(u) + c_1 f^{m-1}(u) + \dots + c_m u = 0,$$

where $m \leq n$ is the largest integer such that the coefficient of $f^m(u)$ is nonzero (m must exist since we have a nontrivial linear dependency). Now, because the field \mathbb{C} is algebraically closed, the polynomial

$$c_0 X^m + c_1 X^{m-1} + \dots + c_m$$

can be written as a product of linear factors as

$$c_0 X^m + c_1 X^{m-1} + \dots + c_m = c_0 (X - \lambda_1) \cdots (X - \lambda_m)$$

for some complex numbers $\lambda_1, \dots, \lambda_m \in \mathbb{C}$, not necessarily distinct. But then, since $c_0 \neq 0$,

$$c_0 f^m(u) + c_1 f^{m-1}(u) + \dots + c_m u = 0$$

is equivalent to

$$(f - \lambda_1 I) \circ \dots \circ (f - \lambda_m I)(u) = 0.$$

If all the linear maps $f - \lambda_i I$ were injective, then $(f - \lambda_1 I) \circ \dots \circ (f - \lambda_m I)$ would be injective, contradicting the fact that $u \neq 0$. Therefore, some linear map $f - \lambda_i I$ must have a nontrivial kernel, which means that there is some $v \neq 0$ so that

$$f(v) = \lambda_i v;$$

that is, λ_i is some eigenvalue of f and v is some eigenvector of f .

As nice as the above argument is, it does not provide a method for *finding* the eigenvalues of f , and even if we prefer avoiding determinants as much as possible, we are forced to deal with the characteristic polynomial $\det(XI - f)$.

Definition 8.3. Let A be an $n \times n$ matrix over a field K . Assume that all the roots of the characteristic polynomial $\chi_A(X) = \det(XI - A)$ of A belong to K , which means that we can write

$$\det(XI - A) = (X - \lambda_1)^{k_1} \cdots (X - \lambda_m)^{k_m},$$

where $\lambda_1, \dots, \lambda_m \in K$ are the distinct roots of $\det(XI - A)$ and $k_1 + \cdots + k_m = n$. The integer k_i is called the *algebraic multiplicity* of the eigenvalue λ_i , and the dimension of the eigenspace $E_{\lambda_i} = \text{Ker}(\lambda_i I - A)$ is called the *geometric multiplicity* of λ_i . We denote the algebraic multiplicity of λ_i by $\text{alg}(\lambda_i)$, and its geometric multiplicity by $\text{geo}(\lambda_i)$.

By definition, the sum of the algebraic multiplicities is equal to n , but the sum of the geometric multiplicities can be strictly smaller.

Proposition 8.2. Let A be an $n \times n$ matrix over a field K and assume that all the roots of the characteristic polynomial $\chi_A(X) = \det(XI - A)$ of A belong to K . For every eigenvalue λ_i of A , the geometric multiplicity of λ_i is always less than or equal to its algebraic multiplicity, that is,

$$\text{geo}(\lambda_i) \leq \text{alg}(\lambda_i).$$

Proof. To see this, if n_i is the dimension of the eigenspace E_{λ_i} associated with the eigenvalue λ_i , we can form a basis of K^n obtained by picking a basis of E_{λ_i} and completing this linearly independent family to a basis of K^n . With respect to this new basis, our matrix is of the form

$$A' = \begin{pmatrix} \lambda_i I_{n_i} & B \\ 0 & D \end{pmatrix}$$

and a simple determinant calculation shows that

$$\det(XI - A) = \det(XI - A') = (X - \lambda_i)^{n_i} \det(XI_{n-n_i} - D).$$

Therefore, $(X - \lambda_i)^{n_i}$ divides the characteristic polynomial of A' , and thus, the characteristic polynomial of A . It follows that n_i is less than or equal to the algebraic multiplicity of λ_i . \square

The following proposition shows an interesting property of eigenspaces.

Proposition 8.3. Let E be any vector space of finite dimension n and let f be any linear map. If u_1, \dots, u_m are eigenvectors associated with pairwise distinct eigenvalues $\lambda_1, \dots, \lambda_m$, then the family (u_1, \dots, u_m) is linearly independent.

Proof. Assume that (u_1, \dots, u_m) is linearly dependent. Then, there exists $\mu_1, \dots, \mu_k \in K$ such that

$$\mu_1 u_{i_1} + \dots + \mu_k u_{i_k} = 0,$$

where $1 \leq k \leq m$, $\mu_i \neq 0$ for all i , $1 \leq i \leq k$, $\{i_1, \dots, i_k\} \subseteq \{1, \dots, m\}$, and no proper subfamily of $(u_{i_1}, \dots, u_{i_k})$ is linearly dependent (in other words, we consider a dependency relation with k minimal). Applying f to this dependency relation, we get

$$\mu_1 \lambda_{i_1} u_{i_1} + \dots + \mu_k \lambda_{i_k} u_{i_k} = 0,$$

and if we multiply the original dependency relation by λ_{i_1} and subtract it from the above, we get

$$\mu_2 (\lambda_{i_2} - \lambda_{i_1}) u_{i_2} + \dots + \mu_k (\lambda_{i_k} - \lambda_{i_1}) u_{i_k} = 0,$$

which is a nontrivial linear dependency among a proper subfamily of $(u_{i_1}, \dots, u_{i_k})$ since the λ_j are all distinct and the μ_i are nonzero, a contradiction. \square

Thus, from Proposition 8.3, if $\lambda_1, \dots, \lambda_m$ are all the pairwise distinct eigenvalues of f (where $m \leq n$), we have a direct sum

$$E_{\lambda_1} \oplus \dots \oplus E_{\lambda_m}$$

of the eigenspaces E_{λ_i} . Unfortunately, it is not always the case that

$$E = E_{\lambda_1} \oplus \dots \oplus E_{\lambda_m}.$$

When

$$E = E_{\lambda_1} \oplus \dots \oplus E_{\lambda_m},$$

we say that f is *diagonalizable* (and similarly for any matrix associated with f). Indeed, picking a basis in each E_{λ_i} , we obtain a matrix which is a diagonal matrix consisting of the eigenvalues, each λ_i occurring a number of times equal to the dimension of E_{λ_i} . This happens if the algebraic multiplicity and the geometric multiplicity of every eigenvalue are equal. In particular, when the characteristic polynomial has n distinct roots, then f is diagonalizable. It can also be shown that symmetric matrices have real eigenvalues and can be diagonalized.

For a negative example, we leave as exercise to show that the matrix

$$M = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$$

cannot be diagonalized, even though 1 is an eigenvalue. The problem is that the eigenspace of 1 only has dimension 1. The matrix

$$A = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$$

cannot be diagonalized either, because it has no real eigenvalues, unless $\theta = k\pi$. However, over the field of complex numbers, it can be diagonalized.

8.2 Reduction to Upper Triangular Form

Unfortunately, not every linear map on a complex vector space can be diagonalized. The next best thing is to “triangularize,” which means to find a basis over which the matrix has zero entries below the main diagonal. Fortunately, such a basis always exist.

We say that a square matrix A is an *upper triangular matrix* if it has the following shape,

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1n-1} & a_{1n} \\ 0 & a_{22} & a_{23} & \cdots & a_{2n-1} & a_{2n} \\ 0 & 0 & a_{33} & \cdots & a_{3n-1} & a_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & a_{n-1n-1} & a_{n-1n} \\ 0 & 0 & 0 & \cdots & 0 & a_{nn} \end{pmatrix},$$

i.e., $a_{ij} = 0$ whenever $j < i$, $1 \leq i, j \leq n$.

Theorem 8.4. *Given any finite dimensional vector space over a field K , for any linear map $f: E \rightarrow E$, there is a basis (u_1, \dots, u_n) with respect to which f is represented by an upper triangular matrix (in $M_n(K)$) iff all the eigenvalues of f belong to K . Equivalently, for every $n \times n$ matrix $A \in M_n(K)$, there is an invertible matrix P and an upper triangular matrix T (both in $M_n(K)$) such that*

$$A = PTP^{-1}$$

iff all the eigenvalues of A belong to K .

Proof. If there is a basis (u_1, \dots, u_n) with respect to which f is represented by an upper triangular matrix T in $M_n(K)$, then since the eigenvalues of f are the diagonal entries of T , all the eigenvalues of f belong to K .

For the converse, we proceed by induction on the dimension n of E . For $n = 1$ the result is obvious. If $n > 1$, since by assumption f has all its eigenvalue in K , pick some eigenvalue $\lambda_1 \in K$ of f , and let u_1 be some corresponding (nonzero) eigenvector. We can find $n - 1$ vectors (v_2, \dots, v_n) such that (u_1, v_2, \dots, v_n) is a basis of E , and let F be the subspace of dimension $n - 1$ spanned by (v_2, \dots, v_n) . In the basis (u_1, v_2, \dots, v_n) , the matrix of f is of the form

$$U = \begin{pmatrix} \lambda_1 & a_{12} & \cdots & a_{1n} \\ 0 & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{n2} & \cdots & a_{nn} \end{pmatrix},$$

since its first column contains the coordinates of $\lambda_1 u_1$ over the basis (u_1, v_2, \dots, v_n) . If we let $p: E \rightarrow F$ be the projection defined such that $p(u_1) = 0$ and $p(v_i) = v_i$ when $2 \leq i \leq n$, the linear map $g: F \rightarrow F$ defined as the restriction of $p \circ f$ to F is represented by the $(n - 1) \times (n - 1)$ matrix $V = (a_{ij})_{2 \leq i, j \leq n}$ over the basis (v_2, \dots, v_n) . We need to prove

that all the eigenvalues of g belong to K . However, since the first column of U has a single nonzero entry, we get

$$\chi_U(X) = \det(XI - U) = (X - \lambda_1) \det(XI - V) = (X - \lambda_1) \chi_V(X),$$

where $\chi_U(X)$ is the characteristic polynomial of U and $\chi_V(X)$ is the characteristic polynomial of V . It follows that $\chi_V(X)$ divides $\chi_U(X)$, and since all the roots of $\chi_U(X)$ are in K , all the roots of $\chi_V(X)$ are also in K . Consequently, we can apply the induction hypothesis, and there is a basis (u_2, \dots, u_n) of F such that g is represented by an upper triangular matrix $(b_{ij})_{1 \leq i, j \leq n-1}$. However,

$$E = Ku_1 \oplus F,$$

and thus (u_1, \dots, u_n) is a basis for E . Since p is the projection from $E = Ku_1 \oplus F$ onto F and $g: F \rightarrow F$ is the restriction of $p \circ f$ to F , we have

$$f(u_1) = \lambda_1 u_1$$

and

$$f(u_{i+1}) = a_{1i} u_1 + \sum_{j=1}^i b_{ij} u_{j+1}$$

for some $a_{1i} \in K$, when $1 \leq i \leq n-1$. But then the matrix of f with respect to (u_1, \dots, u_n) is upper triangular.

For the matrix version, we assume that A is the matrix of f with respect to some basis. Then, we just proved that there is a change of basis matrix P such that $A = PTP^{-1}$ where T is upper triangular. \square

If $A = PTP^{-1}$ where T is upper triangular, note that the diagonal entries of T are the eigenvalues $\lambda_1, \dots, \lambda_n$ of A . Indeed, A and T have the same characteristic polynomial. Also, if A is a real matrix whose eigenvalues are all real, then P can be chosen to real, and if A is a rational matrix whose eigenvalues are all rational, then P can be chosen rational. Since any polynomial over \mathbb{C} has all its roots in \mathbb{C} , Theorem 8.4 implies that every complex $n \times n$ matrix can be triangularized.

If λ is an eigenvalue of the matrix A and if u is an eigenvector associated with λ , from

$$Au = \lambda u,$$

we obtain

$$A^2 u = A(Au) = A(\lambda u) = \lambda Au = \lambda^2 u,$$

which shows that λ^2 is an eigenvalue of A^2 for the eigenvector u . An obvious induction shows that λ^k is an eigenvalue of A^k for the eigenvector u , for all $k \geq 1$. Now, if all eigenvalues $\lambda_1, \dots, \lambda_n$ of A are in K , it follows that $\lambda_1^k, \dots, \lambda_n^k$ are eigenvalues of A^k . However, it is not obvious that A^k does not have other eigenvalues. In fact, this can't happen, and this can be proved using Theorem 8.4.

Proposition 8.5. *Given any $n \times n$ matrix $A \in M_n(K)$ with coefficients in a field K , if all eigenvalues $\lambda_1, \dots, \lambda_n$ of A are in K , then for every polynomial $q(X) \in K[X]$, the eigenvalues of $q(A)$ are exactly $(q(\lambda_1), \dots, q(\lambda_n))$.*

Proof. By Theorem 8.4, there is an upper triangular matrix T and an invertible matrix P (both in $M_n(K)$) such that

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

Since A and T are similar, they have the same eigenvalues (with the same multiplicities), so the diagonal entries of T are the eigenvalues of A . Since

$$A^k = PT^kP^{-1}, \quad k \geq 1,$$

for any polynomial $q(X) = c_0X^m + \dots + c_{m-1}X + c_m$, we have

$$\begin{aligned} q(A) &= c_0A^m + \dots + c_{m-1}A + c_mI \\ &= c_0PT^mP^{-1} + \dots + c_{m-1}PTP^{-1} + c_mPIP^{-1} \\ &= P(c_0T^m + \dots + c_{m-1}T + c_mI)P^{-1} \\ &= Pq(T)P^{-1}. \end{aligned}$$

Furthermore, it is easy to check that $q(T)$ is upper triangular and that its diagonal entries are $q(\lambda_1), \dots, q(\lambda_n)$, where $\lambda_1, \dots, \lambda_n$ are the diagonal entries of T , namely the eigenvalues of A . It follows that $q(\lambda_1), \dots, q(\lambda_n)$ are the eigenvalues of $q(A)$. \square

If E is a Hermitian space (see Chapter 12), the proof of Theorem 8.4 can be easily adapted to prove that there is an *orthonormal* basis (u_1, \dots, u_n) with respect to which the matrix of f is upper triangular. This is usually known as *Schur's lemma*.

Theorem 8.6. (*Schur decomposition*) *Given any linear map $f: E \rightarrow E$ over a complex Hermitian space E , there is an orthonormal basis (u_1, \dots, u_n) with respect to which f is represented by an upper triangular matrix. Equivalently, for every $n \times n$ matrix $A \in M_n(\mathbb{C})$, there is a unitary matrix U and an upper triangular matrix T such that*

$$A = UTU^*.$$

If A is real and if all its eigenvalues are real, then there is an orthogonal matrix Q and a real upper triangular matrix T such that

$$A = QTQ^\top.$$

Proof. During the induction, we choose F to be the orthogonal complement of $\mathbb{C}u_1$ and we pick orthonormal bases (use Propositions 12.10 and 12.9). If E is a real Euclidean space and if the eigenvalues of f are all real, the proof also goes through with real matrices (use Propositions 10.9 and 10.8). \square

Using Theorem 8.6, we can derive the fact that if A is a Hermitian matrix, then there is a unitary matrix U and a real diagonal matrix D such that $A = UDU^*$. Indeed, since $A^* = A$, we get

$$UTU^* = UT^*U^*,$$

which implies that $T = T^*$. Since T is an upper triangular matrix, T^* is a lower triangular matrix, which implies that T is a real diagonal matrix. In fact, applying this result to a (real) symmetric matrix A , we obtain the fact that all the eigenvalues of a symmetric matrix are real, and by applying Theorem 8.6 again, we conclude that $A = QDQ^T$, where Q is orthogonal and D is a real diagonal matrix. We will also prove this in Chapter 13.

When A has complex eigenvalues, there is a version of Theorem 8.6 involving only real matrices provided that we allow T to be block upper-triangular (the diagonal entries may be 2×2 matrices or real entries).

Theorem 8.6 is not a very practical result but it is a useful theoretical result to cope with matrices that cannot be diagonalized. For example, it can be used to prove that *every* complex matrix is the limit of a sequence of diagonalizable matrices that have distinct eigenvalues!

Remark: There is another way to prove Proposition 8.5 that does not use Theorem 8.4, but instead uses the fact that given any field K , there is field extension \overline{K} of K ($K \subseteq \overline{K}$) such that every polynomial $q(X) = c_0X^m + \cdots + c_{m-1}X + c_m$ (of degree $m \geq 1$) with coefficients $c_i \in K$ factors as

$$q(X) = c_0(X - \alpha_1) \cdots (X - \alpha_n), \quad \alpha_i \in \overline{K}, i = 1, \dots, n.$$

The field \overline{K} is called an *algebraically closed field* (and an algebraic closure of K).

Assume that all eigenvalues $\lambda_1, \dots, \lambda_n$ of A belong to K . Let $q(X)$ be any polynomial (in $K[X]$) and let $\mu \in \overline{K}$ be any eigenvalue of $q(A)$ (this means that μ is a zero of the characteristic polynomial $\chi_{q(A)}(X) \in K[X]$ of $q(A)$). Since \overline{K} is algebraically closed, $\chi_{q(A)}(X)$ has all its root in \overline{K} . We claim that $\mu = q(\lambda_i)$ for some eigenvalue λ_i of A .

Proof. (After Lax [71], Chapter 6). Since \overline{K} is algebraically closed, the polynomial $\mu - q(X)$ factors as

$$\mu - q(X) = c_0(X - \alpha_1) \cdots (X - \alpha_n),$$

for some $\alpha_i \in \overline{K}$. Now, $\mu I - q(A)$ is a matrix in $M_n(\overline{K})$, and since μ is an eigenvalue of $q(A)$, it must be singular. We have

$$\mu I - q(A) = c_0(A - \alpha_1 I) \cdots (A - \alpha_n I),$$

and since the left-hand side is singular, so is the right-hand side, which implies that some factor $A - \alpha_i I$ is singular. This means that α_i is an eigenvalue of A , say $\alpha_i = \lambda_i$. As $\alpha_i = \lambda_i$ is a zero of $\mu - q(X)$, we get

$$\mu = q(\lambda_i),$$

which proves that μ is indeed of the form $q(\lambda_i)$ for some eigenvalue λ_i of A . \square

8.3 Location of Eigenvalues

If A is an $n \times n$ complex (or real) matrix A , it would be useful to know, even roughly, where the eigenvalues of A are located in the complex plane \mathbb{C} . The Gershgorin discs provide some precise information about this.

Definition 8.4. For any complex $n \times n$ matrix A , for $i = 1, \dots, n$, let

$$R'_i(A) = \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|$$

and let

$$G(A) = \bigcup_{i=1}^n \{z \in \mathbb{C} \mid |z - a_{ii}| \leq R'_i(A)\}.$$

Each disc $\{z \in \mathbb{C} \mid |z - a_{ii}| \leq R'_i(A)\}$ is called a *Gershgorin disc* and their union $G(A)$ is called the *Gershgorin domain*.

Although easy to prove, the following theorem is very useful:

Theorem 8.7. (*Gershgorin's disc theorem*) For any complex $n \times n$ matrix A , all the eigenvalues of A belong to the Gershgorin domain $G(A)$. Furthermore the following properties hold:

(1) If A is strictly row diagonally dominant, that is

$$|a_{ii}| > \sum_{j=1, j \neq i}^n |a_{ij}|, \quad \text{for } i = 1, \dots, n,$$

then A is invertible.

(2) If A is strictly row diagonally dominant, and if $a_{ii} > 0$ for $i = 1, \dots, n$, then every eigenvalue of A has a strictly positive real part.

Proof. Let λ be any eigenvalue of A and let u be a corresponding eigenvector (recall that we must have $u \neq 0$). Let k be an index such that

$$|u_k| = \max_{1 \leq i \leq n} |u_i|.$$

Since $Au = \lambda u$, we have

$$(\lambda - a_{kk})u_k = \sum_{\substack{j=1 \\ j \neq k}}^n a_{kj}u_j,$$

which implies that

$$|\lambda - a_{kk}| |u_k| \leq \sum_{\substack{j=1 \\ j \neq k}}^n |a_{kj}| |u_j| \leq |u_k| \sum_{\substack{j=1 \\ j \neq k}}^n |a_{kj}|$$

and since $u \neq 0$ and $|u_k| = \max_{1 \leq i \leq n} |u_i|$, we must have $|u_k| \neq 0$, and it follows that

$$|\lambda - a_{kk}| \leq \sum_{\substack{j=1 \\ j \neq k}}^n |a_{kj}| = R'_k(A),$$

and thus

$$\lambda \in \{z \in \mathbb{C} \mid |z - a_{kk}| \leq R'_k(A)\} \subseteq G(A),$$

as claimed.

(1) Strict row diagonal dominance implies that 0 does not belong to any of the Gershgorin discs, so all eigenvalues of A are nonzero, and A is invertible.

(2) If A is strictly row diagonally dominant and $a_{ii} > 0$ for $i = 1, \dots, n$, then each of the Gershgorin discs lies strictly in the right half-plane, so every eigenvalue of A has a strictly positive real part. \square

In particular, Theorem 8.7 implies that if a symmetric matrix is strictly row diagonally dominant and has strictly positive diagonal entries, then it is positive definite. Theorem 8.7 is sometimes called the *Gershgorin–Hadamard theorem*.

Since A and A^\top have the same eigenvalues (even for complex matrices) we also have a version of Theorem 8.7 for the discs of radius

$$C'_j(A) = \sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}|,$$

whose domain is denoted by $G(A^\top)$. Thus we get the following:

Theorem 8.8. *For any complex $n \times n$ matrix A , all the eigenvalues of A belong to the intersection of the Gershgorin domains, $G(A) \cap G(A^\top)$. Furthermore the following properties hold:*

(1) *If A is strictly column diagonally dominant, that is*

$$|a_{ii}| > \sum_{i=1, i \neq j}^n |a_{ij}|, \quad \text{for } j = 1, \dots, n,$$

then A is invertible.

- (2) If A is strictly column diagonally dominant, and if $a_{ii} > 0$ for $i = 1, \dots, n$, then every eigenvalue of A has a strictly positive real part.

There are refinements of Gershgorin's theorem and eigenvalue location results involving other domains besides discs; for more on this subject, see Horn and Johnson [57], Sections 6.1 and 6.2.

Remark: Neither strict row diagonal dominance nor strict column diagonal dominance are necessary for invertibility. Also, if we relax all strict inequalities to inequalities, then row diagonal dominance (or column diagonal dominance) is not a sufficient condition for invertibility.

8.4 Summary

The main concepts and results of this chapter are listed below:

- *Diagonal matrix.*
- *Eigenvalues, eigenvectors; the eigenspace associated with an eigenvalue.*
- *The characteristic polynomial.*
- *The trace.*
- *algebraic and geometric multiplicity.*
- Eigenspaces associated with distinct eigenvalues form a direct sum (Proposition 8.3).
- Reduction of a matrix to an upper-triangular matrix.
- *Schur decomposition.*
- The *Gershgorin's discs* can be used to locate the eigenvalues of a complex matrix; see Theorems 8.7 and 8.8.

Chapter 9

Iterative Methods for Solving Linear Systems

9.1 Convergence of Sequences of Vectors and Matrices

In Chapter 6 we have discussed some of the main methods for solving systems of linear equations. These methods are *direct methods*, in the sense that they yield exact solutions (assuming infinite precision!).

Another class of methods for solving linear systems consists in approximating solutions using *iterative methods*. The basic idea is this: Given a linear system $Ax = b$ (with A a square invertible matrix), find another matrix B and a vector c , such that

1. The matrix $I - B$ is invertible
2. The unique solution \tilde{x} of the system $Ax = b$ is identical to the unique solution \tilde{u} of the system

$$u = Bu + c,$$

and then, starting from any vector u_0 , compute the sequence (u_k) given by

$$u_{k+1} = Bu_k + c, \quad k \in \mathbb{N}.$$

Under certain conditions (to be clarified soon), the sequence (u_k) converges to a limit \tilde{u} which is the unique solution of $u = Bu + c$, and thus of $Ax = b$.

Consequently, it is important to find conditions that ensure the convergence of the above sequences and to have tools to compare the “rate” of convergence of these sequences. Thus, we begin with some general results about the convergence of sequences of vectors and matrices.

Let $(E, \|\cdot\|)$ be a normed vector space. Recall that a sequence (u_k) of vectors $u_k \in E$ *converges to a limit* $u \in E$, if for every $\epsilon > 0$, there some natural number N such that

$$\|u_k - u\| \leq \epsilon, \quad \text{for all } k \geq N.$$

We write

$$u = \lim_{k \rightarrow \infty} u_k.$$

If E is a finite-dimensional vector space and $\dim(E) = n$, we know from Theorem 7.3 that any two norms are equivalent, and if we choose the norm $\|\cdot\|_\infty$, we see that the convergence of the sequence of vectors u_k is equivalent to the convergence of the n sequences of scalars formed by the components of these vectors (over any basis). The same property applies to the finite-dimensional vector space $M_{m,n}(K)$ of $m \times n$ matrices (with $K = \mathbb{R}$ or $K = \mathbb{C}$), which means that the convergence of a sequence of matrices $A_k = (a_{ij}^{(k)})$ is equivalent to the convergence of the $m \times n$ sequences of scalars $(a_{ij}^{(k)})$, with i, j fixed ($1 \leq i \leq m$, $1 \leq j \leq n$).

The first theorem below gives a necessary and sufficient condition for the sequence (B^k) of powers of a matrix B to converge to the zero matrix. Recall that the spectral radius $\rho(B)$ of a matrix B is the maximum of the moduli $|\lambda_i|$ of the eigenvalues of B .

Theorem 9.1. *For any square matrix B , the following conditions are equivalent:*

- (1) $\lim_{k \rightarrow \infty} B^k = 0$,
- (2) $\lim_{k \rightarrow \infty} B^k v = 0$, for all vectors v ,
- (3) $\rho(B) < 1$,
- (4) $\|B\| < 1$, for some subordinate matrix norm $\|\cdot\|$.

Proof. Assume (1) and let $\|\cdot\|$ be a vector norm on E and $\|\cdot\|$ be the corresponding matrix norm. For every vector $v \in E$, because $\|\cdot\|$ is a matrix norm, we have

$$\|B^k v\| \leq \|B^k\| \|v\|,$$

and since $\lim_{k \rightarrow \infty} B^k = 0$ means that $\lim_{k \rightarrow \infty} \|B^k\| = 0$, we conclude that $\lim_{k \rightarrow \infty} \|B^k v\| = 0$, that is, $\lim_{k \rightarrow \infty} B^k v = 0$. This proves that (1) implies (2).

Assume (2). If we had $\rho(B) \geq 1$, then there would be some eigenvector u ($\neq 0$) and some eigenvalue λ such that

$$Bu = \lambda u, \quad |\lambda| = \rho(B) \geq 1,$$

but then the sequence $(B^k u)$ would not converge to 0, because $B^k u = \lambda^k u$ and $|\lambda^k| = |\lambda|^k \geq 1$. It follows that (2) implies (3).

Assume that (3) holds, that is, $\rho(B) < 1$. By Proposition 7.10, we can find $\epsilon > 0$ small enough that $\rho(B) + \epsilon < 1$, and a subordinate matrix norm $\|\cdot\|$ such that

$$\|B\| \leq \rho(B) + \epsilon,$$

which is (4).

Finally, assume (4). Because $\|\cdot\|$ is a matrix norm,

$$\|B^k\| \leq \|B\|^k,$$

and since $\|B\| < 1$, we deduce that (1) holds. \square

The following proposition is needed to study the rate of convergence of iterative methods.

Proposition 9.2. *For every square matrix B and every matrix norm $\|\cdot\|$, we have*

$$\lim_{k \rightarrow \infty} \|B^k\|^{1/k} = \rho(B).$$

Proof. We know from Proposition 7.4 that $\rho(B) \leq \|B\|$, and since $\rho(B) = (\rho(B^k))^{1/k}$, we deduce that

$$\rho(B) \leq \|B^k\|^{1/k} \quad \text{for all } k \geq 1,$$

and so

$$\rho(B) \leq \lim_{k \rightarrow \infty} \|B^k\|^{1/k}.$$

Now, let us prove that for every $\epsilon > 0$, there is some integer $N(\epsilon)$ such that

$$\|B^k\|^{1/k} \leq \rho(B) + \epsilon \quad \text{for all } k \geq N(\epsilon),$$

which proves that

$$\lim_{k \rightarrow \infty} \|B^k\|^{1/k} \leq \rho(B),$$

and our proposition.

For any given $\epsilon > 0$, let B_ϵ be the matrix

$$B_\epsilon = \frac{B}{\rho(B) + \epsilon}.$$

Since $\|B_\epsilon\| < 1$, Theorem 9.1 implies that $\lim_{k \rightarrow \infty} B_\epsilon^k = 0$. Consequently, there is some integer $N(\epsilon)$ such that for all $k \geq N(\epsilon)$, we have

$$\|B^k\| = \frac{\|B^k\|}{(\rho(B) + \epsilon)^k} \leq 1,$$

which implies that

$$\|B^k\|^{1/k} \leq \rho(B) + \epsilon,$$

as claimed. \square

We now apply the above results to the convergence of iterative methods.

9.2 Convergence of Iterative Methods

Recall that iterative methods for solving a linear system $Ax = b$ (with A invertible) consists in finding some matrix B and some vector c , such that $I - B$ is invertible, and the unique solution \tilde{x} of $Ax = b$ is equal to the unique solution \tilde{u} of $u = Bu + c$. Then, starting from any vector u_0 , compute the sequence (u_k) given by

$$u_{k+1} = Bu_k + c, \quad k \in \mathbb{N},$$

and say that the iterative method is *convergent* iff

$$\lim_{k \rightarrow \infty} u_k = \tilde{u},$$

for *every* initial vector u_0 .

Here is a fundamental criterion for the convergence of any iterative methods based on a matrix B , called the *matrix of the iterative method*.

Theorem 9.3. *Given a system $u = Bu + c$ as above, where $I - B$ is invertible, the following statements are equivalent:*

- (1) *The iterative method is convergent.*
- (2) $\rho(B) < 1$.
- (3) $\|B\| < 1$, for some subordinate matrix norm $\|\cdot\|$.

Proof. Define the vector e_k (error vector) by

$$e_k = u_k - \tilde{u},$$

where \tilde{u} is the unique solution of the system $u = Bu + c$. Clearly, the iterative method is convergent iff

$$\lim_{k \rightarrow \infty} e_k = 0.$$

We claim that

$$e_k = B^k e_0, \quad k \geq 0,$$

where $e_0 = u_0 - \tilde{u}$.

This is proved by induction on k . The base case $k = 0$ is trivial. By the induction hypothesis, $e_k = B^k e_0$, and since $u_{k+1} = Bu_k + c$, we get

$$u_{k+1} - \tilde{u} = Bu_k + c - \tilde{u},$$

and because $\tilde{u} = B\tilde{u} + c$ and $e_k = B^k e_0$ (by the induction hypothesis), we obtain

$$u_{k+1} - \tilde{u} = Bu_k - B\tilde{u} = B(u_k - \tilde{u}) = Be_k = BB^k e_0 = B^{k+1} e_0,$$

proving the induction step. Thus, the iterative method converges iff

$$\lim_{k \rightarrow \infty} B^k e_0 = 0.$$

Consequently, our theorem follows by Theorem 9.1. □

The next proposition is needed to compare the rate of convergence of iterative methods. It shows that *asymptotically, the error vector $e_k = B^k e_0$ behaves at worst like $(\rho(B))^k$.*

Proposition 9.4. *Let $\|\cdot\|$ be any vector norm, let B be a matrix such that $I - B$ is invertible, and let \tilde{u} be the unique solution of $u = Bu + c$.*

(1) *If (u_k) is any sequence defined iteratively by*

$$u_{k+1} = Bu_k + c, \quad k \in \mathbb{N},$$

then

$$\lim_{k \rightarrow \infty} \left[\sup_{\|u_0 - \tilde{u}\|=1} \|u_k - \tilde{u}\|^{1/k} \right] = \rho(B).$$

(2) *Let B_1 and B_2 be two matrices such that $I - B_1$ and $I - B_2$ are invertible, assume that both $u = B_1 u + c_1$ and $u = B_2 u + c_2$ have the same unique solution \tilde{u} , and consider any two sequences (u_k) and (v_k) defined inductively by*

$$\begin{aligned} u_{k+1} &= B_1 u_k + c_1 \\ v_{k+1} &= B_2 v_k + c_2, \end{aligned}$$

with $u_0 = v_0$. If $\rho(B_1) < \rho(B_2)$, then for any $\epsilon > 0$, there is some integer $N(\epsilon)$, such that for all $k \geq N(\epsilon)$, we have

$$\sup_{\|u_0 - \tilde{u}\|=1} \left[\frac{\|v_k - \tilde{u}\|}{\|u_k - \tilde{u}\|} \right]^{1/k} \geq \frac{\rho(B_2)}{\rho(B_1) + \epsilon}.$$

Proof. Let $\|\cdot\|$ be the subordinate matrix norm. Recall that

$$u_k - \tilde{u} = B^k e_0,$$

with $e_0 = u_0 - \tilde{u}$. For every $k \in \mathbb{N}$, we have

$$(\rho(B_1))^k = \rho(B_1^k) \leq \|B_1^k\| = \sup_{\|e_0\|=1} \|B_1^k e_0\|,$$

which implies

$$\rho(B_1) = \sup_{\|e_0\|=1} \|B_1^k e_0\|^{1/k} = \|B_1\|^{1/k},$$

and statement (1) follows from Proposition 9.2.

Because $u_0 = v_0$, we have

$$\begin{aligned} u_k - \tilde{u} &= B_1^k e_0 \\ v_k - \tilde{u} &= B_2^k e_0, \end{aligned}$$

with $e_0 = u_0 - \tilde{u} = v_0 - \tilde{u}$. Again, by Proposition 9.2, for every $\epsilon > 0$, there is some natural number $N(\epsilon)$ such that if $k \geq N(\epsilon)$, then

$$\sup_{\|e_0\|=1} \|B_1^k e_0\|^{1/k} \leq \rho(B_1) + \epsilon.$$

Furthermore, for all $k \geq N(\epsilon)$, there exists a vector $e_0 = e_0(k)$ such that

$$\|e_0\| = 1 \quad \text{and} \quad \|B_2^k e_0\|^{1/k} = \|B_2^k\|^{1/k} \geq \rho(B_2),$$

which implies statement (2). □

In light of the above, we see that when we investigate new iterative methods, we have to deal with the following two problems:

1. Given an iterative method with matrix B , determine whether the method is convergent. This involves determining whether $\rho(B) < 1$, or equivalently whether there is a subordinate matrix norm such that $\|B\| < 1$. By Proposition 7.9, this implies that $I - B$ is invertible (since $\| -B \| = \|B\|$, Proposition 7.9 applies).
2. Given two convergent iterative methods, compare them. The iterative method which is faster is that whose matrix has the smaller spectral radius.

We now discuss three iterative methods for solving linear systems:

1. Jacobi's method
2. Gauss-Seidel's method
3. The relaxation method.

9.3 Description of the Methods of Jacobi, Gauss-Seidel, and Relaxation

The methods described in this section are instances of the following scheme: Given a linear system $Ax = b$, with A invertible, suppose we can write A in the form

$$A = M - N,$$

with M invertible, and “easy to invert,” which means that M is close to being a diagonal or a triangular matrix (perhaps by blocks). Then, $Au = b$ is equivalent to

$$Mu = Nu + b,$$

that is,

$$u = M^{-1}Nu + M^{-1}b.$$

Therefore, we are in the situation described in the previous sections with $B = M^{-1}N$ and $c = M^{-1}b$. In fact, since $A = M - N$, we have

$$B = M^{-1}N = M^{-1}(M - A) = I - M^{-1}A,$$

which shows that $I - B = M^{-1}A$ is invertible. The iterative method associated with the matrix $B = M^{-1}N$ is given by

$$u_{k+1} = M^{-1}Nu_k + M^{-1}b, \quad k \geq 0,$$

starting from any arbitrary vector u_0 . From a practical point of view, we do not invert M , and instead we solve iteratively the systems

$$Mu_{k+1} = Nu_k + b, \quad k \geq 0.$$

Various methods correspond to various ways of choosing M and N from A . The first two methods choose M and N as disjoint submatrices of A , but the relaxation method allows some overlapping of M and N .

To describe the various choices of M and N , it is convenient to write A in terms of three submatrices D, E, F , as

$$A = D - E - F,$$

where the only nonzero entries in D are the diagonal entries in A , the only nonzero entries in E are entries in A below the diagonal, and the only nonzero entries in F are entries in A above the diagonal. More explicitly, if

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1n-1} & a_{1n} \\ a_{21} & a_{22} & a_{23} & \cdots & a_{2n-1} & a_{2n} \\ a_{31} & a_{32} & a_{33} & \cdots & a_{3n-1} & a_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{n-11} & a_{n-12} & a_{n-13} & \cdots & a_{n-1n-1} & a_{n-1n} \\ a_{n1} & a_{n2} & a_{n3} & \cdots & a_{nn-1} & a_{nn} \end{pmatrix},$$

then

$$D = \begin{pmatrix} a_{11} & 0 & 0 & \cdots & 0 & 0 \\ 0 & a_{22} & 0 & \cdots & 0 & 0 \\ 0 & 0 & a_{33} & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & a_{n-1n-1} & 0 \\ 0 & 0 & 0 & \cdots & 0 & a_{nn} \end{pmatrix},$$

$$-E = \begin{pmatrix} 0 & 0 & 0 & \cdots & 0 & 0 \\ a_{21} & 0 & 0 & \cdots & 0 & 0 \\ a_{31} & a_{32} & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots & \vdots \\ a_{n-11} & a_{n-12} & a_{n-13} & \ddots & 0 & 0 \\ a_{n1} & a_{n2} & a_{n3} & \cdots & a_{nn-1} & 0 \end{pmatrix}, \quad -F = \begin{pmatrix} 0 & a_{12} & a_{13} & \cdots & a_{1n-1} & a_{1n} \\ 0 & 0 & a_{23} & \cdots & a_{2n-1} & a_{2n} \\ 0 & 0 & 0 & \ddots & a_{3n-1} & a_{3n} \\ \vdots & \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 0 & a_{n-1n} \\ 0 & 0 & 0 & \cdots & 0 & 0 \end{pmatrix}.$$

In *Jacobi's method*, we assume that all diagonal entries in A are nonzero, and we pick

$$M = D \\ N = E + F,$$

so that

$$B = M^{-1}N = D^{-1}(E + F) = I - D^{-1}A.$$

As a matter of notation, we let

$$J = I - D^{-1}A = D^{-1}(E + F),$$

which is called *Jacobi's matrix*. The corresponding method, *Jacobi's iterative method*, computes the sequence (u_k) using the recurrence

$$u_{k+1} = D^{-1}(E + F)u_k + D^{-1}b, \quad k \geq 0.$$

In practice, we iteratively solve the systems

$$Du_{k+1} = (E + F)u_k + b, \quad k \geq 0.$$

If we write $u_k = (u_1^k, \dots, u_n^k)$, we solve iteratively the following system:

$$\begin{array}{rcllclcl} a_{11}u_1^{k+1} & = & & -a_{12}u_2^k & -a_{13}u_3^k & \cdots & -a_{1n}u_n^k & + b_1 \\ a_{22}u_2^{k+1} & = & -a_{21}u_1^k & & -a_{23}u_3^k & \cdots & -a_{2n}u_n^k & + b_2 \\ \vdots & \vdots & \vdots & & & & & \\ a_{n-1n-1}u_{n-1}^{k+1} & = & -a_{n-11}u_1^k & \cdots & -a_{n-1n-2}u_{n-2}^k & & -a_{n-1n}u_n^k & + b_{n-1} \\ a_{nn}u_n^{k+1} & = & -a_{n1}u_1^k & -a_{n2}u_2^k & \cdots & -a_{nn-1}u_{n-1}^k & & + b_n \end{array}.$$

Observe that we can try to “speed up” the method by using the new value u_1^{k+1} instead of u_1^k in solving for u_2^{k+2} using the second equations, and more generally, use $u_1^{k+1}, \dots, u_{i-1}^{k+1}$ instead of u_1^k, \dots, u_{i-1}^k in solving for u_i^{k+1} in the i th equation. This observation leads to the system

$$\begin{array}{rcllclcl} a_{11}u_1^{k+1} & = & & -a_{12}u_2^k & -a_{13}u_3^k & \cdots & -a_{1n}u_n^k & + b_1 \\ a_{22}u_2^{k+1} & = & -a_{21}u_1^{k+1} & & -a_{23}u_3^k & \cdots & -a_{2n}u_n^k & + b_2 \\ \vdots & \vdots & \vdots & & & & & \\ a_{n-1n-1}u_{n-1}^{k+1} & = & -a_{n-11}u_1^{k+1} & \cdots & -a_{n-1n-2}u_{n-2}^{k+1} & & -a_{n-1n}u_n^k & + b_{n-1} \\ a_{nn}u_n^{k+1} & = & -a_{n1}u_1^{k+1} & -a_{n2}u_2^{k+1} & \cdots & -a_{nn-1}u_{n-1}^{k+1} & & + b_n, \end{array}$$

which, in matrix form, is written

$$Du_{k+1} = Eu_{k+1} + Fu_k + b.$$

Because D is invertible and E is lower triangular, the matrix $D - E$ is invertible, so the above equation is equivalent to

$$u_{k+1} = (D - E)^{-1}Fu_k + (D - E)^{-1}b, \quad k \geq 0.$$

The above corresponds to choosing M and N to be

$$\begin{aligned} M &= D - E \\ N &= F, \end{aligned}$$

and the matrix B is given by

$$B = M^{-1}N = (D - E)^{-1}F.$$

Since $M = D - E$ is invertible, we know that $I - B = M^{-1}A$ is also invertible.

The method that we just described is the *iterative method of Gauss-Seidel*, and the matrix B is called the *matrix of Gauss-Seidel* and denoted by \mathcal{L}_1 , with

$$\mathcal{L}_1 = (D - E)^{-1}F.$$

One of the advantages of the method of Gauss-Seidel is that it requires only half of the memory used by Jacobi's method, since we only need

$$u_1^{k+1}, \dots, u_{i-1}^{k+1}, u_{i+1}^k, \dots, u_n^k$$

to compute u_i^{k+1} . We also show that in certain important cases (for example, if A is a tridiagonal matrix), the method of Gauss-Seidel converges faster than Jacobi's method (in this case, they both converge or diverge simultaneously).

The new ingredient in the *relaxation method* is to incorporate part of the matrix D into N : we define M and N by

$$\begin{aligned} M &= \frac{D}{\omega} - E \\ N &= \frac{1 - \omega}{\omega}D + F, \end{aligned}$$

where $\omega \neq 0$ is a real parameter to be suitably chosen. Actually, we show in Section 9.4 that for the relaxation method to converge, we must have $\omega \in (0, 2)$. Note that the case $\omega = 1$ corresponds to the method of Gauss-Seidel.

If we assume that all diagonal entries of D are nonzero, the matrix M is invertible. The matrix B is denoted by \mathcal{L}_ω and called the *matrix of relaxation*, with

$$\mathcal{L}_\omega = \left(\frac{D}{\omega} - E \right)^{-1} \left(\frac{1-\omega}{\omega} D + F \right) = (D - \omega E)^{-1} ((1-\omega)D + \omega F).$$

The number ω is called the *parameter of relaxation*. When $\omega > 1$, the relaxation method is known as *successive overrelaxation*, abbreviated as *SOR*.

At first glance, the relaxation matrix \mathcal{L}_ω seems at lot more complicated than the Gauss-Seidel matrix \mathcal{L}_1 , but the iterative system associated with the relaxation method is very similar to the method of Gauss-Seidel, and is quite simple. Indeed, the system associated with the relaxation method is given by

$$\left(\frac{D}{\omega} - E \right) u_{k+1} = \left(\frac{1-\omega}{\omega} D + F \right) u_k + b,$$

which is equivalent to

$$(D - \omega E) u_{k+1} = ((1-\omega)D + \omega F) u_k + \omega b,$$

and can be written

$$Du_{k+1} = Du_k - \omega(Du_k - Eu_{k+1} - Fu_k - b).$$

Explicitly, this is the system

$$\begin{aligned} a_{11}u_1^{k+1} &= a_{11}u_1^k - \omega(a_{11}u_1^k + a_{12}u_2^k + a_{13}u_3^k + \cdots + a_{1n-2}u_{n-2}^k + a_{1n-1}u_{n-1}^k + a_{1n}u_n^k - b_1) \\ a_{22}u_2^{k+1} &= a_{22}u_2^k - \omega(a_{21}u_1^{k+1} + a_{22}u_2^k + a_{23}u_3^k + \cdots + a_{2n-2}u_{n-2}^k + a_{2n-1}u_{n-1}^k + a_{2n}u_n^k - b_2) \\ &\vdots \\ a_{nn}u_n^{k+1} &= a_{nn}u_n^k - \omega(a_{n1}u_1^{k+1} + a_{n2}u_2^{k+1} + \cdots + a_{nn-2}u_{n-2}^{k+1} + a_{nn-1}u_{n-1}^{k+1} + a_{nn}u_n^k - b_n). \end{aligned}$$

What remains to be done is to find conditions that ensure the convergence of the relaxation method (and the Gauss-Seidel method), that is:

1. Find conditions on ω , namely some interval $I \subseteq \mathbb{R}$ so that $\omega \in I$ implies $\rho(\mathcal{L}_\omega) < 1$; we will prove that $\omega \in (0, 2)$ is a necessary condition.
2. Find if there exist some *optimal value* ω_0 of $\omega \in I$, so that

$$\rho(\mathcal{L}_{\omega_0}) = \inf_{\omega \in I} \rho(\mathcal{L}_\omega).$$

We will give partial answers to the above questions in the next section.

It is also possible to extend the methods of this section by using *block decompositions* of the form $A = D - E - F$, where D, E , and F consist of blocks, and with D an invertible block-diagonal matrix.

9.4 Convergence of the Methods of Jacobi, Gauss-Seidel, and Relaxation

We begin with a general criterion for the convergence of an iterative method associated with a (complex) Hermitian, positive, definite matrix, $A = M - N$. Next, we apply this result to the relaxation method.

Proposition 9.5. *Let A be any Hermitian, positive, definite matrix, written as*

$$A = M - N,$$

with M invertible. Then, $M^ + N$ is Hermitian, and if it is positive, definite, then*

$$\rho(M^{-1}N) < 1,$$

so that the iterative method converges.

Proof. Since $M = A + N$ and A is Hermitian, $A^* = A$, so we get

$$M^* + N = A^* + N^* + N = A + N + N^* = M + N^* = (M^* + N)^*,$$

which shows that $M^* + N$ is indeed Hermitian.

Because A is symmetric, positive, definite, the function

$$v \mapsto (v^*Av)^{1/2}$$

from \mathbb{C}^n to \mathbb{R} is a vector norm $\|\cdot\|$, and let $\|\cdot\|$ also denote its subordinate matrix norm. We prove that

$$\|M^{-1}N\| < 1,$$

which, by Theorem 9.1 proves that $\rho(M^{-1}N) < 1$. By definition

$$\|M^{-1}N\| = \|I - M^{-1}A\| = \sup_{\|v\|=1} \|v - M^{-1}Av\|,$$

which leads us to evaluate $\|v - M^{-1}Av\|$ when $\|v\| = 1$. If we write $w = M^{-1}Av$, using the facts that $\|v\| = 1$, $v = A^{-1}Mw$, $A^* = A$, and $A = M - N$, we have

$$\begin{aligned} \|v - w\|^2 &= (v - w)^*A(v - w) \\ &= \|v\|^2 - v^*Aw - w^*Av + w^*Aw \\ &= 1 - w^*M^*w - w^*Mw + w^*Aw \\ &= 1 - w^*(M^* + N)w. \end{aligned}$$

Now, since we assumed that $M^* + N$ is positive definite, if $w \neq 0$, then $w^*(M^* + N)w > 0$, and we conclude that

$$\text{if } \|v\| = 1 \quad \text{then} \quad \|v - M^{-1}Av\| < 1.$$

Finally, the function

$$v \mapsto \|v - M^{-1}Av\|$$

is continuous as a composition of continuous functions, therefore it achieves its maximum on the compact subset $\{v \in \mathbb{C}^n \mid \|v\| = 1\}$, which proves that

$$\sup_{\|v\|=1} \|v - M^{-1}Av\| < 1,$$

and completes the proof. \square

Now, as in the previous sections, we assume that A is written as $A = D - E - F$, with D invertible, possibly in block form. The next theorem provides a sufficient condition (which turns out to be also necessary) for the relaxation method to converge (and thus, for the method of Gauss-Seidel to converge). This theorem is known as the *Ostrowski-Reich theorem*.

Theorem 9.6. *If $A = D - E - F$ is Hermitian, positive, definite, and if $0 < \omega < 2$, then the relaxation method converges. This also holds for a block decomposition of A .*

Proof. Recall that for the relaxation method, $A = M - N$ with

$$M = \frac{D}{\omega} - E$$

$$N = \frac{1-\omega}{\omega}D + F,$$

and because $D^* = D$, $E^* = F$ (since A is Hermitian) and $\omega \neq 0$ is real, we have

$$M^* + N = \frac{D^*}{\omega} - E^* + \frac{1-\omega}{\omega}D + F = \frac{2-\omega}{\omega}D.$$

If D consists of the diagonal entries of A , then we know from Section 6.3 that these entries are all positive, and since $\omega \in (0, 2)$, we see that the matrix $((2-\omega)/\omega)D$ is positive definite. If D consists of diagonal blocks of A , because A is positive, definite, by choosing vectors z obtained by picking a nonzero vector for each block of D and padding with zeros, we see that each block of D is positive, definite, and thus D itself is positive definite. Therefore, in all cases, $M^* + N$ is positive, definite, and we conclude by using Proposition 9.5. \square

Remark: What if we allow the parameter ω to be a nonzero complex number $\omega \in \mathbb{C}$? In this case, we get

$$M^* + N = \frac{D^*}{\bar{\omega}} - E^* + \frac{1-\omega}{\omega}D + F = \left(\frac{1}{\omega} + \frac{1}{\bar{\omega}} - 1\right)D.$$

But,

$$\frac{1}{\omega} + \frac{1}{\bar{\omega}} - 1 = \frac{\omega + \bar{\omega} - \omega\bar{\omega}}{\omega\bar{\omega}} = \frac{1 - (\omega - 1)(\bar{\omega} - 1)}{|\omega|^2} = \frac{1 - |\omega - 1|^2}{|\omega|^2},$$

so the relaxation method also converges for $\omega \in \mathbb{C}$, provided that

$$|\omega - 1| < 1.$$

This condition reduces to $0 < \omega < 2$ if ω is real.

Unfortunately, Theorem 9.6 does not apply to Jacobi's method, but in special cases, Proposition 9.5 can be used to prove its convergence. On the positive side, if a matrix is strictly column (or row) diagonally dominant, then it can be shown that the method of Jacobi and the method of Gauss-Seidel both converge. The relaxation method also converges if $\omega \in (0, 1]$, but this is not a very useful result because the speed-up of convergence usually occurs for $\omega > 1$.

We now prove that, without any assumption on $A = D - E - F$, other than the fact that A and D are invertible, in order for the relaxation method to converge, we must have $\omega \in (0, 2)$.

Proposition 9.7. *Given any matrix $A = D - E - F$, with A and D invertible, for any $\omega \neq 0$, we have*

$$\rho(\mathcal{L}_\omega) \geq |\omega - 1|.$$

Therefore, the relaxation method (possibly by blocks) does not converge unless $\omega \in (0, 2)$. If we allow ω to be complex, then we must have

$$|\omega - 1| < 1$$

for the relaxation method to converge.

Proof. Observe that the product $\lambda_1 \cdots \lambda_n$ of the eigenvalues of \mathcal{L}_ω , which is equal to $\det(\mathcal{L}_\omega)$, is given by

$$\lambda_1 \cdots \lambda_n = \det(\mathcal{L}_\omega) = \frac{\det\left(\frac{1-\omega}{\omega}D + F\right)}{\det\left(\frac{D}{\omega} - E\right)} = (1-\omega)^n.$$

It follows that

$$\rho(\mathcal{L}_\omega) \geq |\lambda_1 \cdots \lambda_n|^{1/n} = |1 - \omega|.$$

The proof is the same if $\omega \in \mathbb{C}$. □

We now consider the case where A is a *tridiagonal matrix*, possibly by blocks. In this case, we obtain precise results about the spectral radius of J and \mathcal{L}_ω , and as a consequence, about the convergence of these methods. We also obtain some information about the rate of convergence of these methods. We begin with the case $\omega = 1$, which is technically easier to deal with. The following proposition gives us the precise relationship between the spectral radii $\rho(J)$ and $\rho(\mathcal{L}_1)$ of the Jacobi matrix and the Gauss-Seidel matrix.

Proposition 9.8. *Let A be a tridiagonal matrix (possibly by blocks). If $\rho(J)$ is the spectral radius of the Jacobi matrix and $\rho(\mathcal{L}_1)$ is the spectral radius of the Gauss-Seidel matrix, then we have*

$$\rho(\mathcal{L}_1) = (\rho(J))^2.$$

Consequently, the method of Jacobi and the method of Gauss-Seidel both converge or both diverge simultaneously (even when A is tridiagonal by blocks); when they converge, the method of Gauss-Seidel converges faster than Jacobi's method.

Proof. We begin with a preliminary result. Let $A(\mu)$ with a tridiagonal matrix by block of the form

$$A(\mu) = \begin{pmatrix} A_1 & \mu^{-1}C_1 & 0 & 0 & \cdots & 0 \\ \mu B_1 & A_2 & \mu^{-1}C_2 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \ddots & \cdots & \vdots \\ \vdots & \cdots & \ddots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \mu B_{p-2} & A_{p-1} & \mu^{-1}C_{p-1} \\ 0 & \cdots & \cdots & 0 & \mu B_{p-1} & A_p \end{pmatrix},$$

then

$$\det(A(\mu)) = \det(A(1)), \quad \mu \neq 0.$$

To prove this fact, form the block diagonal matrix

$$P(\mu) = \text{diag}(\mu I_1, \mu^2 I_2, \dots, \mu^p I_p),$$

where I_j is the identity matrix of the same dimension as the block A_j . Then, it is easy to see that

$$A(\mu) = P(\mu)A(1)P(\mu)^{-1},$$

and thus,

$$\det(A(\mu)) = \det(P(\mu)A(1)P(\mu)^{-1}) = \det(A(1)).$$

Since the Jacobi matrix is $J = D^{-1}(E + F)$, the eigenvalues of J are the zeros of the characteristic polynomial

$$p_J(\lambda) = \det(\lambda I - D^{-1}(E + F)),$$

and thus, they are also the zeros of the polynomial

$$q_J(\lambda) = \det(\lambda D - E - F) = \det(D)p_J(\lambda).$$

Similarly, since the Gauss-Seidel matrix is $\mathcal{L}_1 = (D - E)^{-1}F$, the zeros of the characteristic polynomial

$$p_{\mathcal{L}_1}(\lambda) = \det(\lambda I - (D - E)^{-1}F)$$

are also the zeros of the polynomial

$$q_{\mathcal{L}_1}(\lambda) = \det(\lambda D - \lambda E - F) = \det(D - E)p_{\mathcal{L}_1}(\lambda).$$

Since A is tridiagonal (or tridiagonal by blocks), using our preliminary result with $\mu = \lambda \neq 0$, we get

$$q_{\mathcal{L}_1}(\lambda^2) = \det(\lambda^2 D - \lambda^2 E - F) = \det(\lambda^2 D - \lambda E - \lambda F) = \lambda^n q_J(\lambda).$$

By continuity, the above equation also holds for $\lambda = 0$. But then, we deduce that:

1. For any $\beta \neq 0$, if β is an eigenvalue of \mathcal{L}_1 , then $\beta^{1/2}$ and $-\beta^{1/2}$ are both eigenvalues of J , where $\beta^{1/2}$ is one of the complex square roots of β .
2. For any $\alpha \neq 0$, if α and $-\alpha$ are both eigenvalues of J , then α^2 is an eigenvalue of \mathcal{L}_1 .

The above immediately implies that $\rho(\mathcal{L}_1) = (\rho(J))^2$. \square

We now consider the more general situation where ω is any real in $(0, 2)$.

Proposition 9.9. *Let A be a tridiagonal matrix (possibly by blocks), and assume that the eigenvalues of the Jacobi matrix are all real. If $\omega \in (0, 2)$, then the method of Jacobi and the method of relaxation both converge or both diverge simultaneously (even when A is tridiagonal by blocks). When they converge, the function $\omega \mapsto \rho(\mathcal{L}_\omega)$ (for $\omega \in (0, 2)$) has a unique minimum equal to $\omega_0 - 1$ for*

$$\omega_0 = \frac{2}{1 + \sqrt{1 - (\rho(J))^2}},$$

where $1 < \omega_0 < 2$ if $\rho(J) > 0$. We also have $\rho(\mathcal{L}_1) = (\rho(J))^2$, as before.

Proof. The proof is very technical and can be found in Serre [96] and Ciarlet [24]. As in the proof of the previous proposition, we begin by showing that the eigenvalues of the matrix \mathcal{L}_ω are the zeros of the polynomial

$$q_{\mathcal{L}_\omega}(\lambda) = \det\left(\frac{\lambda + \omega - 1}{\omega} D - \lambda E - F\right) = \det\left(\frac{D}{\omega} - E\right) p_{\mathcal{L}_\omega}(\lambda),$$

where $p_{\mathcal{L}_\omega}(\lambda)$ is the characteristic polynomial of \mathcal{L}_ω . Then, using the preliminary fact from Proposition 9.8, it is easy to show that

$$q_{\mathcal{L}_\omega}(\lambda^2) = \lambda^n q_J\left(\frac{\lambda^2 + \omega - 1}{\lambda\omega}\right),$$

for all $\lambda \in \mathbb{C}$, with $\lambda \neq 0$. This time, we cannot extend the above equation to $\lambda = 0$. This leads us to consider the equation

$$\frac{\lambda^2 + \omega - 1}{\lambda\omega} = \alpha,$$

which is equivalent to

$$\lambda^2 - \alpha\omega\lambda + \omega - 1 = 0,$$

for all $\lambda \neq 0$. Since $\lambda \neq 0$, the above equivalence does not hold for $\omega = 1$, but this is not a problem since the case $\omega = 1$ has already been considered in the previous proposition. Then, we can show the following:

1. For any $\beta \neq 0$, if β is an eigenvalue of \mathcal{L}_ω , then

$$\frac{\beta + \omega - 1}{\beta^{1/2}\omega}, \quad -\frac{\beta + \omega - 1}{\beta^{1/2}\omega}$$

are eigenvalues of J .

2. For every $\alpha \neq 0$, if α and $-\alpha$ are eigenvalues of J , then $\mu_+(\alpha, \omega)$ and $\mu_-(\alpha, \omega)$ are eigenvalues of \mathcal{L}_ω , where $\mu_+(\alpha, \omega)$ and $\mu_-(\alpha, \omega)$ are the squares of the roots of the equation

$$\lambda^2 - \alpha\omega\lambda + \omega - 1 = 0.$$

It follows that

$$\rho(\mathcal{L}_\omega) = \max_{\lambda \mid p_J(\lambda)=0} \{\max(|\mu_+(\alpha, \omega)|, |\mu_-(\alpha, \omega)|)\},$$

and since we are assuming that J has real roots, we are led to study the function

$$M(\alpha, \omega) = \max\{|\mu_+(\alpha, \omega)|, |\mu_-(\alpha, \omega)|\},$$

where $\alpha \in \mathbb{R}$ and $\omega \in (0, 2)$. Actually, because $M(-\alpha, \omega) = M(\alpha, \omega)$, it is only necessary to consider the case where $\alpha \geq 0$.

Note that for $\alpha \neq 0$, the roots of the equation

$$\lambda^2 - \alpha\omega\lambda + \omega - 1 = 0.$$

are

$$\frac{\alpha\omega \pm \sqrt{\alpha^2\omega^2 - 4\omega + 4}}{2}.$$

In turn, this leads to consider the roots of the equation

$$\omega^2\alpha^2 - 4\omega + 4 = 0,$$

which are

$$\frac{2(1 \pm \sqrt{1 - \alpha^2})}{\alpha^2},$$

for $\alpha \neq 0$. Since we have

$$\frac{2(1 + \sqrt{1 - \alpha^2})}{\alpha^2} = \frac{2(1 + \sqrt{1 - \alpha^2})(1 - \sqrt{1 - \alpha^2})}{\alpha^2(1 - \sqrt{1 - \alpha^2})} = \frac{2}{1 - \sqrt{1 - \alpha^2}}$$

and

$$\frac{2(1 - \sqrt{1 - \alpha^2})}{\alpha^2} = \frac{2(1 + \sqrt{1 - \alpha^2})(1 - \sqrt{1 - \alpha^2})}{\alpha^2(1 + \sqrt{1 - \alpha^2})} = \frac{2}{1 + \sqrt{1 - \alpha^2}},$$

these roots are

$$\omega_0(\alpha) = \frac{2}{1 + \sqrt{1 - \alpha^2}}, \quad \omega_1(\alpha) = \frac{2}{1 - \sqrt{1 - \alpha^2}}.$$

Observe that the expression for $\omega_0(\alpha)$ is exactly the expression in the statement of our proposition! The rest of the proof consists in analyzing the variations of the function $M(\alpha, \omega)$ by considering various cases for α . In the end, we find that the minimum of $\rho(\mathcal{L}_\omega)$ is obtained for $\omega_0(\rho(J))$. The details are tedious and we omit them. The reader will find complete proofs in Serre [96] and Ciarlet [24]. \square

Combining the results of Theorem 9.6 and Proposition 9.9, we obtain the following result which gives precise information about the spectral radii of the matrices J , \mathcal{L}_1 , and \mathcal{L}_ω .

Proposition 9.10. *Let A be a tridiagonal matrix (possibly by blocks) which is Hermitian, positive, definite. Then, the methods of Jacobi, Gauss-Seidel, and relaxation, all converge for $\omega \in (0, 2)$. There is a unique optimal relaxation parameter*

$$\omega_0 = \frac{2}{1 + \sqrt{1 - (\rho(J))^2}},$$

such that

$$\rho(\mathcal{L}_{\omega_0}) = \inf_{0 < \omega < 2} \rho(\mathcal{L}_\omega) = \omega_0 - 1.$$

Furthermore, if $\rho(J) > 0$, then

$$\rho(\mathcal{L}_{\omega_0}) < \rho(\mathcal{L}_1) = (\rho(J))^2 < \rho(J),$$

and if $\rho(J) = 0$, then $\omega_0 = 1$ and $\rho(\mathcal{L}_1) = \rho(J) = 0$.

Proof. In order to apply Proposition 9.9, we have to check that $J = D^{-1}(E + F)$ has real eigenvalues. However, if α is any eigenvalue of J and if u is any corresponding eigenvector, then

$$D^{-1}(E + F)u = \alpha u$$

implies that

$$(E + F)u = \alpha Du,$$

and since $A = D - E - F$, the above shows that $(D - A)u = \alpha Du$, that is,

$$Au = (1 - \alpha)Du.$$

Consequently,

$$u^* Au = (1 - \alpha)u^* Du,$$

and since A and D are hermitian, positive, definite, we have $u^* Au > 0$ and $u^* Du > 0$ if $u \neq 0$, which proves that $\alpha \in \mathbb{R}$. The rest follows from Theorem 9.6 and Proposition 9.9. \square

Remark: It is preferable to overestimate rather than underestimate the relaxation parameter when the optimum relaxation parameter is not known exactly.

9.5 Summary

The main concepts and results of this chapter are listed below:

- Iterative methods. Splitting A as $A = M - N$.
- *Convergence of a sequence of vectors or matrices.*
- A criterion for the convergence of the sequence (B^k) of powers of a matrix B to zero in terms of the spectral radius $\rho(B)$.
- A characterization of the spectral radius $\rho(B)$ as the limit of the sequence $(\|B^k\|^{1/k})$.
- A criterion of the convergence of iterative methods.
- Asymptotic behavior of iterative methods.
- Splitting A as $A = D - E - F$, and the methods of *Jacobi*, *Gauss-Seidel*, and *relaxation* (and *SOR*).
- The *Jacobi matrix*, $J = D^{-1}(E + F)$.
- The *Gauss-Seidel matrix*, $\mathcal{L}_2 = (D - E)^{-1}F$.
- The *matrix of relaxation*, $\mathcal{L}_\omega = (D - \omega E)^{-1}((1 - \omega)D + \omega F)$.
- Convergence of iterative methods: a general result when $A = M - N$ is Hermitian, positive, definite.
- A sufficient condition for the convergence of the methods of Jacobi, Gauss-Seidel, and relaxation. The *Ostrowski-Reich Theorem*: A is symmetric, positive, definite, and $\omega \in (0, 2)$.
- A necessary condition for the convergence of the methods of Jacobi, Gauss-Seidel, and relaxation: $\omega \in (0, 2)$.
- The case of tridiagonal matrices (possibly by blocks). Simultaneous convergence or divergence of Jacobi's method and Gauss-Seidel's method, and comparison of the spectral radii of $\rho(J)$ and $\rho(\mathcal{L}_1)$: $\rho(\mathcal{L}_1) = (\rho(J))^2$.
- The case of tridiagonal, Hermitian, positive, definite matrices (possibly by blocks). The methods of Jacobi, Gauss-Seidel, and relaxation, all converge.
- In the above case, there is a unique optimal relaxation parameter for which $\rho(\mathcal{L}_{\omega_0}) < \rho(\mathcal{L}_1) = (\rho(J))^2 < \rho(J)$ (if $\rho(J) \neq 0$).

Chapter 10

Euclidean Spaces

Rien n'est beau que le vrai.
—Hermann Minkowski

10.1 Inner Products, Euclidean Spaces

So far, the framework of vector spaces allows us to deal with ratios of vectors and linear combinations, but there is no way to express the notion of length of a line segment or to talk about orthogonality of vectors. A Euclidean structure allows us to deal with *metric notions* such as orthogonality and length (or distance).

This chapter covers the bare bones of Euclidean geometry. Deeper aspects of Euclidean geometry are investigated in Chapter 11. One of our main goals is to give the basic properties of the transformations that preserve the Euclidean structure, rotations and reflections, since they play an important role in practice. Euclidean geometry is the study of properties invariant under certain affine maps called *rigid motions*. Rigid motions are the maps that preserve the distance between points.

We begin by defining inner products and Euclidean spaces. The Cauchy–Schwarz inequality and the Minkowski inequality are shown. We define orthogonality of vectors and of subspaces, orthogonal bases, and orthonormal bases. We prove that every finite-dimensional Euclidean space has orthonormal bases. The first proof uses duality, and the second one the Gram–Schmidt orthogonalization procedure. The QR -decomposition for invertible matrices is shown as an application of the Gram–Schmidt procedure. Linear isometries (also called orthogonal transformations) are defined and studied briefly. We conclude with a short section in which some applications of Euclidean geometry are sketched. One of the most important applications, the method of least squares, is discussed in Chapter 17.

For a more detailed treatment of Euclidean geometry, see Berger [8, 9], Snapper and Troyer [99], or any other book on geometry, such as Pedoe [88], Coxeter [26], Fresnel [40], Tisseron [109], or Cagnac, Ramis, and Commeau [19]. Serious readers should consult Emil

Artin's famous book [3], which contains an in-depth study of the orthogonal group, as well as other groups arising in geometry. It is still worth consulting some of the older classics, such as Hadamard [53, 54] and Rouché and de Comberousse [89]. The first edition of [53] was published in 1898, and finally reached its thirteenth edition in 1947! In this chapter it is assumed that all vector spaces are defined over the field \mathbb{R} of real numbers unless specified otherwise (in a few cases, over the complex numbers \mathbb{C}).

First, we define a Euclidean structure on a vector space. Technically, a Euclidean structure over a vector space E is provided by a symmetric bilinear form on the vector space satisfying some extra properties. Recall that a bilinear form $\varphi: E \times E \rightarrow \mathbb{R}$ is *definite* if for every $u \in E$, $u \neq 0$ implies that $\varphi(u, u) \neq 0$, and *positive* if for every $u \in E$, $\varphi(u, u) \geq 0$.

Definition 10.1. A *Euclidean space* is a real vector space E equipped with a symmetric bilinear form $\varphi: E \times E \rightarrow \mathbb{R}$ that is *positive definite*. More explicitly, $\varphi: E \times E \rightarrow \mathbb{R}$ satisfies the following axioms:

$$\begin{aligned}\varphi(u_1 + u_2, v) &= \varphi(u_1, v) + \varphi(u_2, v), \\ \varphi(u, v_1 + v_2) &= \varphi(u, v_1) + \varphi(u, v_2), \\ \varphi(\lambda u, v) &= \lambda \varphi(u, v), \\ \varphi(u, \lambda v) &= \lambda \varphi(u, v), \\ \varphi(u, v) &= \varphi(v, u), \\ u \neq 0 &\text{ implies that } \varphi(u, u) > 0.\end{aligned}$$

The real number $\varphi(u, v)$ is also called the *inner product (or scalar product) of u and v* . We also define the *quadratic form associated with φ* as the function $\Phi: E \rightarrow \mathbb{R}_+$ such that

$$\Phi(u) = \varphi(u, u),$$

for all $u \in E$.

Since φ is bilinear, we have $\varphi(0, 0) = 0$, and since it is positive definite, we have the stronger fact that

$$\varphi(u, u) = 0 \quad \text{iff} \quad u = 0,$$

that is, $\Phi(u) = 0$ iff $u = 0$.

Given an inner product $\varphi: E \times E \rightarrow \mathbb{R}$ on a vector space E , we also denote $\varphi(u, v)$ by

$$u \cdot v \quad \text{or} \quad \langle u, v \rangle \quad \text{or} \quad (u|v),$$

and $\sqrt{\Phi(u)}$ by $\|u\|$.

Example 10.1. The standard example of a Euclidean space is \mathbb{R}^n , under the inner product \cdot defined such that

$$(x_1, \dots, x_n) \cdot (y_1, \dots, y_n) = x_1 y_1 + x_2 y_2 + \dots + x_n y_n.$$

This Euclidean space is denoted by \mathbb{E}^n .

There are other examples.

Example 10.2. For instance, let E be a vector space of dimension 2, and let (e_1, e_2) be a basis of E . If $a > 0$ and $b^2 - ac < 0$, the bilinear form defined such that

$$\varphi(x_1e_1 + y_1e_2, x_2e_1 + y_2e_2) = ax_1x_2 + b(x_1y_2 + x_2y_1) + cy_1y_2$$

yields a Euclidean structure on E . In this case,

$$\Phi(xe_1 + ye_2) = ax^2 + 2bxy + cy^2.$$

Example 10.3. Let $\mathcal{C}[a, b]$ denote the set of continuous functions $f: [a, b] \rightarrow \mathbb{R}$. It is easily checked that $\mathcal{C}[a, b]$ is a vector space of infinite dimension. Given any two functions $f, g \in \mathcal{C}[a, b]$, let

$$\langle f, g \rangle = \int_a^b f(t)g(t)dt.$$

We leave as an easy exercise that $\langle -, - \rangle$ is indeed an inner product on $\mathcal{C}[a, b]$. In the case where $a = -\pi$ and $b = \pi$ (or $a = 0$ and $b = 2\pi$, this makes basically no difference), one should compute

$$\langle \sin px, \sin qx \rangle, \quad \langle \sin px, \cos qx \rangle, \quad \text{and} \quad \langle \cos px, \cos qx \rangle,$$

for all natural numbers $p, q \geq 1$. The outcome of these calculations is what makes Fourier analysis possible!

Example 10.4. Let $E = M_n(\mathbb{R})$ be the vector space of real $n \times n$ matrices. If we view a matrix $A \in M_n(\mathbb{R})$ as a “long” column vector obtained by concatenating together its columns, we can define the inner product of two matrices $A, B \in M_n(\mathbb{R})$ as

$$\langle A, B \rangle = \sum_{i,j=1}^n a_{ij}b_{ij},$$

which can be conveniently written as

$$\langle A, B \rangle = \text{tr}(A^\top B) = \text{tr}(B^\top A).$$

Since this can be viewed as the Euclidean product on \mathbb{R}^{n^2} , it is an inner product on $M_n(\mathbb{R})$. The corresponding norm

$$\|A\|_F = \sqrt{\text{tr}(A^\top A)}$$

is the Frobenius norm (see Section 7.2).

Let us observe that φ can be recovered from Φ . Indeed, by bilinearity and symmetry, we have

$$\begin{aligned} \Phi(u + v) &= \varphi(u + v, u + v) \\ &= \varphi(u, u + v) + \varphi(v, u + v) \\ &= \varphi(u, u) + 2\varphi(u, v) + \varphi(v, v) \\ &= \Phi(u) + 2\varphi(u, v) + \Phi(v). \end{aligned}$$

Thus, we have

$$\varphi(u, v) = \frac{1}{2}[\Phi(u + v) - \Phi(u) - \Phi(v)].$$

We also say that φ is the *polar form* of Φ .

If E is finite-dimensional and if $\varphi: E \times E \rightarrow \mathbb{R}$ is a bilinear form on E , given any basis (e_1, \dots, e_n) of E , we can write $x = \sum_{i=1}^n x_i e_i$ and $y = \sum_{j=1}^n y_j e_j$, and we have

$$\varphi(x, y) = \varphi\left(\sum_{i=1}^n x_i e_i, \sum_{j=1}^n y_j e_j\right) = \sum_{i,j=1}^n x_i y_j \varphi(e_i, e_j).$$

If we let G be the matrix $G = (\varphi(e_i, e_j))$, and if x and y are the column vectors associated with (x_1, \dots, x_n) and (y_1, \dots, y_n) , then we can write

$$\varphi(x, y) = x^\top G y = y^\top G^\top x.$$

Note that we are committing an abuse of notation, since $x = \sum_{i=1}^n x_i e_i$ is a vector in E , but the column vector associated with (x_1, \dots, x_n) belongs to \mathbb{R}^n . To avoid this minor abuse, we could denote the column vector associated with (x_1, \dots, x_n) by \mathbf{x} (and similarly \mathbf{y} for the column vector associated with (y_1, \dots, y_n)), in which case the “correct” expression for $\varphi(x, y)$ is

$$\varphi(x, y) = \mathbf{x}^\top G \mathbf{y}.$$

However, in view of the isomorphism between E and \mathbb{R}^n , to keep notation as simple as possible, we will use x and y instead of \mathbf{x} and \mathbf{y} .

Also observe that φ is symmetric iff $G = G^\top$, and φ is positive definite iff the matrix G is positive definite, that is,

$$x^\top G x > 0 \quad \text{for all } x \in \mathbb{R}^n, x \neq 0.$$

The matrix G associated with an inner product is called the *Gram matrix* of the inner product with respect to the basis (e_1, \dots, e_n) .

Conversely, if A is a symmetric positive definite $n \times n$ matrix, it is easy to check that the bilinear form

$$\langle x, y \rangle = x^\top A y$$

is an inner product. If we make a change of basis from the basis (e_1, \dots, e_n) to the basis (f_1, \dots, f_n) , and if the change of basis matrix is P (where the j th column of P consists of the coordinates of f_j over the basis (e_1, \dots, e_n)), then with respect to coordinates x' and y' over the basis (f_1, \dots, f_n) , we have

$$\langle x, y \rangle = x^\top G y = x'^\top P^\top G P y',$$

so the matrix of our inner product over the basis (f_1, \dots, f_n) is $P^\top G P$. We summarize these facts in the following proposition.

Proposition 10.1. *Let E be a finite-dimensional vector space, and let (e_1, \dots, e_n) be a basis of E .*

1. *For any inner product $\langle -, - \rangle$ on E , if $G = (\langle e_i, e_j \rangle)$ is the Gram matrix of the inner product $\langle -, - \rangle$ w.r.t. the basis (e_1, \dots, e_n) , then G is symmetric positive definite.*
2. *For any change of basis matrix P , the Gram matrix of $\langle -, - \rangle$ with respect to the new basis is $P^\top GP$.*
3. *If A is any $n \times n$ symmetric positive definite matrix, then*

$$\langle x, y \rangle = x^\top Ay$$

is an inner product on E .

We will see later that a symmetric matrix is positive definite iff its eigenvalues are all positive.

One of the very important properties of an inner product φ is that the map $u \mapsto \sqrt{\Phi(u)}$ is a norm.

Proposition 10.2. *Let E be a Euclidean space with inner product φ , and let Φ be the corresponding quadratic form. For all $u, v \in E$, we have the Cauchy–Schwarz inequality*

$$\varphi(u, v)^2 \leq \Phi(u)\Phi(v),$$

the equality holding iff u and v are linearly dependent.

We also have the Minkowski inequality

$$\sqrt{\Phi(u+v)} \leq \sqrt{\Phi(u)} + \sqrt{\Phi(v)},$$

the equality holding iff u and v are linearly dependent, where in addition if $u \neq 0$ and $v \neq 0$, then $u = \lambda v$ for some $\lambda > 0$.

Proof. For any vectors $u, v \in E$, we define the function $T: \mathbb{R} \rightarrow \mathbb{R}$ such that

$$T(\lambda) = \Phi(u + \lambda v),$$

for all $\lambda \in \mathbb{R}$. Using bilinearity and symmetry, we have

$$\begin{aligned} \Phi(u + \lambda v) &= \varphi(u + \lambda v, u + \lambda v) \\ &= \varphi(u, u + \lambda v) + \lambda \varphi(v, u + \lambda v) \\ &= \varphi(u, u) + 2\lambda \varphi(u, v) + \lambda^2 \varphi(v, v) \\ &= \Phi(u) + 2\lambda \varphi(u, v) + \lambda^2 \Phi(v). \end{aligned}$$

Since φ is positive definite, Φ is nonnegative, and thus $T(\lambda) \geq 0$ for all $\lambda \in \mathbb{R}$. If $\Phi(v) = 0$, then $v = 0$, and we also have $\varphi(u, v) = 0$. In this case, the Cauchy–Schwarz inequality is trivial, and $v = 0$ and u are linearly dependent.

Now, assume $\Phi(v) > 0$. Since $T(\lambda) \geq 0$, the quadratic equation

$$\lambda^2\Phi(v) + 2\lambda\varphi(u, v) + \Phi(u) = 0$$

cannot have distinct real roots, which means that its discriminant

$$\Delta = 4(\varphi(u, v)^2 - \Phi(u)\Phi(v))$$

is null or negative, which is precisely the Cauchy–Schwarz inequality

$$\varphi(u, v)^2 \leq \Phi(u)\Phi(v).$$

If

$$\varphi(u, v)^2 = \Phi(u)\Phi(v)$$

then there are two cases. If $\Phi(v) = 0$, then $v = 0$ and u and v are linearly dependent. If $\Phi(v) \neq 0$, then the above quadratic equation has a double root λ_0 , and we have $\Phi(u + \lambda_0 v) = 0$. Since φ is positive definite, $\Phi(u + \lambda_0 v) = 0$ implies that $u + \lambda_0 v = 0$, which shows that u and v are linearly dependent. Conversely, it is easy to check that we have equality when u and v are linearly dependent.

The Minkowski inequality

$$\sqrt{\Phi(u + v)} \leq \sqrt{\Phi(u)} + \sqrt{\Phi(v)}$$

is equivalent to

$$\Phi(u + v) \leq \Phi(u) + \Phi(v) + 2\sqrt{\Phi(u)\Phi(v)}.$$

However, we have shown that

$$2\varphi(u, v) = \Phi(u + v) - \Phi(u) - \Phi(v),$$

and so the above inequality is equivalent to

$$\varphi(u, v) \leq \sqrt{\Phi(u)\Phi(v)},$$

which is trivial when $\varphi(u, v) \leq 0$, and follows from the Cauchy–Schwarz inequality when $\varphi(u, v) \geq 0$. Thus, the Minkowski inequality holds. Finally, assume that $u \neq 0$ and $v \neq 0$, and that

$$\sqrt{\Phi(u + v)} = \sqrt{\Phi(u)} + \sqrt{\Phi(v)}.$$

When this is the case, we have

$$\varphi(u, v) = \sqrt{\Phi(u)\Phi(v)},$$

and we know from the discussion of the Cauchy–Schwarz inequality that the equality holds iff u and v are linearly dependent. The Minkowski inequality is an equality when u or v is null. Otherwise, if $u \neq 0$ and $v \neq 0$, then $u = \lambda v$ for some $\lambda \neq 0$, and since

$$\varphi(u, v) = \lambda\varphi(v, v) = \sqrt{\Phi(u)\Phi(v)},$$

by positivity, we must have $\lambda > 0$. □

Note that the Cauchy–Schwarz inequality can also be written as

$$|\varphi(u, v)| \leq \sqrt{\Phi(u)}\sqrt{\Phi(v)}.$$

Remark: It is easy to prove that the Cauchy–Schwarz and the Minkowski inequalities still hold for a symmetric bilinear form that is positive, but not necessarily definite (i.e., $\varphi(u, v) \geq 0$ for all $u, v \in E$). However, u and v need not be linearly dependent when the equality holds.

The Minkowski inequality

$$\sqrt{\Phi(u+v)} \leq \sqrt{\Phi(u)} + \sqrt{\Phi(v)}$$

shows that the map $u \mapsto \sqrt{\Phi(u)}$ satisfies the convexity inequality (also known as triangle inequality), condition (N3) of Definition 7.1, and since φ is bilinear and positive definite, it also satisfies conditions (N1) and (N2) of Definition 7.1, and thus it is a *norm* on E . The norm induced by φ is called the *Euclidean norm induced by φ* .

Note that the Cauchy–Schwarz inequality can be written as

$$|u \cdot v| \leq \|u\|\|v\|,$$

and the Minkowski inequality as

$$\|u+v\| \leq \|u\| + \|v\|.$$

Remark: One might wonder if every norm on a vector space is induced by some Euclidean inner product. In general, this is false, but remarkably, there is a simple necessary and sufficient condition, which is that the norm must satisfy the *parallelogram law*:

$$\|u+v\|^2 + \|u-v\|^2 = 2(\|u\|^2 + \|v\|^2).$$

If $\langle -, - \rangle$ is an inner product, then we have

$$\begin{aligned}\|u+v\|^2 &= \|u\|^2 + \|v\|^2 + 2\langle u, v \rangle \\ \|u-v\|^2 &= \|u\|^2 + \|v\|^2 - 2\langle u, v \rangle,\end{aligned}$$

and by adding and subtracting these identities, we get the parallelogram law and the equation

$$\langle u, v \rangle = \frac{1}{4}(\|u+v\|^2 - \|u-v\|^2),$$

which allows us to recover $\langle -, - \rangle$ from the norm.

Conversely, if $\| \cdot \|$ is a norm satisfying the parallelogram law, and if it comes from an inner product, then this inner product must be given by

$$\langle u, v \rangle = \frac{1}{4}(\|u + v\|^2 - \|u - v\|^2).$$

We need to prove that the above form is indeed symmetric and bilinear.

Symmetry holds because $\|u - v\| = \|(u - v)\| = \|v - u\|$. Let us prove additivity in the variable u . By the parallelogram law, we have

$$2(\|x + z\|^2 + \|y\|^2) = \|x + y + z\|^2 + \|x - y + z\|^2$$

which yields

$$\begin{aligned} \|x + y + z\|^2 &= 2(\|x + z\|^2 + \|y\|^2) - \|x - y + z\|^2 \\ \|x + y + z\|^2 &= 2(\|y + z\|^2 + \|x\|^2) - \|y - x + z\|^2, \end{aligned}$$

where the second formula is obtained by swapping x and y . Then by adding up these equations, we get

$$\|x + y + z\|^2 = \|x\|^2 + \|y\|^2 + \|x + z\|^2 + \|y + z\|^2 - \frac{1}{2}\|x - y + z\|^2 - \frac{1}{2}\|y - x + z\|^2.$$

Replacing z by $-z$ in the above equation, we get

$$\|x + y - z\|^2 = \|x\|^2 + \|y\|^2 + \|x - z\|^2 + \|y - z\|^2 - \frac{1}{2}\|x - y - z\|^2 - \frac{1}{2}\|y - x - z\|^2,$$

Since $\|x - y + z\| = \|(x - y + z)\| = \|y - x - z\|$ and $\|y - x + z\| = \|(y - x + z)\| = \|x - y - z\|$, by subtracting the last two equations, we get

$$\begin{aligned} \langle x + y, z \rangle &= \frac{1}{4}(\|x + y + z\|^2 - \|x + y - z\|^2) \\ &= \frac{1}{4}(\|x + z\|^2 - \|x - z\|^2) + \frac{1}{4}(\|y + z\|^2 - \|y - z\|^2) \\ &= \langle x, z \rangle + \langle y, z \rangle, \end{aligned}$$

as desired.

Proving that

$$\langle \lambda x, y \rangle = \lambda \langle x, y \rangle \quad \text{for all } \lambda \in \mathbb{R}$$

is a little tricky. The strategy is to prove the identity for $\lambda \in \mathbb{Z}$, then to promote it to \mathbb{Q} , and then to \mathbb{R} by continuity.

Since

$$\begin{aligned} \langle -u, v \rangle &= \frac{1}{4}(\| -u + v \|^2 - \| -u - v \|^2) \\ &= \frac{1}{4}(\| u - v \|^2 - \| u + v \|^2) \\ &= -\langle u, v \rangle, \end{aligned}$$

the property holds for $\lambda = -1$. By linearity and by induction, for any $n \in \mathbb{N}$ with $n \geq 1$, writing $n = n - 1 + 1$, we get

$$\langle \lambda x, y \rangle = \lambda \langle x, y \rangle \quad \text{for all } \lambda \in \mathbb{N},$$

and since the above also holds for $\lambda = -1$, it holds for all $\lambda \in \mathbb{Z}$. For $\lambda = p/q$ with $p, q \in \mathbb{Z}$ and $q \neq 0$, we have

$$q \langle (p/q)u, v \rangle = \langle pu, v \rangle = p \langle u, v \rangle,$$

which shows that

$$\langle (p/q)u, v \rangle = (p/q) \langle u, v \rangle,$$

and thus

$$\langle \lambda x, y \rangle = \lambda \langle x, y \rangle \quad \text{for all } \lambda \in \mathbb{Q}.$$

To finish the proof, we use the fact that a norm is a continuous map $x \mapsto \|x\|$. Then, the continuous function $t \mapsto \frac{1}{t} \langle tu, v \rangle$ defined on $\mathbb{R} - \{0\}$ agrees with $\langle u, v \rangle$ on $\mathbb{Q} - \{0\}$, so it is equal to $\langle u, v \rangle$ on $\mathbb{R} - \{0\}$. The case $\lambda = 0$ is trivial, so we are done.

We now define orthogonality.

10.2 Orthogonality, Duality, Adjoint of a Linear Map

An inner product on a vector space gives the ability to define the notion of orthogonality. Families of nonnull pairwise orthogonal vectors must be linearly independent. They are called orthogonal families. In a vector space of finite dimension it is always possible to find orthogonal bases. This is very useful theoretically and practically. Indeed, in an orthogonal basis, finding the coordinates of a vector is very cheap: It takes an inner product. Fourier series make crucial use of this fact. When E has finite dimension, we prove that the inner product on E induces a natural isomorphism between E and its dual space E^* . This allows us to define the adjoint of a linear map in an intrinsic fashion (i.e., independently of bases). It is also possible to orthonormalize any basis (certainly when the dimension is finite). We give two proofs, one using duality, the other more constructive using the Gram–Schmidt orthonormalization procedure.

Definition 10.2. Given a Euclidean space E , any two vectors $u, v \in E$ are *orthogonal*, or *perpendicular*, if $u \cdot v = 0$. Given a family $(u_i)_{i \in I}$ of vectors in E , we say that $(u_i)_{i \in I}$ is *orthogonal* if $u_i \cdot u_j = 0$ for all $i, j \in I$, where $i \neq j$. We say that the family $(u_i)_{i \in I}$ is *orthonormal* if $u_i \cdot u_j = 0$ for all $i, j \in I$, where $i \neq j$, and $\|u_i\| = u_i \cdot u_i = 1$, for all $i \in I$. For any subset F of E , the set

$$F^\perp = \{v \in E \mid u \cdot v = 0, \text{ for all } u \in F\},$$

of all vectors orthogonal to all vectors in F , is called the *orthogonal complement* of F .

Since inner products are positive definite, observe that for any vector $u \in E$, we have

$$u \cdot v = 0 \quad \text{for all } v \in E \quad \text{iff} \quad u = 0.$$

It is immediately verified that the orthogonal complement F^\perp of F is a subspace of E .

Example 10.5. Going back to Example 10.3 and to the inner product

$$\langle f, g \rangle = \int_{-\pi}^{\pi} f(t)g(t)dt$$

on the vector space $\mathcal{C}[-\pi, \pi]$, it is easily checked that

$$\langle \sin px, \sin qx \rangle = \begin{cases} \pi & \text{if } p = q, p, q \geq 1, \\ 0 & \text{if } p \neq q, p, q \geq 1, \end{cases}$$

$$\langle \cos px, \cos qx \rangle = \begin{cases} \pi & \text{if } p = q, p, q \geq 1, \\ 0 & \text{if } p \neq q, p, q \geq 0, \end{cases}$$

and

$$\langle \sin px, \cos qx \rangle = 0,$$

for all $p \geq 1$ and $q \geq 0$, and of course, $\langle 1, 1 \rangle = \int_{-\pi}^{\pi} dx = 2\pi$.

As a consequence, the family $(\sin px)_{p \geq 1} \cup (\cos qx)_{q \geq 0}$ is orthogonal. It is not orthonormal, but becomes so if we divide every trigonometric function by $\sqrt{\pi}$, and 1 by $\sqrt{2\pi}$.

We leave the following simple two results as exercises.

Proposition 10.3. *Given a Euclidean space E , for any family $(u_i)_{i \in I}$ of nonnull vectors in E , if $(u_i)_{i \in I}$ is orthogonal, then it is linearly independent.*

Proposition 10.4. *Given a Euclidean space E , any two vectors $u, v \in E$ are orthogonal iff*

$$\|u + v\|^2 = \|u\|^2 + \|v\|^2.$$

One of the most useful features of orthonormal bases is that they afford a very simple method for computing the coordinates of a vector over any basis vector. Indeed, assume that (e_1, \dots, e_m) is an orthonormal basis. For any vector

$$x = x_1 e_1 + \dots + x_m e_m,$$

if we compute the inner product $x \cdot e_i$, we get

$$x \cdot e_i = x_1 e_1 \cdot e_i + \dots + x_i e_i \cdot e_i + \dots + x_m e_m \cdot e_i = x_i,$$

since

$$e_i \cdot e_j = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j \end{cases}$$

is the property characterizing an orthonormal family. Thus,

$$x_i = x \cdot e_i,$$

which means that $x_i e_i = (x \cdot e_i) e_i$ is the orthogonal projection of x onto the subspace generated by the basis vector e_i . If the basis is orthogonal but not necessarily orthonormal, then

$$x_i = \frac{x \cdot e_i}{e_i \cdot e_i} = \frac{x \cdot e_i}{\|e_i\|^2}.$$

All this is true even for an infinite orthonormal (or orthogonal) basis $(e_i)_{i \in I}$.



However, remember that every vector x is expressed as a linear combination

$$x = \sum_{i \in I} x_i e_i$$

where the family of scalars $(x_i)_{i \in I}$ has **finite support**, which means that $x_i = 0$ for all $i \in I - J$, where J is a finite set. Thus, even though the family $(\sin px)_{p \geq 1} \cup (\cos qx)_{q \geq 0}$ is orthogonal (it is not orthonormal, but becomes so if we divide every trigonometric function by $\sqrt{\pi}$, and 1 by $\sqrt{2\pi}$; we won't because it looks messy!), the fact that a function $f \in \mathcal{C}^0[-\pi, \pi]$ can be written as a Fourier series as

$$f(x) = a_0 + \sum_{k=1}^{\infty} (a_k \cos kx + b_k \sin kx)$$

does not mean that $(\sin px)_{p \geq 1} \cup (\cos qx)_{q \geq 0}$ is a basis of this vector space of functions, because in general, the families (a_k) and (b_k) **do not** have finite support! In order for this infinite linear combination to make sense, it is necessary to prove that the partial sums

$$a_0 + \sum_{k=1}^n (a_k \cos kx + b_k \sin kx)$$

of the series converge to a limit when n goes to infinity. This requires a topology on the space.

A very important property of Euclidean spaces of finite dimension is that the inner product induces a canonical bijection (i.e., independent of the choice of bases) between the vector space E and its dual E^* .

Given a Euclidean space E , for any vector $u \in E$, let $\varphi_u: E \rightarrow \mathbb{R}$ be the map defined such that

$$\varphi_u(v) = u \cdot v,$$

for all $v \in E$.

Since the inner product is bilinear, the map φ_u is a linear form in E^* . Thus, we have a map $\flat: E \rightarrow E^*$, defined such that

$$\flat(u) = \varphi_u.$$

Theorem 10.5. *Given a Euclidean space E , the map $\flat: E \rightarrow E^*$ defined such that*

$$\flat(u) = \varphi_u$$

is linear and injective. When E is also of finite dimension, the map $\flat: E \rightarrow E^$ is a canonical isomorphism.*

Proof. That $\flat: E \rightarrow E^*$ is a linear map follows immediately from the fact that the inner product is bilinear. If $\varphi_u = \varphi_v$, then $\varphi_u(w) = \varphi_v(w)$ for all $w \in E$, which by definition of φ_u means that

$$u \cdot w = v \cdot w$$

for all $w \in E$, which by bilinearity is equivalent to

$$(v - u) \cdot w = 0$$

for all $w \in E$, which implies that $u = v$, since the inner product is positive definite. Thus, $\flat: E \rightarrow E^*$ is injective. Finally, when E is of finite dimension n , we know that E^* is also of dimension n , and then $\flat: E \rightarrow E^*$ is bijective. \square

The inverse of the isomorphism $\flat: E \rightarrow E^*$ is denoted by $\sharp: E^* \rightarrow E$.

As a consequence of Theorem 10.5, if E is a Euclidean space of finite dimension, every linear form $f \in E^*$ corresponds to a unique $u \in E$ such that

$$f(v) = u \cdot v,$$

for every $v \in E$. In particular, if f is not the null form, the kernel of f , which is a hyperplane H , is precisely the set of vectors that are orthogonal to u .

Remarks:

- (1) The “musical map” $\flat: E \rightarrow E^*$ is not surjective when E has infinite dimension. The result can be salvaged by restricting our attention to continuous linear maps, and by assuming that the vector space E is a *Hilbert space* (i.e., E is a complete normed vector space w.r.t. the Euclidean norm). This is the famous “little” Riesz theorem (or Riesz representation theorem).

- (2) Theorem 10.5 still holds if the inner product on E is replaced by a nondegenerate symmetric bilinear form φ . We say that a symmetric bilinear form $\varphi: E \times E \rightarrow \mathbb{R}$ is *nondegenerate* if for every $u \in E$,

$$\text{if } \varphi(u, v) = 0 \text{ for all } v \in E, \text{ then } u = 0.$$

For example, the symmetric bilinear form on \mathbb{R}^4 defined such that

$$\varphi((x_1, x_2, x_3, x_4), (y_1, y_2, y_3, y_4)) = x_1y_1 + x_2y_2 + x_3y_3 - x_4y_4$$

is nondegenerate. However, there are nonnull vectors $u \in \mathbb{R}^4$ such that $\varphi(u, u) = 0$, which is impossible in a Euclidean space. Such vectors are called *isotropic*.

The existence of the isomorphism $b: E \rightarrow E^*$ is crucial to the existence of adjoint maps. The importance of adjoint maps stems from the fact that the linear maps arising in physical problems are often self-adjoint, which means that $f = f^*$. Moreover, self-adjoint maps can be diagonalized over orthonormal bases of eigenvectors. This is the key to the solution of many problems in mechanics, and engineering in general (see Strang [104]).

Let E be a Euclidean space of finite dimension n , and let $f: E \rightarrow E$ be a linear map. For every $u \in E$, the map

$$v \mapsto u \cdot f(v)$$

is clearly a linear form in E^* , and by Theorem 10.5, there is a unique vector in E denoted by $f^*(u)$ such that

$$f^*(u) \cdot v = u \cdot f(v),$$

for every $v \in E$. The following simple proposition shows that the map f^* is linear.

Proposition 10.6. *Given a Euclidean space E of finite dimension, for every linear map $f: E \rightarrow E$, there is a unique linear map $f^*: E \rightarrow E$ such that*

$$f^*(u) \cdot v = u \cdot f(v),$$

for all $u, v \in E$. The map f^ is called the adjoint of f (w.r.t. to the inner product).*

Proof. Given $u_1, u_2 \in E$, since the inner product is bilinear, we have

$$(u_1 + u_2) \cdot f(v) = u_1 \cdot f(v) + u_2 \cdot f(v),$$

for all $v \in E$, and

$$(f^*(u_1) + f^*(u_2)) \cdot v = f^*(u_1) \cdot v + f^*(u_2) \cdot v,$$

for all $v \in E$, and since by assumption,

$$f^*(u_1) \cdot v = u_1 \cdot f(v)$$

and

$$f^*(u_2) \cdot v = u_2 \cdot f(v),$$

for all $v \in E$, we get

$$(f^*(u_1) + f^*(u_2)) \cdot v = (u_1 + u_2) \cdot f(v),$$

for all $v \in E$. Since \flat is bijective, this implies that

$$f^*(u_1 + u_2) = f^*(u_1) + f^*(u_2).$$

Similarly,

$$(\lambda u) \cdot f(v) = \lambda(u \cdot f(v)),$$

for all $v \in E$, and

$$(\lambda f^*(u)) \cdot v = \lambda(f^*(u) \cdot v),$$

for all $v \in E$, and since by assumption,

$$f^*(u) \cdot v = u \cdot f(v),$$

for all $v \in E$, we get

$$(\lambda f^*(u)) \cdot v = (\lambda u) \cdot f(v),$$

for all $v \in E$. Since \flat is bijective, this implies that

$$f^*(\lambda u) = \lambda f^*(u).$$

Thus, f^* is indeed a linear map, and it is unique, since \flat is a bijection. □

Linear maps $f: E \rightarrow E$ such that $f = f^*$ are called *self-adjoint* maps. They play a very important role because they have real eigenvalues, and because orthonormal bases arise from their eigenvectors. Furthermore, many physical problems lead to self-adjoint linear maps (in the form of symmetric matrices).

Remark: Proposition 10.6 still holds if the inner product on E is replaced by a nondegenerate symmetric bilinear form φ .

Linear maps such that $f^{-1} = f^*$, or equivalently

$$f^* \circ f = f \circ f^* = \text{id},$$

also play an important role. They are *linear isometries*, or *isometries*. Rotations are special kinds of isometries. Another important class of linear maps are the linear maps satisfying the property

$$f^* \circ f = f \circ f^*,$$

called *normal linear maps*. We will see later on that normal maps can always be diagonalized over orthonormal bases of eigenvectors, but this will require using a Hermitian inner product (over \mathbb{C}).

Given two Euclidean spaces E and F , where the inner product on E is denoted by $\langle -, - \rangle_1$ and the inner product on F is denoted by $\langle -, - \rangle_2$, given any linear map $f: E \rightarrow F$, it is immediately verified that the proof of Proposition 10.6 can be adapted to show that there is a unique linear map $f^*: F \rightarrow E$ such that

$$\langle f(u), v \rangle_2 = \langle u, f^*(v) \rangle_1$$

for all $u \in E$ and all $v \in F$. The linear map f^* is also called the *adjoint of f* .

Remark: Given any basis for E and any basis for F , it is possible to characterize the matrix of the adjoint f^* of f in terms of the matrix of f , and the symmetric matrices defining the inner products. We will do so with respect to orthonormal bases. Also, since inner products are symmetric, the adjoint f^* of f is also characterized by

$$f(u) \cdot v = u \cdot f^*(v),$$

for all $u, v \in E$.

We can also use Theorem 10.5 to show that any Euclidean space of finite dimension has an orthonormal basis.

Proposition 10.7. *Given any nontrivial Euclidean space E of finite dimension $n \geq 1$, there is an orthonormal basis (u_1, \dots, u_n) for E .*

Proof. We proceed by induction on n . When $n = 1$, take any nonnull vector $v \in E$, which exists, since we assumed E nontrivial, and let

$$u = \frac{v}{\|v\|}.$$

If $n \geq 2$, again take any nonnull vector $v \in E$, and let

$$u_1 = \frac{v}{\|v\|}.$$

Consider the linear form φ_{u_1} associated with u_1 . Since $u_1 \neq 0$, by Theorem 10.5, the linear form φ_{u_1} is nonnull, and its kernel is a hyperplane H . Since $\varphi_{u_1}(w) = 0$ iff $u_1 \cdot w = 0$, the hyperplane H is the orthogonal complement of $\{u_1\}$. Furthermore, since $u_1 \neq 0$ and the inner product is positive definite, $u_1 \cdot u_1 \neq 0$, and thus, $u_1 \notin H$, which implies that $E = H \oplus \mathbb{R}u_1$. However, since E is of finite dimension n , the hyperplane H has dimension $n - 1$, and by the induction hypothesis, we can find an orthonormal basis (u_2, \dots, u_n) for H . Now, because H and the one dimensional space $\mathbb{R}u_1$ are orthogonal and $E = H \oplus \mathbb{R}u_1$, it is clear that (u_1, \dots, u_n) is an orthonormal basis for E . \square

There is a more constructive way of proving Proposition 10.7, using a procedure known as the *Gram–Schmidt orthonormalization procedure*. Among other things, the Gram–Schmidt orthonormalization procedure yields the *QR-decomposition for matrices*, an important tool in numerical methods.

Proposition 10.8. *Given any nontrivial Euclidean space E of finite dimension $n \geq 1$, from any basis (e_1, \dots, e_n) for E we can construct an orthonormal basis (u_1, \dots, u_n) for E , with the property that for every k , $1 \leq k \leq n$, the families (e_1, \dots, e_k) and (u_1, \dots, u_k) generate the same subspace.*

Proof. We proceed by induction on n . For $n = 1$, let

$$u_1 = \frac{e_1}{\|e_1\|}.$$

For $n \geq 2$, we also let

$$u_1 = \frac{e_1}{\|e_1\|},$$

and assuming that (u_1, \dots, u_k) is an orthonormal system that generates the same subspace as (e_1, \dots, e_k) , for every k with $1 \leq k < n$, we note that the vector

$$u'_{k+1} = e_{k+1} - \sum_{i=1}^k (e_{k+1} \cdot u_i) u_i$$

is nonnull, since otherwise, because (u_1, \dots, u_k) and (e_1, \dots, e_k) generate the same subspace, (e_1, \dots, e_{k+1}) would be linearly dependent, which is absurd, since (e_1, \dots, e_n) is a basis. Thus, the norm of the vector u'_{k+1} being nonzero, we use the following construction of the vectors u_k and u'_k :

$$u'_1 = e_1, \quad u_1 = \frac{u'_1}{\|u'_1\|},$$

and for the inductive step

$$u'_{k+1} = e_{k+1} - \sum_{i=1}^k (e_{k+1} \cdot u_i) u_i, \quad u_{k+1} = \frac{u'_{k+1}}{\|u'_{k+1}\|},$$

where $1 \leq k \leq n-1$. It is clear that $\|u_{k+1}\| = 1$, and since (u_1, \dots, u_k) is an orthonormal system, we have

$$u'_{k+1} \cdot u_i = e_{k+1} \cdot u_i - (e_{k+1} \cdot u_i) u_i \cdot u_i = e_{k+1} \cdot u_i - e_{k+1} \cdot u_i = 0,$$

for all i with $1 \leq i \leq k$. This shows that the family (u_1, \dots, u_{k+1}) is orthonormal, and since (u_1, \dots, u_k) and (e_1, \dots, e_k) generates the same subspace, it is clear from the definition of u_{k+1} that (u_1, \dots, u_{k+1}) and (e_1, \dots, e_{k+1}) generate the same subspace. This completes the induction step and the proof of the proposition. \square

Note that u'_{k+1} is obtained by subtracting from e_{k+1} the projection of e_{k+1} itself onto the orthonormal vectors u_1, \dots, u_k that have already been computed. Then, u'_{k+1} is normalized.

Remarks:

(1) The QR -decomposition can now be obtained very easily, but we postpone this until Section 10.4.

(2) We could compute u'_{k+1} using the formula

$$u'_{k+1} = e_{k+1} - \sum_{i=1}^k \left(\frac{e_{k+1} \cdot u'_i}{\|u'_i\|^2} \right) u'_i,$$

and normalize the vectors u'_k at the end. This time, we are subtracting from e_{k+1} the projection of e_{k+1} itself onto the orthogonal vectors u'_1, \dots, u'_k . This might be preferable when writing a computer program.

(3) The proof of Proposition 10.8 also works for a countably infinite basis for E , producing a countably infinite orthonormal basis.

Example 10.6. If we consider polynomials and the inner product

$$\langle f, g \rangle = \int_{-1}^1 f(t)g(t)dt,$$

applying the Gram–Schmidt orthonormalization procedure to the polynomials

$$1, x, x^2, \dots, x^n, \dots,$$

which form a basis of the polynomials in one variable with real coefficients, we get a family of orthonormal polynomials $Q_n(x)$ related to the *Legendre polynomials*.

The Legendre polynomials $P_n(x)$ have many nice properties. They are orthogonal, but their norm is not always 1. The Legendre polynomials $P_n(x)$ can be defined as follows. Letting f_n be the function

$$f_n(x) = (x^2 - 1)^n,$$

we define $P_n(x)$ as follows:

$$P_0(x) = 1, \quad \text{and} \quad P_n(x) = \frac{1}{2^n n!} f_n^{(n)}(x),$$

where $f_n^{(n)}$ is the n th derivative of f_n .

They can also be defined inductively as follows:

$$\begin{aligned} P_0(x) &= 1, \\ P_1(x) &= x, \\ P_{n+1}(x) &= \frac{2n+1}{n+1} x P_n(x) - \frac{n}{n+1} P_{n-1}(x). \end{aligned}$$

The polynomials Q_n are related to the Legendre polynomials P_n as follows:

$$Q_n(x) = \sqrt{\frac{2n+1}{2}} P_n(x).$$

Example 10.7. Consider polynomials over $[-1, 1]$, with the symmetric bilinear form

$$\langle f, g \rangle = \int_{-1}^1 \frac{1}{\sqrt{1-t^2}} f(t)g(t)dt.$$

We leave it as an exercise to prove that the above defines an inner product. It can be shown that the polynomials $T_n(x)$ given by

$$T_n(x) = \cos(n \arccos x), \quad n \geq 0,$$

(equivalently, with $x = \cos \theta$, we have $T_n(\cos \theta) = \cos(n\theta)$) are orthogonal with respect to the above inner product. These polynomials are the *Chebyshev polynomials*. Their norm is not equal to 1. Instead, we have

$$\langle T_n, T_n \rangle = \begin{cases} \frac{\pi}{2} & \text{if } n > 0, \\ \pi & \text{if } n = 0. \end{cases}$$

Using the identity $(\cos \theta + i \sin \theta)^n = \cos n\theta + i \sin n\theta$ and the binomial formula, we obtain the following expression for $T_n(x)$:

$$T_n(x) = \sum_{k=0}^{\lfloor n/2 \rfloor} \binom{n}{2k} (x^2 - 1)^k x^{n-2k}.$$

The Chebyshev polynomials are defined inductively as follows:

$$\begin{aligned} T_0(x) &= 1 \\ T_1(x) &= x \\ T_{n+1}(x) &= 2xT_n(x) - T_{n-1}(x), \quad n \geq 1. \end{aligned}$$

Using these recurrence equations, we can show that

$$T_n(x) = \frac{(x - \sqrt{x^2 - 1})^n + (x + \sqrt{x^2 - 1})^n}{2}.$$

The polynomial T_n has n distinct roots in the interval $[-1, 1]$. The Chebyshev polynomials play an important role in approximation theory. They are used as an approximation to a best polynomial approximation of a continuous function under the sup-norm (∞ -norm).

The inner products of the last two examples are special cases of an inner product of the form

$$\langle f, g \rangle = \int_{-1}^1 W(t)f(t)g(t)dt,$$

where $W(t)$ is a *weight function*. If W is a nonzero continuous function such that $W(x) \geq 0$ on $(-1, 1)$, then the above bilinear form is indeed positive definite. Families of orthogonal polynomials used in approximation theory and in physics arise by a suitable choice of the weight function W . Besides the previous two examples, the *Hermite polynomials* correspond to $W(x) = e^{-x^2}$, the *Laguerre polynomials* to $W(x) = e^{-x}$, and the *Jacobi polynomials* to $W(x) = (1-x)^\alpha(1+x)^\beta$, with $\alpha, \beta > -1$. Comprehensive treatments of orthogonal polynomials can be found in Lebedev [72], Sansone [91], and Andrews, Askey and Roy [2].

As a consequence of Proposition 10.7 (or Proposition 10.8), given any Euclidean space of finite dimension n , if (e_1, \dots, e_n) is an orthonormal basis for E , then for any two vectors $u = u_1e_1 + \dots + u_ne_n$ and $v = v_1e_1 + \dots + v_ne_n$, the inner product $u \cdot v$ is expressed as

$$u \cdot v = (u_1e_1 + \dots + u_ne_n) \cdot (v_1e_1 + \dots + v_ne_n) = \sum_{i=1}^n u_i v_i,$$

and the norm $\|u\|$ as

$$\|u\| = \|u_1e_1 + \dots + u_ne_n\| = \left(\sum_{i=1}^n u_i^2 \right)^{1/2}.$$

The fact that a Euclidean space always has an orthonormal basis implies that any Gram matrix G can be written as

$$G = Q^\top Q,$$

for some invertible matrix Q . Indeed, we know that in a change of basis matrix, a Gram matrix G becomes $G' = P^\top G P$. If the basis corresponding to G' is orthonormal, then $G' = I$, so $G = (P^{-1})^\top P^{-1}$.

We can also prove the following proposition regarding orthogonal spaces.

Proposition 10.9. *Given any nontrivial Euclidean space E of finite dimension $n \geq 1$, for any subspace F of dimension k , the orthogonal complement F^\perp of F has dimension $n - k$, and $E = F \oplus F^\perp$. Furthermore, we have $F^{\perp\perp} = F$.*

Proof. From Proposition 10.7, the subspace F has some orthonormal basis (u_1, \dots, u_k) . This linearly independent family (u_1, \dots, u_k) can be extended to a basis $(u_1, \dots, u_k, v_{k+1}, \dots, v_n)$,

and by Proposition 10.8, it can be converted to an orthonormal basis (u_1, \dots, u_n) , which contains (u_1, \dots, u_k) as an orthonormal basis of F . Now, any vector $w = w_1u_1 + \dots + w_nu_n \in E$ is orthogonal to F iff $w \cdot u_i = 0$, for every i , where $1 \leq i \leq k$, iff $w_i = 0$ for every i , where $1 \leq i \leq k$. Clearly, this shows that (u_{k+1}, \dots, u_n) is a basis of F^\perp , and thus $E = F \oplus F^\perp$, and F^\perp has dimension $n - k$. Similarly, any vector $w = w_1u_1 + \dots + w_nu_n \in E$ is orthogonal to F^\perp iff $w \cdot u_i = 0$, for every i , where $k+1 \leq i \leq n$, iff $w_i = 0$ for every i , where $k+1 \leq i \leq n$. Thus, (u_1, \dots, u_k) is a basis of $F^{\perp\perp}$, and $F^{\perp\perp} = F$. \square

10.3 Linear Isometries (Orthogonal Transformations)

In this section we consider linear maps between Euclidean spaces that preserve the Euclidean norm. These transformations, sometimes called *rigid motions*, play an important role in geometry.

Definition 10.3. Given any two nontrivial Euclidean spaces E and F of the same finite dimension n , a function $f: E \rightarrow F$ is an *orthogonal transformation*, or a *linear isometry*, if it is linear and

$$\|f(u)\| = \|u\|, \quad \text{for all } u \in E.$$

Remarks:

- (1) A linear isometry is often defined as a linear map such that

$$\|f(v) - f(u)\| = \|v - u\|,$$

for all $u, v \in E$. Since the map f is linear, the two definitions are equivalent. The second definition just focuses on preserving the distance between vectors.

- (2) Sometimes, a linear map satisfying the condition of Definition 10.3 is called a *metric map*, and a linear isometry is defined as a *bijective* metric map.

An isometry (without the word linear) is sometimes defined as a function $f: E \rightarrow F$ (not necessarily linear) such that

$$\|f(v) - f(u)\| = \|v - u\|,$$

for all $u, v \in E$, i.e., as a function that preserves the distance. This requirement turns out to be very strong. Indeed, the next proposition shows that all these definitions are equivalent when E and F are of finite dimension, and for functions such that $f(0) = 0$.

Proposition 10.10. *Given any two nontrivial Euclidean spaces E and F of the same finite dimension n , for every function $f: E \rightarrow F$, the following properties are equivalent:*

- (1) *f is a linear map and $\|f(u)\| = \|u\|$, for all $u \in E$;*

(2) $\|f(v) - f(u)\| = \|v - u\|$, for all $u, v \in E$, and $f(0) = 0$;

(3) $f(u) \cdot f(v) = u \cdot v$, for all $u, v \in E$.

Furthermore, such a map is bijective.

Proof. Clearly, (1) implies (2), since in (1) it is assumed that f is linear.

Assume that (2) holds. In fact, we shall prove a slightly stronger result. We prove that if

$$\|f(v) - f(u)\| = \|v - u\|$$

for all $u, v \in E$, then for any vector $\tau \in E$, the function $g: E \rightarrow F$ defined such that

$$g(u) = f(\tau + u) - f(\tau)$$

for all $u \in E$ is a linear map such that $g(0) = 0$ and (3) holds. Clearly, $g(0) = f(\tau) - f(\tau) = 0$.

Note that from the hypothesis

$$\|f(v) - f(u)\| = \|v - u\|$$

for all $u, v \in E$, we conclude that

$$\begin{aligned} \|g(v) - g(u)\| &= \|f(\tau + v) - f(\tau) - (f(\tau + u) - f(\tau))\|, \\ &= \|f(\tau + v) - f(\tau + u)\|, \\ &= \|\tau + v - (\tau + u)\|, \\ &= \|v - u\|, \end{aligned}$$

for all $u, v \in E$. Since $g(0) = 0$, by setting $u = 0$ in

$$\|g(v) - g(u)\| = \|v - u\|,$$

we get

$$\|g(v)\| = \|v\|$$

for all $v \in E$. In other words, g preserves both the distance and the norm.

To prove that g preserves the inner product, we use the simple fact that

$$2u \cdot v = \|u\|^2 + \|v\|^2 - \|u - v\|^2$$

for all $u, v \in E$. Then, since g preserves distance and norm, we have

$$\begin{aligned} 2g(u) \cdot g(v) &= \|g(u)\|^2 + \|g(v)\|^2 - \|g(u) - g(v)\|^2 \\ &= \|u\|^2 + \|v\|^2 - \|u - v\|^2 \\ &= 2u \cdot v, \end{aligned}$$

and thus $g(u) \cdot g(v) = u \cdot v$, for all $u, v \in E$, which is (3). In particular, if $f(0) = 0$, by letting $\tau = 0$, we have $g = f$, and f preserves the scalar product, i.e., (3) holds.

Now assume that (3) holds. Since E is of finite dimension, we can pick an orthonormal basis (e_1, \dots, e_n) for E . Since f preserves inner products, $(f(e_1), \dots, f(e_n))$ is also orthonormal, and since F also has dimension n , it is a basis of F . Then note that for any $u = u_1 e_1 + \dots + u_n e_n$, we have

$$u_i = u \cdot e_i,$$

for all i , $1 \leq i \leq n$. Thus, we have

$$f(u) = \sum_{i=1}^n (f(u) \cdot f(e_i)) f(e_i),$$

and since f preserves inner products, this shows that

$$f(u) = \sum_{i=1}^n (u \cdot e_i) f(e_i) = \sum_{i=1}^n u_i f(e_i),$$

which shows that f is linear. Obviously, f preserves the Euclidean norm, and (3) implies (1).

Finally, if $f(u) = f(v)$, then by linearity $f(v - u) = 0$, so that $\|f(v - u)\| = 0$, and since f preserves norms, we must have $\|v - u\| = 0$, and thus $u = v$. Thus, f is injective, and since E and F have the same finite dimension, f is bijective. \square

Remarks:

- (i) The dimension assumption is needed only to prove that (3) implies (1) when f is not known to be linear, and to prove that f is surjective, but the proof shows that (1) implies that f is injective.
- (ii) The implication that (3) implies (1) holds if we also assume that f is surjective, even if E has infinite dimension.

In (2), when f does not satisfy the condition $f(0) = 0$, the proof shows that f is an affine map. Indeed, taking any vector τ as an origin, the map g is linear, and

$$f(\tau + u) = f(\tau) + g(u) \quad \text{for all } u \in E.$$

From section 19.7, this shows that f is affine with associated linear map g .

This fact is worth recording as the following proposition.

Proposition 10.11. *Given any two nontrivial Euclidean spaces E and F of the same finite dimension n , for every function $f: E \rightarrow F$, if*

$$\|f(v) - f(u)\| = \|v - u\| \quad \text{for all } u, v \in E,$$

then f is an affine map, and its associated linear map g is an isometry.

In view of Proposition 10.10, we will drop the word “linear” in “linear isometry,” unless we wish to emphasize that we are dealing with a map between vector spaces.

We are now going to take a closer look at the isometries $f: E \rightarrow E$ of a Euclidean space of finite dimension.

10.4 The Orthogonal Group, Orthogonal Matrices

In this section we explore some of the basic properties of the orthogonal group and of orthogonal matrices.

Proposition 10.12. *Let E be any Euclidean space of finite dimension n , and let $f: E \rightarrow E$ be any linear map. The following properties hold:*

(1) *The linear map $f: E \rightarrow E$ is an isometry iff*

$$f \circ f^* = f^* \circ f = \text{id}.$$

(2) *For every orthonormal basis (e_1, \dots, e_n) of E , if the matrix of f is A , then the matrix of f^* is the transpose A^\top of A , and f is an isometry iff A satisfies the identities*

$$A A^\top = A^\top A = I_n,$$

where I_n denotes the identity matrix of order n , iff the columns of A form an orthonormal basis of E , iff the rows of A form an orthonormal basis of E .

Proof. (1) The linear map $f: E \rightarrow E$ is an isometry iff

$$f(u) \cdot f(v) = u \cdot v,$$

for all $u, v \in E$, iff

$$f^*(f(u)) \cdot v = f(u) \cdot f(v) = u \cdot v$$

for all $u, v \in E$, which implies

$$(f^*(f(u)) - u) \cdot v = 0$$

for all $u, v \in E$. Since the inner product is positive definite, we must have

$$f^*(f(u)) - u = 0$$

for all $u \in E$, that is,

$$f^* \circ f = f \circ f^* = \text{id}.$$

The converse is established by doing the above steps backward.

(2) If (e_1, \dots, e_n) is an orthonormal basis for E , let $A = (a_{ij})$ be the matrix of f , and let $B = (b_{ij})$ be the matrix of f^* . Since f^* is characterized by

$$f^*(u) \cdot v = u \cdot f(v)$$

for all $u, v \in E$, using the fact that if $w = w_1 e_1 + \dots + w_n e_n$ we have $w_k = w \cdot e_k$ for all k , $1 \leq k \leq n$, letting $u = e_i$ and $v = e_j$, we get

$$b_{ji} = f^*(e_i) \cdot e_j = e_i \cdot f(e_j) = a_{ij},$$

for all i, j , $1 \leq i, j \leq n$. Thus, $B = A^\top$. Now, if X and Y are arbitrary matrices over the basis (e_1, \dots, e_n) , denoting as usual the j th column of X by X^j , and similarly for Y , a simple calculation shows that

$$X^\top Y = (X^i \cdot Y^j)_{1 \leq i, j \leq n}.$$

Then it is immediately verified that if $X = Y = A$, then

$$A^\top A = A A^\top = I_n$$

iff the column vectors (A^1, \dots, A^n) form an orthonormal basis. Thus, from (1), we see that (2) is clear (also because the rows of A are the columns of A^\top). \square

Proposition 10.12 shows that the inverse of an isometry f is its adjoint f^* . Recall that the set of all real $n \times n$ matrices is denoted by $M_n(\mathbb{R})$. Proposition 10.12 also motivates the following definition.

Definition 10.4. A real $n \times n$ matrix is an *orthogonal matrix* if

$$A A^\top = A^\top A = I_n.$$

Remark: It is easy to show that the conditions $A A^\top = I_n$, $A^\top A = I_n$, and $A^{-1} = A^\top$, are equivalent. Given any two orthonormal bases (u_1, \dots, u_n) and (v_1, \dots, v_n) , if P is the change of basis matrix from (u_1, \dots, u_n) to (v_1, \dots, v_n) , since the columns of P are the coordinates of the vectors v_j with respect to the basis (u_1, \dots, u_n) , and since (v_1, \dots, v_n) is orthonormal, the columns of P are orthonormal, and by Proposition 10.12 (2), the matrix P is orthogonal.

The proof of Proposition 10.10 (3) also shows that if f is an isometry, then the image of an orthonormal basis (u_1, \dots, u_n) is an orthonormal basis. Students often ask why *orthogonal* matrices are not called *orthonormal* matrices, since their columns (and rows) are orthonormal bases! I have no good answer, but isometries do preserve orthogonality, and orthogonal matrices correspond to isometries.

Recall that the determinant $\det(f)$ of a linear map $f: E \rightarrow E$ is independent of the choice of a basis in E . Also, for every matrix $A \in M_n(\mathbb{R})$, we have $\det(A) = \det(A^\top)$, and for any two $n \times n$ matrices A and B , we have $\det(AB) = \det(A)\det(B)$. Then, if f is an isometry, and A is its matrix with respect to any orthonormal basis, $AA^\top = A^\top A = I_n$ implies that $\det(A)^2 = 1$, that is, either $\det(A) = 1$, or $\det(A) = -1$. It is also clear that the isometries of a Euclidean space of dimension n form a group, and that the isometries of determinant $+1$ form a subgroup. This leads to the following definition.

Definition 10.5. Given a Euclidean space E of dimension n , the set of isometries $f: E \rightarrow E$ forms a subgroup of $\mathbf{GL}(E)$ denoted by $\mathbf{O}(E)$, or $\mathbf{O}(n)$ when $E = \mathbb{R}^n$, called the *orthogonal group (of E)*. For every isometry f , we have $\det(f) = \pm 1$, where $\det(f)$ denotes the determinant of f . The isometries such that $\det(f) = 1$ are called *rotations, or proper isometries, or proper orthogonal transformations*, and they form a subgroup of the special linear group $\mathbf{SL}(E)$ (and of $\mathbf{O}(E)$), denoted by $\mathbf{SO}(E)$, or $\mathbf{SO}(n)$ when $E = \mathbb{R}^n$, called the *special orthogonal group (of E)*. The isometries such that $\det(f) = -1$ are called *improper isometries, or improper orthogonal transformations, or flip transformations*.

As an immediate corollary of the Gram–Schmidt orthonormalization procedure, we obtain the QR -decomposition for invertible matrices.

10.5 QR-Decomposition for Invertible Matrices

Now that we have the definition of an orthogonal matrix, we can explain how the Gram–Schmidt orthonormalization procedure immediately yields the QR -decomposition for matrices.

Proposition 10.13. *Given any real $n \times n$ matrix A , if A is invertible, then there is an orthogonal matrix Q and an upper triangular matrix R with positive diagonal entries such that $A = QR$.*

Proof. We can view the columns of A as vectors A^1, \dots, A^n in \mathbb{E}^n . If A is invertible, then they are linearly independent, and we can apply Proposition 10.8 to produce an orthonormal basis using the Gram–Schmidt orthonormalization procedure. Recall that we construct vectors Q^k and Q'^k as follows:

$$Q'^1 = A^1, \quad Q^1 = \frac{Q'^1}{\|Q'^1\|},$$

and for the inductive step

$$Q'^{k+1} = A^{k+1} - \sum_{i=1}^k (A^{k+1} \cdot Q^i) Q^i, \quad Q^{k+1} = \frac{Q'^{k+1}}{\|Q'^{k+1}\|},$$

where $1 \leq k \leq n-1$. If we express the vectors A^k in terms of the Q^i and Q'^i , we get the triangular system

$$\begin{aligned} A^1 &= \|Q'^1\| Q^1, \\ &\vdots \\ A^j &= (A^j \cdot Q^1) Q^1 + \cdots + (A^j \cdot Q^i) Q^i + \cdots + \|Q'^j\| Q^j, \\ &\vdots \\ A^n &= (A^n \cdot Q^1) Q^1 + \cdots + (A^n \cdot Q^{n-1}) Q^{n-1} + \|Q'^n\| Q^n. \end{aligned}$$

Letting $r_{kk} = \|Q'^k\|$, and $r_{ij} = A^j \cdot Q^i$ (the reversal of i and j on the right-hand side *is* intentional!), where $1 \leq k \leq n$, $2 \leq j \leq n$, and $1 \leq i \leq j-1$, and letting q_{ij} be the i th component of Q^j , we note that a_{ij} , the i th component of A^j , is given by

$$a_{ij} = r_{1j}q_{i1} + \cdots + r_{ij}q_{ii} + \cdots + r_{jj}q_{ij} = q_{i1}r_{1j} + \cdots + q_{ii}r_{ij} + \cdots + q_{ij}r_{jj}.$$

If we let $Q = (q_{ij})$, the matrix whose columns are the components of the Q^j , and $R = (r_{ij})$, the above equations show that $A = QR$, where R is upper triangular. The diagonal entries $r_{kk} = \|Q'^k\| = A^k \cdot Q^k$ are indeed positive. \square

The reader should try the above procedure on some concrete examples for 2×2 and 3×3 matrices.

Remarks:

- (1) Because the diagonal entries of R are positive, it can be shown that Q and R are unique.
- (2) The QR -decomposition holds even when A is not invertible. In this case, R has some zero on the diagonal. However, a different proof is needed. We will give a nice proof using Householder matrices (see Proposition 11.3, and also Strang [104, 105], Golub and Van Loan [49], Trefethen and Bau [110], Demmel [27], Kincaid and Cheney [63], or Ciarlet [24]).

Example 10.8. Consider the matrix

$$A = \begin{pmatrix} 0 & 0 & 5 \\ 0 & 4 & 1 \\ 1 & 1 & 1 \end{pmatrix}.$$

We leave as an exercise to show that $A = QR$, with

$$Q = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix} \quad \text{and} \quad R = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 4 & 1 \\ 0 & 0 & 5 \end{pmatrix}.$$

Example 10.9. Another example of QR -decomposition is

$$A = \begin{pmatrix} 1 & 1 & 2 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix} = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 0 & 0 & 1 \\ 1/\sqrt{2} & -1/\sqrt{2} & 0 \end{pmatrix} \begin{pmatrix} \sqrt{2} & 1/\sqrt{2} & \sqrt{2} \\ 0 & 1/\sqrt{2} & \sqrt{2} \\ 0 & 0 & 1 \end{pmatrix}.$$

The QR -decomposition yields a rather efficient and numerically stable method for solving systems of linear equations. Indeed, given a system $Ax = b$, where A is an $n \times n$ invertible matrix, writing $A = QR$, since Q is orthogonal, we get

$$Rx = Q^\top b,$$

and since R is upper triangular, we can solve it by Gaussian elimination, by solving for the last variable x_n first, substituting its value into the system, then solving for x_{n-1} , etc. The QR -decomposition is also very useful in solving least squares problems (we will come back to this later on), and for finding eigenvalues. It can be easily adapted to the case where A is a rectangular $m \times n$ matrix with independent columns (thus, $n \leq m$). In this case, Q is not quite orthogonal. It is an $m \times n$ matrix whose columns are orthogonal, and R is an invertible $n \times n$ upper triangular matrix with positive diagonal entries. For more on QR , see Strang [104, 105], Golub and Van Loan [49], Demmel [27], Trefethen and Bau [110], or Serre [96].

It should also be said that the Gram–Schmidt orthonormalization procedure that we have presented is not very stable numerically, and instead, one should use the *modified Gram–Schmidt method*. To compute Q^{k+1} , instead of projecting A^{k+1} onto Q^1, \dots, Q^k in a single step, it is better to perform k projections. We compute $Q_1^{k+1}, Q_2^{k+1}, \dots, Q_k^{k+1}$ as follows:

$$\begin{aligned} Q_1^{k+1} &= A^{k+1} - (A^{k+1} \cdot Q^1) Q^1, \\ Q_{i+1}^{k+1} &= Q_i^{k+1} - (Q_i^{k+1} \cdot Q^{i+1}) Q^{i+1}, \end{aligned}$$

where $1 \leq i \leq k-1$. It is easily shown that $Q'^{k+1} = Q_k^{k+1}$. The reader is urged to code this method.

A somewhat surprising consequence of the QR -decomposition is a famous determinantal inequality due to Hadamard.

Proposition 10.14. (*Hadamard*) For any real $n \times n$ matrix $A = (a_{ij})$, we have

$$|\det(A)| \leq \prod_{i=1}^n \left(\sum_{j=1}^n a_{ij}^2 \right)^{1/2} \quad \text{and} \quad |\det(A)| \leq \prod_{j=1}^n \left(\sum_{i=1}^n a_{ij}^2 \right)^{1/2}.$$

Moreover, equality holds iff either A has a zero column in the left inequality or a zero row in the right inequality, or A is orthogonal.

Proof. If $\det(A) = 0$, then the inequality is trivial. In addition, if the righthand side is also 0, then either some column or some row is zero. If $\det(A) \neq 0$, then we can factor A as $A = QR$, with Q is orthogonal and $R = (r_{ij})$ upper triangular with positive diagonal entries. Then, since Q is orthogonal $\det(Q) = \pm 1$, so

$$|\det(A)| = |\det(Q)| |\det(R)| = \prod_{j=1}^n r_{jj}.$$

Now, as Q is orthogonal, it preserves the Euclidean norm, so

$$\sum_{i=1}^n a_{ij}^2 = \|A^j\|_2^2 = \|QR^j\|_2^2 = \|R^j\|_2^2 = \sum_{i=1}^n r_{ij}^2 \geq r_{jj}^2,$$

which implies that

$$|\det(A)| = \prod_{j=1}^n r_{jj} \leq \prod_{j=1}^n \|R^j\|_2 \leq \prod_{j=1}^n \left(\sum_{i=1}^n a_{ij}^2 \right)^{1/2}.$$

The other inequality is obtained by replacing A by A^\top . Finally, if $\det(A) \neq 0$ and equality holds, then we must have

$$r_{jj} = \|A^j\|_2, \quad 1 \leq j \leq n,$$

which can only occur is R is orthogonal. □

Another version of Hadamard's inequality applies to symmetric positive semidefinite matrices.

Proposition 10.15. (*Hadamard*) *For any real $n \times n$ matrix $A = (a_{ij})$, if A is symmetric positive semidefinite, then we have*

$$\det(A) \leq \prod_{i=1}^n a_{ii}.$$

Moreover, if A is positive definite, then equality holds iff A is a diagonal matrix.

Proof. If $\det(A) = 0$, the inequality is trivial. Otherwise, A is positive definite, and by Theorem 6.10 (the Cholesky Factorization), there is a unique upper triangular matrix B with positive diagonal entries such that

$$A = B^\top B.$$

Thus, $\det(A) = \det(B^\top B) = \det(B^\top) \det(B) = \det(B)^2$. If we apply the Hadamard inequality (Proposition 10.15) to B , we obtain

$$\det(B) \leq \prod_{i=1}^n \left(\sum_{j=1}^n b_{ij}^2 \right)^{1/2}. \quad (*)$$

However, the diagonal entries a_{ii} of $A = B^\top B$ are precisely the square norms $\|B^i\|_2^2 = \sum_{j=1}^n b_{ij}^2$, so by squaring (*), we obtain

$$\det(A) = \det(B)^2 \leq \prod_{i=1}^n \left(\sum_{j=1}^n b_{ij}^2 \right) = \prod_{i=1}^n a_{ii}.$$

If $\det(A) \neq 0$ and equality holds, then B must be orthogonal, which implies that B is a diagonal matrix, and so is A . \square

We derived the second Hadamard inequality (Proposition 10.15) from the first (Proposition 10.14). We leave it as an exercise to prove that the first Hadamard inequality can be deduced from the second Hadamard inequality.

10.6 Some Applications of Euclidean Geometry

Euclidean geometry has applications in computational geometry, in particular Voronoi diagrams and Delaunay triangulations. In turn, Voronoi diagrams have applications in motion planning (see O'Rourke [87]).

Euclidean geometry also has applications to matrix analysis. Recall that a real $n \times n$ matrix A is *symmetric* if it is equal to its transpose A^\top . One of the most important properties of symmetric matrices is that they have real eigenvalues and that they can be diagonalized by an orthogonal matrix (see Chapter 13). This means that for every symmetric matrix A , there is a diagonal matrix D and an orthogonal matrix P such that

$$A = PDP^\top.$$

Even though it is not always possible to diagonalize an arbitrary matrix, there are various decompositions involving orthogonal matrices that are of great practical interest. For example, for every real matrix A , there is the *QR-decomposition*, which says that a real matrix A can be expressed as

$$A = QR,$$

where Q is orthogonal and R is an upper triangular matrix. This can be obtained from the Gram–Schmidt orthonormalization procedure, as we saw in Section 10.5, or better, using Householder matrices, as shown in Section 11.2. There is also the *polar decomposition*, which says that a real matrix A can be expressed as

$$A = QS,$$

where Q is orthogonal and S is symmetric positive semidefinite (which means that the eigenvalues of S are nonnegative). Such a decomposition is important in continuum mechanics and in robotics, since it separates stretching from rotation. Finally, there is the wonderful

singular value decomposition, abbreviated as SVD, which says that a real matrix A can be expressed as

$$A = VDU^\top,$$

where U and V are orthogonal and D is a diagonal matrix with nonnegative entries (see Chapter 16). This decomposition leads to the notion of *pseudo-inverse*, which has many applications in engineering (least squares solutions, etc). For an excellent presentation of all these notions, we highly recommend Strang [105, 104], Golub and Van Loan [49], Demmel [27], Serre [96], and Trefethen and Bau [110].

The method of least squares, invented by Gauss and Legendre around 1800, is another great application of Euclidean geometry. Roughly speaking, the method is used to solve inconsistent linear systems $Ax = b$, where the number of equations is greater than the number of variables. Since this is generally impossible, the method of least squares consists in finding a solution x minimizing the Euclidean norm $\|Ax - b\|^2$, that is, the sum of the squares of the “errors.” It turns out that there is always a unique solution x^+ of smallest norm minimizing $\|Ax - b\|^2$, and that it is a solution of the square system

$$A^\top Ax = A^\top b,$$

called the system of *normal equations*. The solution x^+ can be found either by using the QR -decomposition in terms of Householder transformations, or by using the notion of pseudo-inverse of a matrix. The pseudo-inverse can be computed using the SVD decomposition. Least squares methods are used extensively in computer vision. More details on the method of least squares and pseudo-inverses can be found in Chapter 17.

10.7 Summary

The main concepts and results of this chapter are listed below:

- Bilinear forms; *positive definite* bilinear forms.
- *inner products*, *scalar products*, *Euclidean spaces*.
- *quadratic form* associated with a bilinear form.
- The Euclidean space \mathbb{E}^n .
- The *polar form* of a quadratic form.
- *Gram matrix* associated with an inner product.
- The *Cauchy–Schwarz inequality*; the *Minkowski inequality*.
- The *parallelogram law*.

- *Orthogonality, orthogonal complement F^\perp ; orthonormal family.*
- The *musical isomorphisms* $\flat: E \rightarrow E^*$ and $\sharp: E^* \rightarrow E$ (when E is finite-dimensional); Theorem 10.5.
- The *adjoint* of a linear map (with respect to an inner product).
- Existence of an orthonormal basis in a finite-dimensional Euclidean space (Proposition 10.7).
- The *Gram–Schmidt orthonormalization procedure* (Proposition 10.8).
- The *Legendre* and the *Chebyshev* polynomials.
- *Linear isometries (orthogonal transformations, rigid motions).*
- The *orthogonal group, orthogonal matrices.*
- The matrix representing the adjoint f^* of a linear map f is the transpose of the matrix representing f .
- The *orthogonal group* $\mathbf{O}(n)$ and the *special orthogonal group* $\mathbf{SO}(n)$.
- *QR-decomposition* for invertible matrices.
- The *Hadamard inequality* for arbitrary real matrices.
- The *Hadamard inequality* for symmetric positive semidefinite matrices.

Chapter 11

QR -Decomposition for Arbitrary Matrices

11.1 Orthogonal Reflections

Hyperplane reflections are represented by matrices called Householder matrices. These matrices play an important role in numerical methods, for instance for solving systems of linear equations, solving least squares problems, for computing eigenvalues, and for transforming a symmetric matrix into a tridiagonal matrix. We prove a simple geometric lemma that immediately yields the QR -decomposition of arbitrary matrices in terms of Householder matrices.

Orthogonal symmetries are a very important example of isometries. First let us review the definition of projections. Given a vector space E , let F and G be subspaces of E that form a direct sum $E = F \oplus G$. Since every $u \in E$ can be written uniquely as $u = v + w$, where $v \in F$ and $w \in G$, we can define the two *projections* $p_F: E \rightarrow F$ and $p_G: E \rightarrow G$ such that $p_F(u) = v$ and $p_G(u) = w$. It is immediately verified that p_G and p_F are linear maps, and that $p_F^2 = p_F$, $p_G^2 = p_G$, $p_F \circ p_G = p_G \circ p_F = 0$, and $p_F + p_G = \text{id}$.

Definition 11.1. Given a vector space E , for any two subspaces F and G that form a direct sum $E = F \oplus G$, the *symmetry (or reflection) with respect to F and parallel to G* is the linear map $s: E \rightarrow E$ defined such that

$$s(u) = 2p_F(u) - u,$$

for every $u \in E$.

Because $p_F + p_G = \text{id}$, note that we also have

$$s(u) = p_F(u) - p_G(u)$$

and

$$s(u) = u - 2p_G(u),$$

$s^2 = \text{id}$, s is the identity on F , and $s = -\text{id}$ on G . We now assume that E is a Euclidean space of finite dimension.

Definition 11.2. Let E be a Euclidean space of finite dimension n . For any two subspaces F and G , if F and G form a direct sum $E = F \oplus G$ and F and G are orthogonal, i.e., $F = G^\perp$, the *orthogonal symmetry (or reflection) with respect to F and parallel to G* is the linear map $s: E \rightarrow E$ defined such that

$$s(u) = 2p_F(u) - u,$$

for every $u \in E$. When F is a hyperplane, we call s a *hyperplane symmetry with respect to F (or reflection about F)*, and when G is a plane (and thus $\dim(F) = n - 2$), we call s a *flip about F* .

For any two vectors $u, v \in E$, it is easily verified using the bilinearity of the inner product that

$$\|u + v\|^2 - \|u - v\|^2 = 4(u \cdot v).$$

Then, since

$$u = p_F(u) + p_G(u)$$

and

$$s(u) = p_F(u) - p_G(u),$$

since F and G are orthogonal, it follows that

$$p_F(u) \cdot p_G(v) = 0,$$

and thus,

$$\|s(u)\| = \|u\|,$$

so that s is an isometry.

Using Proposition 10.8, it is possible to find an orthonormal basis (e_1, \dots, e_n) of E consisting of an orthonormal basis of F and an orthonormal basis of G . Assume that F has dimension p , so that G has dimension $n - p$. With respect to the orthonormal basis (e_1, \dots, e_n) , the symmetry s has a matrix of the form

$$\begin{pmatrix} I_p & 0 \\ 0 & -I_{n-p} \end{pmatrix}.$$

Thus, $\det(s) = (-1)^{n-p}$, and s is a rotation iff $n - p$ is even. In particular, when F is a hyperplane H , we have $p = n - 1$ and $n - p = 1$, so that s is an improper orthogonal transformation. When $F = \{0\}$, we have $s = -\text{id}$, which is called the *symmetry with respect to the origin*. The symmetry with respect to the origin is a rotation iff n is even, and an improper orthogonal transformation iff n is odd. When n is odd, we observe that every improper orthogonal transformation is the composition of a rotation with the symmetry

with respect to the origin. When G is a plane, $p = n - 2$, and $\det(s) = (-1)^2 = 1$, so that a flip about F is a rotation. In particular, when $n = 3$, F is a line, and a flip about the line F is indeed a rotation of measure π .

Remark: Given any two orthogonal subspaces F, G forming a direct sum $E = F \oplus G$, let f be the symmetry with respect to F and parallel to G , and let g be the symmetry with respect to G and parallel to F . We leave as an exercise to show that

$$f \circ g = g \circ f = -\text{id}.$$

When $F = H$ is a hyperplane, we can give an explicit formula for $s(u)$ in terms of any nonnull vector w orthogonal to H . Indeed, from

$$u = p_H(u) + p_G(u),$$

since $p_G(u) \in G$ and G is spanned by w , which is orthogonal to H , we have

$$p_G(u) = \lambda w$$

for some $\lambda \in \mathbb{R}$, and we get

$$u \cdot w = \lambda \|w\|^2,$$

and thus

$$p_G(u) = \frac{(u \cdot w)}{\|w\|^2} w.$$

Since

$$s(u) = u - 2p_G(u),$$

we get

$$s(u) = u - 2 \frac{(u \cdot w)}{\|w\|^2} w.$$

Such reflections are represented by matrices called *Householder matrices*, and they play an important role in numerical matrix analysis (see Kincaid and Cheney [63] or Ciarlet [24]). Householder matrices are symmetric and orthogonal. It is easily checked that over an orthonormal basis (e_1, \dots, e_n) , a hyperplane reflection about a hyperplane H orthogonal to a nonnull vector w is represented by the matrix

$$H = I_n - 2 \frac{WW^\top}{\|W\|^2} = I_n - 2 \frac{WW^\top}{W^\top W},$$

where W is the column vector of the coordinates of w over the basis (e_1, \dots, e_n) , and I_n is the identity $n \times n$ matrix. Since

$$p_G(u) = \frac{(u \cdot w)}{\|w\|^2} w,$$

the matrix representing p_G is

$$\frac{WW^\top}{W^\top W},$$

and since $p_H + p_G = \text{id}$, the matrix representing p_H is

$$I_n - \frac{WW^\top}{W^\top W}.$$

These formulae can be used to derive a formula for a rotation of \mathbb{R}^3 , given the direction w of its axis of rotation and given the angle θ of rotation.

The following fact is the key to the proof that every isometry can be decomposed as a product of reflections.

Proposition 11.1. *Let E be any nontrivial Euclidean space. For any two vectors $u, v \in E$, if $\|u\| = \|v\|$, then there is a hyperplane H such that the reflection s about H maps u to v , and if $u \neq v$, then this reflection is unique.*

Proof. If $u = v$, then any hyperplane containing u does the job. Otherwise, we must have $H = \{v - u\}^\perp$, and by the above formula,

$$s(u) = u - 2 \frac{(u \cdot (v - u))}{\|(v - u)\|^2} (v - u) = u + \frac{2\|u\|^2 - 2u \cdot v}{\|(v - u)\|^2} (v - u),$$

and since

$$\|(v - u)\|^2 = \|u\|^2 + \|v\|^2 - 2u \cdot v$$

and $\|u\| = \|v\|$, we have

$$\|(v - u)\|^2 = 2\|u\|^2 - 2u \cdot v,$$

and thus, $s(u) = v$. □



If E is a complex vector space and the inner product is Hermitian, Proposition 11.1 is false. The problem is that the vector $v - u$ does not work unless the inner product $u \cdot v$ is real! The proposition can be salvaged enough to yield the QR -decomposition in terms of Householder transformations; see Gallier [44].

We now show that hyperplane reflections can be used to obtain another proof of the QR -decomposition.

11.2 QR -Decomposition Using Householder Matrices

First, we state the result geometrically. When translated in terms of Householder matrices, we obtain the fact advertised earlier that every matrix (not necessarily invertible) has a QR -decomposition.

Proposition 11.2. *Let E be a nontrivial Euclidean space of dimension n . For any orthonormal basis (e_1, \dots, e_n) and for any n -tuple of vectors (v_1, \dots, v_n) , there is a sequence of n isometries h_1, \dots, h_n such that h_i is a hyperplane reflection or the identity, and if (r_1, \dots, r_n) are the vectors given by*

$$r_j = h_n \circ \dots \circ h_2 \circ h_1(v_j),$$

then every r_j is a linear combination of the vectors (e_1, \dots, e_j) , $1 \leq j \leq n$. Equivalently, the matrix R whose columns are the components of the r_j over the basis (e_1, \dots, e_n) is an upper triangular matrix. Furthermore, the h_i can be chosen so that the diagonal entries of R are nonnegative.

Proof. We proceed by induction on n . For $n = 1$, we have $v_1 = \lambda e_1$ for some $\lambda \in \mathbb{R}$. If $\lambda \geq 0$, we let $h_1 = \text{id}$, else if $\lambda < 0$, we let $h_1 = -\text{id}$, the reflection about the origin.

For $n \geq 2$, we first have to find h_1 . Let

$$r_{1,1} = \|v_1\|.$$

If $v_1 = r_{1,1}e_1$, we let $h_1 = \text{id}$. Otherwise, there is a unique hyperplane reflection h_1 such that

$$h_1(v_1) = r_{1,1}e_1,$$

defined such that

$$h_1(u) = u - 2 \frac{(u \cdot w_1)}{\|w_1\|^2} w_1$$

for all $u \in E$, where

$$w_1 = r_{1,1}e_1 - v_1.$$

The map h_1 is the reflection about the hyperplane H_1 orthogonal to the vector $w_1 = r_{1,1}e_1 - v_1$. Letting

$$r_1 = h_1(v_1) = r_{1,1}e_1,$$

it is obvious that r_1 belongs to the subspace spanned by e_1 , and $r_{1,1} = \|v_1\|$ is nonnegative.

Next, assume that we have found k linear maps h_1, \dots, h_k , hyperplane reflections or the identity, where $1 \leq k \leq n - 1$, such that if (r_1, \dots, r_k) are the vectors given by

$$r_j = h_k \circ \dots \circ h_2 \circ h_1(v_j),$$

then every r_j is a linear combination of the vectors (e_1, \dots, e_j) , $1 \leq j \leq k$. The vectors (e_1, \dots, e_k) form a basis for the subspace denoted by U'_k , the vectors (e_{k+1}, \dots, e_n) form a basis for the subspace denoted by U''_k , the subspaces U'_k and U''_k are orthogonal, and $E = U'_k \oplus U''_k$. Let

$$u_{k+1} = h_k \circ \dots \circ h_2 \circ h_1(v_{k+1}).$$

We can write

$$u_{k+1} = u'_{k+1} + u''_{k+1},$$

where $u'_{k+1} \in U'_k$ and $u''_{k+1} \in U''_k$. Let

$$r_{k+1,k+1} = \|u''_{k+1}\|.$$

If $u''_{k+1} = r_{k+1,k+1} e_{k+1}$, we let $h_{k+1} = \text{id}$. Otherwise, there is a unique hyperplane reflection h_{k+1} such that

$$h_{k+1}(u''_{k+1}) = r_{k+1,k+1} e_{k+1},$$

defined such that

$$h_{k+1}(u) = u - 2 \frac{(u \cdot w_{k+1})}{\|w_{k+1}\|^2} w_{k+1}$$

for all $u \in E$, where

$$w_{k+1} = r_{k+1,k+1} e_{k+1} - u''_{k+1}.$$

The map h_{k+1} is the reflection about the hyperplane H_{k+1} orthogonal to the vector $w_{k+1} = r_{k+1,k+1} e_{k+1} - u''_{k+1}$. However, since $u''_{k+1}, e_{k+1} \in U''_k$ and U'_k is orthogonal to U''_k , the subspace U'_k is contained in H_{k+1} , and thus, the vectors (r_1, \dots, r_k) and u'_{k+1} , which belong to U'_k , are invariant under h_{k+1} . This proves that

$$h_{k+1}(u_{k+1}) = h_{k+1}(u'_{k+1}) + h_{k+1}(u''_{k+1}) = u'_{k+1} + r_{k+1,k+1} e_{k+1}$$

is a linear combination of (e_1, \dots, e_{k+1}) . Letting

$$r_{k+1} = h_{k+1}(u_{k+1}) = u'_{k+1} + r_{k+1,k+1} e_{k+1},$$

since $u_{k+1} = h_k \circ \dots \circ h_2 \circ h_1(v_{k+1})$, the vector

$$r_{k+1} = h_{k+1} \circ \dots \circ h_2 \circ h_1(v_{k+1})$$

is a linear combination of (e_1, \dots, e_{k+1}) . The coefficient of r_{k+1} over e_{k+1} is $r_{k+1,k+1} = \|u''_{k+1}\|$, which is nonnegative. This concludes the induction step, and thus the proof. \square

Remarks:

- (1) Since every h_i is a hyperplane reflection or the identity,

$$\rho = h_n \circ \dots \circ h_2 \circ h_1$$

is an isometry.

- (2) If we allow negative diagonal entries in R , the last isometry h_n may be omitted.

- (3) Instead of picking $r_{k,k} = \|u_k''\|$, which means that

$$w_k = r_{k,k} e_k - u_k'',$$

where $1 \leq k \leq n$, it might be preferable to pick $r_{k,k} = -\|u_k''\|$ if this makes $\|w_k\|^2$ larger, in which case

$$w_k = r_{k,k} e_k + u_k''.$$

Indeed, since the definition of h_k involves division by $\|w_k\|^2$, it is desirable to avoid division by very small numbers.

- (4) The method also applies to any m -tuple of vectors (v_1, \dots, v_m) , where m is not necessarily equal to n (the dimension of E). In this case, R is an upper triangular $n \times m$ matrix we leave the minor adjustments to the method as an exercise to the reader (if $m > n$, the last $m - n$ vectors are unchanged).

Proposition 11.2 directly yields the QR -decomposition in terms of Householder transformations (see Strang [104, 105], Golub and Van Loan [49], Trefethen and Bau [110], Kincaid and Cheney [63], or Ciarlet [24]).

Theorem 11.3. *For every real $n \times n$ matrix A , there is a sequence H_1, \dots, H_n of matrices, where each H_i is either a Householder matrix or the identity, and an upper triangular matrix R such that*

$$R = H_n \cdots H_2 H_1 A.$$

As a corollary, there is a pair of matrices Q, R , where Q is orthogonal and R is upper triangular, such that $A = QR$ (a QR -decomposition of A). Furthermore, R can be chosen so that its diagonal entries are nonnegative.

Proof. The j th column of A can be viewed as a vector v_j over the canonical basis (e_1, \dots, e_n) of \mathbb{E}^n (where $(e_j)_i = 1$ if $i = j$, and 0 otherwise, $1 \leq i, j \leq n$). Applying Proposition 11.2 to (v_1, \dots, v_n) , there is a sequence of n isometries h_1, \dots, h_n such that h_i is a hyperplane reflection or the identity, and if (r_1, \dots, r_n) are the vectors given by

$$r_j = h_n \circ \cdots \circ h_2 \circ h_1(v_j),$$

then every r_j is a linear combination of the vectors (e_1, \dots, e_j) , $1 \leq j \leq n$. Letting R be the matrix whose columns are the vectors r_j , and H_i the matrix associated with h_i , it is clear that

$$R = H_n \cdots H_2 H_1 A,$$

where R is upper triangular and every H_i is either a Householder matrix or the identity. However, $h_i \circ h_i = \text{id}$ for all i , $1 \leq i \leq n$, and so

$$v_j = h_1 \circ h_2 \circ \cdots \circ h_n(r_j)$$

for all j , $1 \leq j \leq n$. But $\rho = h_1 \circ h_2 \circ \cdots \circ h_n$ is an isometry represented by the orthogonal matrix $Q = H_1 H_2 \cdots H_n$. It is clear that $A = QR$, where R is upper triangular. As we noted in Proposition 11.2, the diagonal entries of R can be chosen to be nonnegative. \square

Remarks:

(1) Letting

$$A_{k+1} = H_k \cdots H_2 H_1 A,$$

with $A_1 = A$, $1 \leq k \leq n$, the proof of Proposition 11.2 can be interpreted in terms of the computation of the sequence of matrices $A_1, \dots, A_{n+1} = R$. The matrix A_{k+1} has the shape

$$A_{k+1} = \begin{pmatrix} \times & \times & \times & u_1^{k+1} & \times & \times & \times & \times \\ 0 & \times & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \times & u_k^{k+1} & \times & \times & \times & \times \\ 0 & 0 & 0 & u_{k+1}^{k+1} & \times & \times & \times & \times \\ 0 & 0 & 0 & u_{k+2}^{k+1} & \times & \times & \times & \times \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & u_{n-1}^{k+1} & \times & \times & \times & \times \\ 0 & 0 & 0 & u_n^{k+1} & \times & \times & \times & \times \end{pmatrix},$$

where the $(k+1)$ th column of the matrix is the vector

$$u_{k+1} = h_k \circ \cdots \circ h_2 \circ h_1(v_{k+1}),$$

and thus

$$u'_{k+1} = (u_1^{k+1}, \dots, u_k^{k+1})$$

and

$$u''_{k+1} = (u_{k+1}^{k+1}, u_{k+2}^{k+1}, \dots, u_n^{k+1}).$$

If the last $n - k - 1$ entries in column $k+1$ are all zero, there is nothing to do, and we let $H_{k+1} = I$. Otherwise, we kill these $n - k - 1$ entries by multiplying A_{k+1} on the left by the Householder matrix H_{k+1} sending

$$(0, \dots, 0, u_{k+1}^{k+1}, \dots, u_n^{k+1}) \quad \text{to} \quad (0, \dots, 0, r_{k+1,k+1}, 0, \dots, 0),$$

where $r_{k+1,k+1} = \|(u_{k+1}^{k+1}, \dots, u_n^{k+1})\|$.

(2) If A is invertible and the diagonal entries of R are positive, it can be shown that Q and R are unique.

(3) If we allow negative diagonal entries in R , the matrix H_n may be omitted ($H_n = I$).

(4) The method allows the computation of the determinant of A . We have

$$\det(A) = (-1)^m r_{1,1} \cdots r_{n,n},$$

where m is the number of Householder matrices (not the identity) among the H_i .

- (5) The “condition number” of the matrix A is preserved (see Strang [105], Golub and Van Loan [49], Trefethen and Bau [110], Kincaid and Cheney [63], or Ciarlet [24]). This is very good for numerical stability.
- (6) The method also applies to a rectangular $m \times n$ matrix. In this case, R is also an $m \times n$ matrix (and it is upper triangular).

11.3 Summary

The main concepts and results of this chapter are listed below:

- *Symmetry (or reflection) with respect to F and parallel to G .*
- *Orthogonal symmetry (or reflection) with respect to F and parallel to G ; reflections, flips.*
- *Hyperplane reflections and Householder matrices.*
- *A key fact about reflections (Proposition 11.1).*
- *QR -decomposition in terms of Householder transformations (Theorem 11.3).*

Chapter 12

Hermitian Spaces

12.1 Sesquilinear and Hermitian Forms, Pre-Hilbert Spaces and Hermitian Spaces

In this chapter we generalize the basic results of Euclidean geometry presented in Chapter 10 to vector spaces over the complex numbers. Such a generalization is inevitable, and not simply a luxury. For example, linear maps may not have real eigenvalues, but they always have complex eigenvalues. Furthermore, some very important classes of linear maps can be diagonalized if they are extended to the complexification of a real vector space. This is the case for orthogonal matrices, and, more generally, normal matrices. Also, complex vector spaces are often the natural framework in physics or engineering, and they are more convenient for dealing with Fourier series. However, some complications arise due to complex conjugation.

Recall that for any complex number $z \in \mathbb{C}$, if $z = x + iy$ where $x, y \in \mathbb{R}$, we let $\Re z = x$, the real part of z , and $\Im z = y$, the imaginary part of z . We also denote the conjugate of $z = x + iy$ by $\bar{z} = x - iy$, and the absolute value (or length, or modulus) of z by $|z|$. Recall that $|z|^2 = z\bar{z} = x^2 + y^2$.

There are many natural situations where a map $\varphi: E \times E \rightarrow \mathbb{C}$ is linear in its first argument and only semilinear in its second argument, which means that $\varphi(u, \mu v) = \bar{\mu}\varphi(u, v)$, as opposed to $\varphi(u, \mu v) = \mu\varphi(u, v)$. For example, the natural inner product to deal with functions $f: \mathbb{R} \rightarrow \mathbb{C}$, especially Fourier series, is

$$\langle f, g \rangle = \int_{-\pi}^{\pi} f(x) \overline{g(x)} dx,$$

which is semilinear (but not linear) in g . Thus, when generalizing a result from the real case of a Euclidean space to the complex case, we always have to check very carefully that our proofs do not rely on linearity in the second argument. Otherwise, we need to revise our proofs, and sometimes the result is simply wrong!

Before defining the natural generalization of an inner product, it is convenient to define semilinear maps.

Definition 12.1. Given two vector spaces E and F over the complex field \mathbb{C} , a function $f: E \rightarrow F$ is *semilinear* if

$$\begin{aligned} f(u + v) &= f(u) + f(v), \\ f(\lambda u) &= \bar{\lambda}f(u), \end{aligned}$$

for all $u, v \in E$ and all $\lambda \in \mathbb{C}$.

Remark: Instead of defining semilinear maps, we could have defined the vector space \bar{E} as the vector space with the same carrier set E whose addition is the same as that of E , but whose multiplication by a complex number is given by

$$(\lambda, u) \mapsto \bar{\lambda}u.$$

Then it is easy to check that a function $f: E \rightarrow \mathbb{C}$ is semilinear iff $f: \bar{E} \rightarrow \mathbb{C}$ is linear.

We can now define sesquilinear forms and Hermitian forms.

Definition 12.2. Given a complex vector space E , a function $\varphi: E \times E \rightarrow \mathbb{C}$ is a *sesquilinear form* if it is linear in its first argument and semilinear in its second argument, which means that

$$\begin{aligned} \varphi(u_1 + u_2, v) &= \varphi(u_1, v) + \varphi(u_2, v), \\ \varphi(u, v_1 + v_2) &= \varphi(u, v_1) + \varphi(u, v_2), \\ \varphi(\lambda u, v) &= \lambda\varphi(u, v), \\ \varphi(u, \mu v) &= \bar{\mu}\varphi(u, v), \end{aligned}$$

for all $u, v, u_1, u_2, v_1, v_2 \in E$, and all $\lambda, \mu \in \mathbb{C}$. A function $\varphi: E \times E \rightarrow \mathbb{C}$ is a *Hermitian form* if it is sesquilinear and if

$$\varphi(v, u) = \overline{\varphi(u, v)}$$

for all $u, v \in E$.

Obviously, $\varphi(0, v) = \varphi(u, 0) = 0$. Also note that if $\varphi: E \times E \rightarrow \mathbb{C}$ is sesquilinear, we have

$$\varphi(\lambda u + \mu v, \lambda u + \mu v) = |\lambda|^2\varphi(u, u) + \lambda\bar{\mu}\varphi(u, v) + \bar{\lambda}\mu\varphi(v, u) + |\mu|^2\varphi(v, v),$$

and if $\varphi: E \times E \rightarrow \mathbb{C}$ is Hermitian, we have

$$\varphi(\lambda u + \mu v, \lambda u + \mu v) = |\lambda|^2\varphi(u, u) + 2\Re(\lambda\bar{\mu}\varphi(u, v)) + |\mu|^2\varphi(v, v).$$

Note that restricted to real coefficients, a sesquilinear form is bilinear (we sometimes say \mathbb{R} -bilinear). The function $\Phi: E \rightarrow \mathbb{C}$ defined such that $\Phi(u) = \varphi(u, u)$ for all $u \in E$ is called the *quadratic form* associated with φ .

The standard example of a Hermitian form on \mathbb{C}^n is the map φ defined such that

$$\varphi((x_1, \dots, x_n), (y_1, \dots, y_n)) = x_1 \overline{y_1} + x_2 \overline{y_2} + \dots + x_n \overline{y_n}.$$

This map is also positive definite, but before dealing with these issues, we show the following useful proposition.

Proposition 12.1. *Given a complex vector space E , the following properties hold:*

- (1) *A sesquilinear form $\varphi: E \times E \rightarrow \mathbb{C}$ is a Hermitian form iff $\varphi(u, u) \in \mathbb{R}$ for all $u \in E$.*
- (2) *If $\varphi: E \times E \rightarrow \mathbb{C}$ is a sesquilinear form, then*

$$\begin{aligned} 4\varphi(u, v) &= \varphi(u + v, u + v) - \varphi(u - v, u - v) \\ &\quad + i\varphi(u + iv, u + iv) - i\varphi(u - iv, u - iv), \end{aligned}$$

and

$$2\varphi(u, v) = (1 + i)(\varphi(u, u) + \varphi(v, v)) - \varphi(u - v, u - v) - i\varphi(u - iv, u - iv).$$

These are called polarization identities.

Proof. (1) If φ is a Hermitian form, then

$$\varphi(v, u) = \overline{\varphi(u, v)}$$

implies that

$$\varphi(u, u) = \overline{\varphi(u, u)},$$

and thus $\varphi(u, u) \in \mathbb{R}$. If φ is sesquilinear and $\varphi(u, u) \in \mathbb{R}$ for all $u \in E$, then

$$\varphi(u + v, u + v) = \varphi(u, u) + \varphi(u, v) + \varphi(v, u) + \varphi(v, v),$$

which proves that

$$\varphi(u, v) + \varphi(v, u) = \alpha,$$

where α is real, and changing u to iu , we have

$$i(\varphi(u, v) - \varphi(v, u)) = \beta,$$

where β is real, and thus

$$\varphi(u, v) = \frac{\alpha - i\beta}{2} \quad \text{and} \quad \varphi(v, u) = \frac{\alpha + i\beta}{2},$$

proving that φ is Hermitian.

(2) These identities are verified by expanding the right-hand side, and we leave them as an exercise. \square

Proposition 12.1 shows that a sesquilinear form is completely determined by the quadratic form $\Phi(u) = \varphi(u, u)$, even if φ is not Hermitian. This is false for a real bilinear form, unless it is symmetric. For example, the bilinear form $\varphi: \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ defined such that

$$\varphi((x_1, y_1), (x_2, y_2)) = x_1 y_2 - x_2 y_1$$

is not identically zero, and yet it is null on the diagonal. However, a real symmetric bilinear form is indeed determined by its values on the diagonal, as we saw in Chapter 10.

As in the Euclidean case, Hermitian forms for which $\varphi(u, u) \geq 0$ play an important role.

Definition 12.3. Given a complex vector space E , a Hermitian form $\varphi: E \times E \rightarrow \mathbb{C}$ is *positive* if $\varphi(u, u) \geq 0$ for all $u \in E$, and *positive definite* if $\varphi(u, u) > 0$ for all $u \neq 0$. A pair $\langle E, \varphi \rangle$ where E is a complex vector space and φ is a Hermitian form on E is called a *pre-Hilbert space* if φ is positive, and a *Hermitian (or unitary) space* if φ is positive definite.

We warn our readers that some authors, such as Lang [69], define a pre-Hilbert space as what we define as a Hermitian space. We prefer following the terminology used in Schwartz [93] and Bourbaki [16]. The quantity $\varphi(u, v)$ is usually called the *Hermitian product* of u and v . We will occasionally call it the inner product of u and v .

Given a pre-Hilbert space $\langle E, \varphi \rangle$, as in the case of a Euclidean space, we also denote $\varphi(u, v)$ by

$$u \cdot v \quad \text{or} \quad \langle u, v \rangle \quad \text{or} \quad (u|v),$$

and $\sqrt{\Phi(u)}$ by $\|u\|$.

Example 12.1. The complex vector space \mathbb{C}^n under the Hermitian form

$$\varphi((x_1, \dots, x_n), (y_1, \dots, y_n)) = x_1 \overline{y_1} + x_2 \overline{y_2} + \dots + x_n \overline{y_n}$$

is a Hermitian space.

Example 12.2. Let l^2 denote the set of all countably infinite sequences $x = (x_i)_{i \in \mathbb{N}}$ of complex numbers such that $\sum_{i=0}^{\infty} |x_i|^2$ is defined (i.e., the sequence $\sum_{i=0}^n |x_i|^2$ converges as $n \rightarrow \infty$). It can be shown that the map $\varphi: l^2 \times l^2 \rightarrow \mathbb{C}$ defined such that

$$\varphi((x_i)_{i \in \mathbb{N}}, (y_i)_{i \in \mathbb{N}}) = \sum_{i=0}^{\infty} x_i \overline{y_i}$$

is well defined, and l^2 is a Hermitian space under φ . Actually, l^2 is even a Hilbert space.

Example 12.3. Let $\mathcal{C}_{\text{piece}}[a, b]$ be the set of piecewise bounded continuous functions $f: [a, b] \rightarrow \mathbb{C}$ under the Hermitian form

$$\langle f, g \rangle = \int_a^b f(x) \overline{g(x)} dx.$$

It is easy to check that this Hermitian form is positive, but it is not definite. Thus, under this Hermitian form, $\mathcal{C}_{\text{piece}}[a, b]$ is only a pre-Hilbert space.

Example 12.4. Let $\mathcal{C}[a, b]$ be the set of complex-valued continuous functions $f: [a, b] \rightarrow \mathbb{C}$ under the Hermitian form

$$\langle f, g \rangle = \int_a^b f(x) \overline{g(x)} dx.$$

It is easy to check that this Hermitian form is positive definite. Thus, $\mathcal{C}[a, b]$ is a Hermitian space.

Example 12.5. Let $E = M_n(\mathbb{C})$ be the vector space of complex $n \times n$ matrices. If we view a matrix $A \in M_n(\mathbb{C})$ as a “long” column vector obtained by concatenating together its columns, we can define the Hermitian product of two matrices $A, B \in M_n(\mathbb{C})$ as

$$\langle A, B \rangle = \sum_{i,j=1}^n a_{ij} \bar{b}_{ij},$$

which can be conveniently written as

$$\langle A, B \rangle = \operatorname{tr}(A^\top \bar{B}) = \operatorname{tr}(B^* A).$$

Since this can be viewed as the standard Hermitian product on \mathbb{C}^{n^2} , it is a Hermitian product on $M_n(\mathbb{C})$. The corresponding norm

$$\|A\|_F = \sqrt{\operatorname{tr}(A^* A)}$$

is the Frobenius norm (see Section 7.2).

If E is finite-dimensional and if $\varphi: E \times E \rightarrow \mathbb{R}$ is a sequilinear form on E , given any basis (e_1, \dots, e_n) of E , we can write $x = \sum_{i=1}^n x_i e_i$ and $y = \sum_{j=1}^n y_j e_j$, and we have

$$\varphi(x, y) = \varphi\left(\sum_{i=1}^n x_i e_i, \sum_{j=1}^n y_j e_j\right) = \sum_{i,j=1}^n x_i \bar{y}_j \varphi(e_i, e_j).$$

If we let G be the matrix $G = (\varphi(e_i, e_j))$, and if x and y are the column vectors associated with (x_1, \dots, x_n) and (y_1, \dots, y_n) , then we can write

$$\varphi(x, y) = x^\top G \bar{y} = y^* G^\top x,$$

where \bar{y} corresponds to $(\bar{y}_1, \dots, \bar{y}_n)$. As in Section 10.1, we are committing the slight abuse of notation of letting x denote both the vector $x = \sum_{i=1}^n x_i e_i$ and the column vector associated with (x_1, \dots, x_n) (and similarly for y). The “correct” expression for $\varphi(x, y)$ is

$$\varphi(x, y) = \mathbf{y}^* G^\top \mathbf{x} = \mathbf{x}^\top G \bar{\mathbf{y}}.$$



Observe that in $\varphi(x, y) = y^* G^\top x$, the matrix involved is the transpose of $G = (\varphi(e_i, e_j))$.

Furthermore, observe that φ is Hermitian iff $G = G^*$, and φ is positive definite iff the matrix G is positive definite, that is,

$$x^\top Gx > 0 \quad \text{for all } x \in \mathbb{C}^n, x \neq 0.$$

The matrix G associated with a Hermitian product is called the *Gram matrix* of the Hermitian product with respect to the basis (e_1, \dots, e_n) .

Remark: To avoid the transposition in the expression for $\varphi(x, y) = y^* G^\top x$, some authors (such as Hoffman and Kunze [62]), define the Gram matrix $G' = (g'_{ij})$ associated with $\langle -, - \rangle$ so that $(g'_{ij}) = (\varphi(e_j, e_i))$; that is, $G' = G^\top$.

Conversely, if A is a Hermitian positive definite $n \times n$ matrix, it is easy to check that the Hermitian form

$$\langle x, y \rangle = y^* Ax$$

is positive definite. If we make a change of basis from the basis (e_1, \dots, e_n) to the basis (f_1, \dots, f_n) , and if the change of basis matrix is P (where the j th column of P consists of the coordinates of f_j over the basis (e_1, \dots, e_n)), then with respect to coordinates x' and y' over the basis (f_1, \dots, f_n) , we have

$$x^\top G \bar{y} = x'^\top P^\top G \bar{P} \bar{y}',$$

so the matrix of our inner product over the basis (f_1, \dots, f_n) is $P^\top G \bar{P} = (\bar{P})^* G \bar{P}$. We summarize these facts in the following proposition.

Proposition 12.2. *Let E be a finite-dimensional vector space, and let (e_1, \dots, e_n) be a basis of E .*

1. *For any Hermitian inner product $\langle -, - \rangle$ on E , if $G = (\langle e_i, e_j \rangle)$ is the Gram matrix of the Hermitian product $\langle -, - \rangle$ w.r.t. the basis (e_1, \dots, e_n) , then G is Hermitian positive definite.*
2. *For any change of basis matrix P , the Gram matrix of $\langle -, - \rangle$ with respect to the new basis is $(\bar{P})^* G \bar{P}$.*
3. *If A is any $n \times n$ Hermitian positive definite matrix, then*

$$\langle x, y \rangle = y^* Ax$$

is a Hermitian product on E .

We will see later that a Hermitian matrix is positive definite iff its eigenvalues are all positive.

The following result reminiscent of the first polarization identity of Proposition 12.1 can be used to prove that two linear maps are identical.

Proposition 12.3. *Given any Hermitian space E with Hermitian product $\langle -, - \rangle$, for any linear map $f: E \rightarrow E$, if $\langle f(x), x \rangle = 0$ for all $x \in E$, then $f = 0$.*

Proof. Compute $\langle f(x+y), x+y \rangle$ and $\langle f(x-y), x-y \rangle$:

$$\begin{aligned}\langle f(x+y), x+y \rangle &= \langle f(x), x \rangle + \langle f(x), y \rangle + \langle f(y), x \rangle + \langle y, y \rangle \\ \langle f(x-y), x-y \rangle &= \langle f(x), x \rangle - \langle f(x), y \rangle - \langle f(y), x \rangle + \langle y, y \rangle;\end{aligned}$$

then, subtract the second equation from the first, to obtain

$$\langle f(x+y), x+y \rangle - \langle f(x-y), x-y \rangle = 2(\langle f(x), y \rangle + \langle f(y), x \rangle).$$

If $\langle f(u), u \rangle = 0$ for all $u \in E$, we get

$$\langle f(x), y \rangle + \langle f(y), x \rangle = 0 \quad \text{for all } x, y \in E.$$

Then, the above equation also holds if we replace x by ix , and we obtain

$$i\langle f(x), y \rangle - i\langle f(y), x \rangle = 0, \quad \text{for all } x, y \in E,$$

so we have

$$\begin{aligned}\langle f(x), y \rangle + \langle f(y), x \rangle &= 0 \\ \langle f(x), y \rangle - \langle f(y), x \rangle &= 0,\end{aligned}$$

which implies that $\langle f(x), y \rangle = 0$ for all $x, y \in E$. Since $\langle -, - \rangle$ is positive definite, we have $f(x) = 0$ for all $x \in E$; that is, $f = 0$. \square

One should be careful not to apply Proposition 12.3 to a linear map on a real Euclidean space, because it is false! The reader should find a counterexample.

The Cauchy–Schwarz inequality and the Minkowski inequalities extend to pre-Hilbert spaces and to Hermitian spaces.

Proposition 12.4. *Let $\langle E, \varphi \rangle$ be a pre-Hilbert space with associated quadratic form Φ . For all $u, v \in E$, we have the Cauchy–Schwarz inequality*

$$|\varphi(u, v)| \leq \sqrt{\Phi(u)}\sqrt{\Phi(v)}.$$

Furthermore, if $\langle E, \varphi \rangle$ is a Hermitian space, the equality holds iff u and v are linearly dependent.

We also have the Minkowski inequality

$$\sqrt{\Phi(u+v)} \leq \sqrt{\Phi(u)} + \sqrt{\Phi(v)}.$$

Furthermore, if $\langle E, \varphi \rangle$ is a Hermitian space, the equality holds iff u and v are linearly dependent, where in addition, if $u \neq 0$ and $v \neq 0$, then $u = \lambda v$ for some real λ such that $\lambda > 0$.

Proof. For all $u, v \in E$ and all $\mu \in \mathbb{C}$, we have observed that

$$\varphi(u + \mu v, u + \mu v) = \varphi(u, u) + 2\Re(\bar{\mu}\varphi(u, v)) + |\mu|^2\varphi(v, v).$$

Let $\varphi(u, v) = \rho e^{i\theta}$, where $|\varphi(u, v)| = \rho$ ($\rho \geq 0$). Let $F: \mathbb{R} \rightarrow \mathbb{R}$ be the function defined such that

$$F(t) = \Phi(u + te^{i\theta}v),$$

for all $t \in \mathbb{R}$. The above shows that

$$F(t) = \varphi(u, u) + 2t|\varphi(u, v)| + t^2\varphi(v, v) = \Phi(u) + 2t|\varphi(u, v)| + t^2\Phi(v).$$

Since φ is assumed to be positive, we have $F(t) \geq 0$ for all $t \in \mathbb{R}$. If $\Phi(v) = 0$, we must have $\varphi(u, v) = 0$, since otherwise, $F(t)$ could be made negative by choosing t negative and small enough. If $\Phi(v) > 0$, in order for $F(t)$ to be nonnegative, the equation

$$\Phi(u) + 2t|\varphi(u, v)| + t^2\Phi(v) = 0$$

must not have distinct real roots, which is equivalent to

$$|\varphi(u, v)|^2 \leq \Phi(u)\Phi(v).$$

Taking the square root on both sides yields the Cauchy–Schwarz inequality.

For the second part of the claim, if φ is positive definite, we argue as follows. If u and v are linearly dependent, it is immediately verified that we get an equality. Conversely, if

$$|\varphi(u, v)|^2 = \Phi(u)\Phi(v),$$

then there are two cases. If $\Phi(v) = 0$, since φ is positive definite, we must have $v = 0$, so u and v are linearly dependent. Otherwise, the equation

$$\Phi(u) + 2t|\varphi(u, v)| + t^2\Phi(v) = 0$$

has a double root t_0 , and thus

$$\Phi(u + t_0 e^{i\theta}v) = 0.$$

Since φ is positive definite, we must have

$$u + t_0 e^{i\theta}v = 0,$$

which shows that u and v are linearly dependent.

If we square the Minkowski inequality, we get

$$\Phi(u + v) \leq \Phi(u) + \Phi(v) + 2\sqrt{\Phi(u)}\sqrt{\Phi(v)}.$$

However, we observed earlier that

$$\Phi(u + v) = \Phi(u) + \Phi(v) + 2\Re(\varphi(u, v)).$$

Thus, it is enough to prove that

$$\Re(\varphi(u, v)) \leq \sqrt{\Phi(u)}\sqrt{\Phi(v)},$$

but this follows from the Cauchy–Schwarz inequality

$$|\varphi(u, v)| \leq \sqrt{\Phi(u)}\sqrt{\Phi(v)}$$

and the fact that $\Re z \leq |z|$.

If φ is positive definite and u and v are linearly dependent, it is immediately verified that we get an equality. Conversely, if equality holds in the Minkowski inequality, we must have

$$\Re(\varphi(u, v)) = \sqrt{\Phi(u)}\sqrt{\Phi(v)},$$

which implies that

$$|\varphi(u, v)| = \sqrt{\Phi(u)}\sqrt{\Phi(v)},$$

since otherwise, by the Cauchy–Schwarz inequality, we would have

$$\Re(\varphi(u, v)) \leq |\varphi(u, v)| < \sqrt{\Phi(u)}\sqrt{\Phi(v)}.$$

Thus, equality holds in the Cauchy–Schwarz inequality, and

$$\Re(\varphi(u, v)) = |\varphi(u, v)|.$$

But then, we proved in the Cauchy–Schwarz case that u and v are linearly dependent. Since we also just proved that $\varphi(u, v)$ is real and nonnegative, the coefficient of proportionality between u and v is indeed nonnegative. \square

As in the Euclidean case, if $\langle E, \varphi \rangle$ is a Hermitian space, the Minkowski inequality

$$\sqrt{\Phi(u+v)} \leq \sqrt{\Phi(u)} + \sqrt{\Phi(v)}$$

shows that the map $u \mapsto \sqrt{\Phi(u)}$ is a *norm* on E . The norm induced by φ is called the *Hermitian norm induced by φ* . We usually denote $\sqrt{\Phi(u)}$ by $\|u\|$, and the Cauchy–Schwarz inequality is written as

$$|u \cdot v| \leq \|u\|\|v\|.$$

Since a Hermitian space is a normed vector space, it is a topological space under the topology induced by the norm (a basis for this topology is given by the open balls $B_0(u, \rho)$ of center u and radius $\rho > 0$, where

$$B_0(u, \rho) = \{v \in E \mid \|v - u\| < \rho\}.$$

If E has finite dimension, every linear map is continuous; see Chapter 7 (or Lang [69, 70], Dixmier [29], or Schwartz [93, 94]). The Cauchy–Schwarz inequality

$$|u \cdot v| \leq \|u\|\|v\|$$

shows that $\varphi: E \times E \rightarrow \mathbb{C}$ is continuous, and thus, that $\|\cdot\|$ is continuous.

If $\langle E, \varphi \rangle$ is only pre-Hilbertian, $\|u\|$ is called a *seminorm*. In this case, the condition

$$\|u\| = 0 \quad \text{implies} \quad u = 0$$

is not necessarily true. However, the Cauchy–Schwarz inequality shows that if $\|u\| = 0$, then $u \cdot v = 0$ for all $v \in E$.

Remark: As in the case of real vector spaces, a norm on a complex vector space is induced by some positive definite Hermitian product $\langle -, - \rangle$ iff it satisfies the *parallelogram law*:

$$\|u + v\|^2 + \|u - v\|^2 = 2(\|u\|^2 + \|v\|^2).$$

This time, the Hermitian product is recovered using the polarization identity from Proposition 12.1:

$$4\langle u, v \rangle = \|u + v\|^2 - \|u - v\|^2 + i\|u + iv\|^2 - i\|u - iv\|^2.$$

It is easy to check that $\langle u, u \rangle = \|u\|^2$, and

$$\begin{aligned} \langle v, u \rangle &= \overline{\langle u, v \rangle} \\ \langle iu, v \rangle &= i\langle u, v \rangle, \end{aligned}$$

so it is enough to check linearity in the variable u , and only for real scalars. This is easily done by applying the proof from Section 10.1 to the real and imaginary part of $\langle u, v \rangle$; the details are left as an exercise.

We will now basically mirror the presentation of Euclidean geometry given in Chapter 10 rather quickly, leaving out most proofs, except when they need to be seriously amended.

12.2 Orthogonality, Duality, Adjoint of a Linear Map

In this section we assume that we are dealing with Hermitian spaces. We denote the Hermitian inner product by $u \cdot v$ or $\langle u, v \rangle$. The concepts of orthogonality, orthogonal family of vectors, orthonormal family of vectors, and orthogonal complement of a set of vectors are unchanged from the Euclidean case (Definition 10.2).

For example, the set $\mathcal{C}[-\pi, \pi]$ of continuous functions $f: [-\pi, \pi] \rightarrow \mathbb{C}$ is a Hermitian space under the product

$$\langle f, g \rangle = \int_{-\pi}^{\pi} f(x) \overline{g(x)} dx,$$

and the family $(e^{ikx})_{k \in \mathbb{Z}}$ is orthogonal.

Proposition 10.3 and 10.4 hold without any changes. It is easy to show that

$$\left\| \sum_{i=1}^n u_i \right\|^2 = \sum_{i=1}^n \|u_i\|^2 + \sum_{1 \leq i < j \leq n} 2\Re(u_i \cdot u_j).$$

Analogously to the case of Euclidean spaces of finite dimension, the Hermitian product induces a canonical bijection (i.e., independent of the choice of bases) between the vector space E and the space E^* . This is one of the places where conjugation shows up, but in this case, troubles are minor.

Given a Hermitian space E , for any vector $u \in E$, let $\varphi_u^l: E \rightarrow \mathbb{C}$ be the map defined such that

$$\varphi_u^l(v) = \overline{u \cdot v}, \quad \text{for all } v \in E.$$

Similarly, for any vector $v \in E$, let $\varphi_v^r: E \rightarrow \mathbb{C}$ be the map defined such that

$$\varphi_v^r(u) = u \cdot v, \quad \text{for all } u \in E.$$

Since the Hermitian product is linear in its first argument u , the map φ_v^r is a linear form in E^* , and since it is semilinear in its second argument v , the map φ_u^l is also a linear form in E^* . Thus, we have two maps $\flat^l: E \rightarrow E^*$ and $\flat^r: E \rightarrow E^*$, defined such that

$$\flat^l(u) = \varphi_u^l, \quad \text{and} \quad \flat^r(v) = \varphi_v^r.$$

Actually, $\varphi_u^l = \varphi_u^r$ and $\flat^l = \flat^r$. Indeed, for all $u, v \in E$, we have

$$\begin{aligned} \flat^l(u)(v) &= \varphi_u^l(v) \\ &= \overline{u \cdot v} \\ &= v \cdot u \\ &= \varphi_u^r(v) \\ &= \flat^r(u)(v). \end{aligned}$$

Therefore, we use the notation φ_u for both φ_u^l and φ_u^r , and \flat for both \flat^l and \flat^r .

Theorem 12.5. *let E be a Hermitian space E . The map $\flat: E \rightarrow E^*$ defined such that*

$$\flat(u) = \varphi_u^l = \varphi_u^r \quad \text{for all } u \in E$$

is semilinear and injective. When E is also of finite dimension, the map $\flat: \overline{E} \rightarrow E^$ is a canonical isomorphism.*

Proof. That $\flat: E \rightarrow E^*$ is a semilinear map follows immediately from the fact that $\flat = \flat^r$, and that the Hermitian product is semilinear in its second argument. If $\varphi_u = \varphi_v$, then $\varphi_u(w) = \varphi_v(w)$ for all $w \in E$, which by definition of φ_u and φ_v means that

$$w \cdot u = w \cdot v$$

for all $w \in E$, which by semilinearity on the right is equivalent to

$$w \cdot (v - u) = 0 \quad \text{for all } w \in E,$$

which implies that $u = v$, since the Hermitian product is positive definite. Thus, $\flat: E \rightarrow E^*$ is injective. Finally, when E is of finite dimension n , E^* is also of dimension n , and then $\flat: E \rightarrow E^*$ is bijective. Since \flat is semilinear, the map $\flat: \overline{E} \rightarrow E^*$ is an isomorphism. \square

The inverse of the isomorphism $\flat: \overline{E} \rightarrow E^*$ is denoted by $\sharp: E^* \rightarrow \overline{E}$.

As a corollary of the isomorphism $\flat: \overline{E} \rightarrow E^*$, if E is a Hermitian space of finite dimension, then every linear form $f \in E^*$ corresponds to a unique $v \in E$, such that

$$f(u) = u \cdot v, \quad \text{for every } u \in E.$$

In particular, if f is not the null form, the kernel of f , which is a hyperplane H , is precisely the set of vectors that are orthogonal to v .

Remarks:

1. The “musical map” $\flat: \overline{E} \rightarrow E^*$ is not surjective when E has infinite dimension. This result can be salvaged by restricting our attention to continuous linear maps, and by assuming that the vector space E is a *Hilbert space*.
2. *Dirac’s “bra-ket” notation.* Dirac invented a notation widely used in quantum mechanics for denoting the linear form $\varphi_u = \flat(u)$ associated to the vector $u \in E$ via the duality induced by a Hermitian inner product. Dirac’s proposal is to denote the vectors u in E by $|u\rangle$, and call them *kets*; the notation $|u\rangle$ is pronounced “ket u .” Given two kets (vectors) $|u\rangle$ and $|v\rangle$, their inner product is denoted by

$$\langle u|v\rangle$$

(instead of $|u\rangle \cdot |v\rangle$). The notation $\langle u|v\rangle$ for the inner product of $|u\rangle$ and $|v\rangle$ anticipates duality. Indeed, we define the dual (usually called adjoint) *bra* u of ket u , denoted by $\langle u|$, as the linear form whose value on any ket v is given by the inner product, so

$$\langle u|(|v\rangle) = \langle u|v\rangle.$$

Thus, bra $u = \langle u|$ is Dirac’s notation for our $\flat(u)$. Since the map \flat is semi-linear, we have

$$\langle \lambda u| = \overline{\lambda} \langle u|.$$

Using the bra-ket notation, given an orthonormal basis $(|u_1\rangle, \dots, |u_n\rangle)$, ket v (a vector) is written as

$$|v\rangle = \sum_{i=1}^n \langle v|u_i\rangle |u_i\rangle,$$

and the corresponding linear form bra v is written as

$$\langle v| = \sum_{i=1}^n \overline{\langle v|u_i\rangle} \langle u_i| = \sum_{i=1}^n \langle u_i|v\rangle \langle u_i|$$

over the dual basis $(\langle u_1|, \dots, \langle u_n|)$. As cute as it looks, we do not recommend using the Dirac notation.

The existence of the isomorphism $\flat: \overline{E} \rightarrow E^*$ is crucial to the existence of adjoint maps. Indeed, Theorem 12.5 allows us to define the adjoint of a linear map on a Hermitian space. Let E be a Hermitian space of finite dimension n , and let $f: E \rightarrow E$ be a linear map. For every $u \in E$, the map

$$v \mapsto \overline{u \cdot f(v)}$$

is clearly a linear form in E^* , and by Theorem 12.5, there is a unique vector in E denoted by $f^*(u)$, such that

$$\overline{f^*(u) \cdot v} = \overline{u \cdot f(v)},$$

that is,

$$f^*(u) \cdot v = u \cdot f(v), \quad \text{for every } v \in E.$$

The following proposition shows that the map f^* is linear.

Proposition 12.6. *Given a Hermitian space E of finite dimension, for every linear map $f: E \rightarrow E$ there is a unique linear map $f^*: E \rightarrow E$ such that*

$$f^*(u) \cdot v = u \cdot f(v),$$

for all $u, v \in E$. The map f^ is called the adjoint of f (w.r.t. to the Hermitian product).*

Proof. Careful inspection of the proof of Proposition 10.6 reveals that it applies unchanged. The only potential problem is in proving that $f^*(\lambda u) = \lambda f^*(u)$, but everything takes place in the first argument of the Hermitian product, and there, we have linearity. \square

The fact that

$$v \cdot u = \overline{u \cdot v}$$

implies that the adjoint f^* of f is also characterized by

$$f(u) \cdot v = u \cdot f^*(v),$$

for all $u, v \in E$. It is also obvious that $f^{**} = f$.

Given two Hermitian spaces E and F , where the Hermitian product on E is denoted by $\langle -, - \rangle_1$ and the Hermitian product on F is denoted by $\langle -, - \rangle_2$, given any linear map $f: E \rightarrow F$, it is immediately verified that the proof of Proposition 12.6 can be adapted to show that there is a unique linear map $f^*: F \rightarrow E$ such that

$$\langle f(u), v \rangle_2 = \langle u, f^*(v) \rangle_1$$

for all $u \in E$ and all $v \in F$. The linear map f^* is also called the *adjoint* of f .

As in the Euclidean case, a linear map $f: E \rightarrow E$ (where E is a finite-dimensional Hermitian space) is *self-adjoint* if $f = f^*$. The map f is *positive semidefinite* iff

$$\langle f(x), x \rangle \geq 0 \quad \text{all } x \in E;$$

positive definite iff

$$\langle f(x), x \rangle > 0 \quad \text{all } x \in E, x \neq 0.$$

An interesting corollary of Proposition 12.3 is that a positive semidefinite linear map must be self-adjoint. In fact, we can prove a slightly more general result.

Proposition 12.7. *Given any finite-dimensional Hermitian space E with Hermitian product $\langle -, - \rangle$, for any linear map $f: E \rightarrow E$, if $\langle f(x), x \rangle \in \mathbb{R}$ for all $x \in E$, then f is self-adjoint. In particular, any positive semidefinite linear map $f: E \rightarrow E$ is self-adjoint.*

Proof. Since $\langle f(x), x \rangle \in \mathbb{R}$ for all $x \in E$, we have

$$\begin{aligned} \langle f(x), x \rangle &= \overline{\langle f(x), x \rangle} \\ &= \langle x, f(x) \rangle \\ &= \langle f^*(x), x \rangle, \end{aligned}$$

so we have

$$\langle (f - f^*)(x), x \rangle = 0 \quad \text{all } x \in E,$$

and Proposition 12.3 implies that $f - f^* = 0$. □

Beware that Proposition 12.7 is false if E is a real Euclidean space.

As in the Euclidean case, Theorem 12.5 can be used to show that any Hermitian space of finite dimension has an orthonormal basis. The proof is unchanged.

Proposition 12.8. *Given any nontrivial Hermitian space E of finite dimension $n \geq 1$, there is an orthonormal basis (u_1, \dots, u_n) for E .*

The *Gram-Schmidt orthonormalization procedure* also applies to Hermitian spaces of finite dimension, without any changes from the Euclidean case!

Proposition 12.9. *Given a nontrivial Hermitian space E of finite dimension $n \geq 1$, from any basis (e_1, \dots, e_n) for E we can construct an orthonormal basis (u_1, \dots, u_n) for E with the property that for every k , $1 \leq k \leq n$, the families (e_1, \dots, e_k) and (u_1, \dots, u_k) generate the same subspace.*

Remark: The remarks made after Proposition 10.8 also apply here, except that in the QR -decomposition, Q is a unitary matrix.

As a consequence of Proposition 10.7 (or Proposition 12.9), given any Hermitian space of finite dimension n , if (e_1, \dots, e_n) is an orthonormal basis for E , then for any two vectors $u = u_1 e_1 + \dots + u_n e_n$ and $v = v_1 e_1 + \dots + v_n e_n$, the Hermitian product $u \cdot v$ is expressed as

$$u \cdot v = (u_1 e_1 + \dots + u_n e_n) \cdot (v_1 e_1 + \dots + v_n e_n) = \sum_{i=1}^n u_i \overline{v_i},$$

and the norm $\|u\|$ as

$$\|u\| = \|u_1e_1 + \cdots + u_ne_n\| = \left(\sum_{i=1}^n |u_i|^2 \right)^{1/2}.$$

The fact that a Hermitian space always has an orthonormal basis implies that any Gram matrix G can be written as

$$G = Q^*Q,$$

for some invertible matrix Q . Indeed, we know that in a change of basis matrix, a Gram matrix G becomes $G' = (\bar{P})^*G\bar{P}$. If the basis corresponding to G' is orthonormal, then $G' = I$, so $G = (\bar{P}^{-1})^*P^{-1}$.

Proposition 10.9 also holds unchanged.

Proposition 12.10. *Given any nontrivial Hermitian space E of finite dimension $n \geq 1$, for any subspace F of dimension k , the orthogonal complement F^\perp of F has dimension $n - k$, and $E = F \oplus F^\perp$. Furthermore, we have $F^{\perp\perp} = F$.*

12.3 Linear Isometries (Also Called Unitary Transformations)

In this section we consider linear maps between Hermitian spaces that preserve the Hermitian norm. All definitions given for Euclidean spaces in Section 10.3 extend to Hermitian spaces, except that orthogonal transformations are called unitary transformation, but Proposition 10.10 extends only with a modified condition (2). Indeed, the old proof that (2) implies (3) does not work, and the implication is in fact false! It can be repaired by strengthening condition (2). For the sake of completeness, we state the Hermitian version of Definition 10.3.

Definition 12.4. Given any two nontrivial Hermitian spaces E and F of the same finite dimension n , a function $f: E \rightarrow F$ is a *unitary transformation*, or a *linear isometry*, if it is linear and

$$\|f(u)\| = \|u\|, \quad \text{for all } u \in E.$$

Proposition 10.10 can be salvaged by strengthening condition (2).

Proposition 12.11. *Given any two nontrivial Hermitian spaces E and F of the same finite dimension n , for every function $f: E \rightarrow F$, the following properties are equivalent:*

- (1) f is a linear map and $\|f(u)\| = \|u\|$, for all $u \in E$;
- (2) $\|f(v) - f(u)\| = \|v - u\|$ and $f(iu) = if(u)$, for all $u, v \in E$.

(3) $f(u) \cdot f(v) = u \cdot v$, for all $u, v \in E$.

Furthermore, such a map is bijective.

Proof. The proof that (2) implies (3) given in Proposition 10.10 needs to be revised as follows. We use the polarization identity

$$2\varphi(u, v) = (1 + i)(\|u\|^2 + \|v\|^2) - \|u - v\|^2 - i\|u - iv\|^2.$$

Since $f(iv) = if(v)$, we get $f(0) = 0$ by setting $v = 0$, so the function f preserves distance and norm, and we get

$$\begin{aligned} 2\varphi(f(u), f(v)) &= (1 + i)(\|f(u)\|^2 + \|f(v)\|^2) - \|f(u) - f(v)\|^2 \\ &\quad - i\|f(u) - if(v)\|^2 \\ &= (1 + i)(\|f(u)\|^2 + \|f(v)\|^2) - \|f(u) - f(v)\|^2 \\ &\quad - i\|f(u) - f(iv)\|^2 \\ &= (1 + i)(\|u\|^2 + \|v\|^2) - \|u - v\|^2 - i\|u - iv\|^2 \\ &= 2\varphi(u, v), \end{aligned}$$

which shows that f preserves the Hermitian inner product, as desired. The rest of the proof is unchanged. \square

Remarks:

(i) In the Euclidean case, we proved that the assumption

$$\|f(v) - f(u)\| = \|v - u\| \quad \text{for all } u, v \in E \text{ and } f(0) = 0 \quad (2')$$

implies (3). For this we used the polarization identity

$$2u \cdot v = \|u\|^2 + \|v\|^2 - \|u - v\|^2.$$

In the Hermitian case the polarization identity involves the complex number i . In fact, the implication (2') implies (3) is false in the Hermitian case! Conjugation $z \mapsto \bar{z}$ satisfies (2') since

$$|\bar{z}_2 - \bar{z}_1| = |\overline{z_2 - z_1}| = |z_2 - z_1|,$$

and yet, it is not linear!

(ii) If we modify (2) by changing the second condition by now requiring that there be some $\tau \in E$ such that

$$f(\tau + iu) = f(\tau) + i(f(\tau + u) - f(\tau))$$

for all $u \in E$, then the function $g: E \rightarrow E$ defined such that

$$g(u) = f(\tau + u) - f(\tau)$$

satisfies the old conditions of (2), and the implications (2) \rightarrow (3) and (3) \rightarrow (1) prove that g is linear, and thus that f is affine. In view of the first remark, some condition involving i is needed on f , in addition to the fact that f is distance-preserving.

12.4 The Unitary Group, Unitary Matrices

In this section, as a mirror image of our treatment of the isometries of a Euclidean space, we explore some of the fundamental properties of the unitary group and of unitary matrices. As an immediate corollary of the Gram–Schmidt orthonormalization procedure, we obtain the QR -decomposition for invertible matrices. In the Hermitian framework, the matrix of the adjoint of a linear map is not given by the transpose of the original matrix, but by its conjugate.

Definition 12.5. Given a complex $m \times n$ matrix A , the *transpose* A^\top of A is the $n \times m$ matrix $A^\top = (a_{ij}^\top)$ defined such that

$$a_{ij}^\top = a_{ji},$$

and the *conjugate* \bar{A} of A is the $m \times n$ matrix $\bar{A} = (\bar{a}_{ij})$ defined such that

$$b_{ij} = \bar{a}_{ij}$$

for all i, j , $1 \leq i \leq m$, $1 \leq j \leq n$. The *adjoint* A^* of A is the matrix defined such that

$$A^* = \overline{(A^\top)} = (\bar{A})^\top.$$

Proposition 12.12. *Let E be any Hermitian space of finite dimension n , and let $f: E \rightarrow E$ be any linear map. The following properties hold:*

(1) *The linear map $f: E \rightarrow E$ is an isometry iff*

$$f \circ f^* = f^* \circ f = \text{id}.$$

(2) *For every orthonormal basis (e_1, \dots, e_n) of E , if the matrix of f is A , then the matrix of f^* is the adjoint A^* of A , and f is an isometry iff A satisfies the identities*

$$A A^* = A^* A = I_n,$$

where I_n denotes the identity matrix of order n , iff the columns of A form an orthonormal basis of E , iff the rows of A form an orthonormal basis of E .

Proof. (1) The proof is identical to that of Proposition 10.12 (1).

(2) If (e_1, \dots, e_n) is an orthonormal basis for E , let $A = (a_{ij})$ be the matrix of f , and let $B = (b_{ij})$ be the matrix of f^* . Since f^* is characterized by

$$f^*(u) \cdot v = u \cdot f(v)$$

for all $u, v \in E$, using the fact that if $w = w_1 e_1 + \cdots + w_n e_n$, we have $w_k = w \cdot e_k$, for all k , $1 \leq k \leq n$; letting $u = e_i$ and $v = e_j$, we get

$$b_{ji} = f^*(e_i) \cdot e_j = e_i \cdot f(e_j) = \overline{f(e_j) \cdot e_i} = \overline{a_{ij}},$$

for all i, j , $1 \leq i, j \leq n$. Thus, $B = A^*$. Now, if X and Y are arbitrary matrices over the basis (e_1, \dots, e_n) , denoting as usual the j th column of X by X^j , and similarly for Y , a simple calculation shows that

$$Y^* X = (X^j \cdot Y^i)_{1 \leq i, j \leq n}.$$

Then it is immediately verified that if $X = Y = A$, then $A^* A = A A^* = I_n$ iff the column vectors (A^1, \dots, A^n) form an orthonormal basis. Thus, from (1), we see that (2) is clear. \square

Proposition 10.12 shows that the inverse of an isometry f is its adjoint f^* . Proposition 10.12 also motivates the following definition.

Definition 12.6. A complex $n \times n$ matrix is a *unitary matrix* if

$$A A^* = A^* A = I_n.$$

Remarks:

- (1) The conditions $A A^* = I_n$, $A^* A = I_n$, and $A^{-1} = A^*$ are equivalent. Given any two orthonormal bases (u_1, \dots, u_n) and (v_1, \dots, v_n) , if P is the change of basis matrix from (u_1, \dots, u_n) to (v_1, \dots, v_n) , it is easy to show that the matrix P is unitary. The proof of Proposition 12.11 (3) also shows that if f is an isometry, then the image of an orthonormal basis (u_1, \dots, u_n) is an orthonormal basis.
- (2) Using the explicit formula for the determinant, we see immediately that

$$\det(\overline{A}) = \overline{\det(A)}.$$

If f is unitary and A is its matrix with respect to any orthonormal basis, from $A A^* = I$, we get

$$\det(A A^*) = \det(A) \det(A^*) = \det(A) \overline{\det(A)} = \det(A) \overline{\det(A)} = |\det(A)|^2,$$

and so $|\det(A)| = 1$. It is clear that the isometries of a Hermitian space of dimension n form a group, and that the isometries of determinant $+1$ form a subgroup.

This leads to the following definition.

Definition 12.7. Given a Hermitian space E of dimension n , the set of isometries $f: E \rightarrow E$ forms a subgroup of $\mathbf{GL}(E, \mathbb{C})$ denoted by $\mathbf{U}(E)$, or $\mathbf{U}(n)$ when $E = \mathbb{C}^n$, called the *unitary group (of E)*. For every isometry f we have $|\det(f)| = 1$, where $\det(f)$ denotes the determinant of f . The isometries such that $\det(f) = 1$ are called *rotations, or proper isometries, or proper unitary transformations*, and they form a subgroup of the special linear group $\mathbf{SL}(E, \mathbb{C})$ (and of $\mathbf{U}(E)$), denoted by $\mathbf{SU}(E)$, or $\mathbf{SU}(n)$ when $E = \mathbb{C}^n$, called the *special unitary group (of E)*. The isometries such that $\det(f) \neq 1$ are called *improper isometries, or improper unitary transformations, or flip transformations*.

A very important example of unitary matrices is provided by Fourier matrices (up to a factor of \sqrt{n}), matrices that arise in the various versions of the discrete Fourier transform. For more on this topic, see the problems, and Strang [104, 106].

Now that we have the definition of a unitary matrix, we can explain how the Gram–Schmidt orthonormalization procedure immediately yields the QR -decomposition for matrices.

Proposition 12.13. *Given any $n \times n$ complex matrix A , if A is invertible, then there is a unitary matrix Q and an upper triangular matrix R with positive diagonal entries such that $A = QR$.*

The proof is absolutely the same as in the real case!

We have the following version of the Hadamard inequality for complex matrices. The proof is essentially the same as in the Euclidean case but it uses Proposition 12.13 instead of Proposition 10.13.

Proposition 12.14. (Hadamard) *For any complex $n \times n$ matrix $A = (a_{ij})$, we have*

$$|\det(A)| \leq \prod_{i=1}^n \left(\sum_{j=1}^n |a_{ij}|^2 \right)^{1/2} \quad \text{and} \quad |\det(A)| \leq \prod_{j=1}^n \left(\sum_{i=1}^n |a_{ij}|^2 \right)^{1/2}.$$

Moreover, equality holds iff either A has a zero column in the left inequality or a zero row in the right inequality, or A is unitary.

We also have the following version of Proposition 10.15 for Hermitian matrices. The proof of Proposition 10.15 goes through because the Cholesky decomposition for a Hermitian positive definite A matrix holds in the form $A = B^*B$, where B is upper triangular with positive diagonal entries. The details are left to the reader.

Proposition 12.15. (Hadamard) *For any complex $n \times n$ matrix $A = (a_{ij})$, if A is Hermitian positive semidefinite, then we have*

$$\det(A) \leq \prod_{i=1}^n a_{ii}.$$

Moreover, if A is positive definite, then equality holds iff A is a diagonal matrix.

Due to space limitations, we will not study the isometries of a Hermitian space in this chapter. However, the reader will find such a study in the supplements on the web site (see <http://www.cis.upenn.edu/~jean/gbooks/geom2.html>).

12.5 Orthogonal Projections and Involutions

In this section, we assume that the field K is not a field of characteristic 2. Recall that a linear map $f: E \rightarrow E$ is an *involution* iff $f^2 = \text{id}$, and is *idempotent* iff $f^2 = f$. We know from Proposition 4.7 that if f is idempotent, then

$$E = \text{Im}(f) \oplus \text{Ker}(f),$$

and that the restriction of f to its image is the identity. For this reason, a linear involution is called a *projection*. The connection between involutions and projections is given by the following simple proposition.

Proposition 12.16. *For any linear map $f: E \rightarrow E$, we have $f^2 = \text{id}$ iff $\frac{1}{2}(\text{id} - f)$ is a projection iff $\frac{1}{2}(\text{id} + f)$ is a projection; in this case, f is equal to the difference of the two projections $\frac{1}{2}(\text{id} + f)$ and $\frac{1}{2}(\text{id} - f)$.*

Proof. We have

$$\left(\frac{1}{2}(\text{id} - f)\right)^2 = \frac{1}{4}(\text{id} - 2f + f^2)$$

so

$$\left(\frac{1}{2}(\text{id} - f)\right)^2 = \frac{1}{2}(\text{id} - f) \quad \text{iff} \quad f^2 = \text{id}.$$

We also have

$$\left(\frac{1}{2}(\text{id} + f)\right)^2 = \frac{1}{4}(\text{id} + 2f + f^2),$$

so

$$\left(\frac{1}{2}(\text{id} + f)\right)^2 = \frac{1}{2}(\text{id} + f) \quad \text{iff} \quad f^2 = \text{id}.$$

Obviously, $f = \frac{1}{2}(\text{id} + f) - \frac{1}{2}(\text{id} - f)$. □

Let $U^+ = \text{Ker}(\frac{1}{2}(\text{id} - f))$ and let $U^- = \text{Im}(\frac{1}{2}(\text{id} - f))$. If $f^2 = \text{id}$, then

$$(\text{id} + f) \circ (\text{id} - f) = \text{id} - f^2 = \text{id} - \text{id} = 0,$$

which implies that

$$\text{Im}\left(\frac{1}{2}(\text{id} + f)\right) \subseteq \text{Ker}\left(\frac{1}{2}(\text{id} - f)\right).$$

Conversely, if $u \in \text{Ker} \left(\frac{1}{2}(\text{id} - f) \right)$, then $f(u) = u$, so

$$\frac{1}{2}(\text{id} + f)(u) = \frac{1}{2}(u + u) = u,$$

and thus

$$\text{Ker} \left(\frac{1}{2}(\text{id} - f) \right) \subseteq \text{Im} \left(\frac{1}{2}(\text{id} + f) \right).$$

Therefore,

$$U^+ = \text{Ker} \left(\frac{1}{2}(\text{id} - f) \right) = \text{Im} \left(\frac{1}{2}(\text{id} + f) \right),$$

and so, $f(u) = u$ on U^+ and $f(u) = -u$ on U^- . The involutions of E that are unitary transformations are characterized as follows.

Proposition 12.17. *Let $f \in \mathbf{GL}(E)$ be an involution. The following properties are equivalent:*

- (a) *The map f is unitary; that is, $f \in \mathbf{U}(E)$.*
 - (b) *The subspaces $U^- = \text{Im}(\frac{1}{2}(\text{id} - f))$ and $U^+ = \text{Im}(\frac{1}{2}(\text{id} + f))$ are orthogonal.*
- Furthermore, if E is finite-dimensional, then (a) and (b) are equivalent to*
- (c) *The map is self-adjoint; that is, $f = f^*$.*

Proof. If f is unitary, then from $\langle f(u), f(v) \rangle = \langle u, v \rangle$ for all $u, v \in E$, we see that if $u \in U^+$ and $v \in U^-$, we get

$$\langle u, v \rangle = \langle f(u), f(v) \rangle = \langle u, -v \rangle = -\langle u, v \rangle,$$

so $2\langle u, v \rangle = 0$, which implies $\langle u, v \rangle = 0$, that is, U^+ and U^- are orthogonal. Thus, (a) implies (b).

Conversely, if (b) holds, since $f(u) = u$ on U^+ and $f(u) = -u$ on U^- , we see that $\langle f(u), f(v) \rangle = \langle u, v \rangle$ if $u, v \in U^+$ or if $u, v \in U^-$. Since $E = U^+ \oplus U^-$ and since U^+ and U^- are orthogonal, we also have $\langle f(u), f(v) \rangle = \langle u, v \rangle$ for all $u, v \in E$, and (b) implies (a).

If E is finite-dimensional, the adjoint f^* of f exists, and we know that $f^{-1} = f^*$. Since f is an involution, $f^2 = \text{id}$, which implies that $f^* = f^{-1} = f$. \square

A unitary involution is the identity on $U^+ = \text{Im}(\frac{1}{2}(\text{id} + f))$, and $f(v) = -v$ for all $v \in U^- = \text{Im}(\frac{1}{2}(\text{id} - f))$. Furthermore, E is an orthogonal direct sum $E = U^+ \oplus U^-$. We say that f is an *orthogonal reflection* about U^+ . In the special case where U^+ is a hyperplane, we say that f is a *hyperplane reflection*. We already studied hyperplane reflections in the Euclidean case; see Chapter 11.

If $f: E \rightarrow E$ is a projection ($f^2 = f$), then

$$(\text{id} - 2f)^2 = \text{id} - 4f + 4f^2 = \text{id} - 4f + 4f = \text{id},$$

so $\text{id} - 2f$ is an involution. As a consequence, we get the following result.

Proposition 12.18. *If $f: E \rightarrow E$ is a projection ($f^2 = f$), then $\text{Ker}(f)$ and $\text{Im}(f)$ are orthogonal iff $f^* = f$.*

Proof. Apply Proposition 12.17 to $g = \text{id} - 2f$. Since $\text{id} - g = 2f$ we have

$$U^+ = \text{Ker} \left(\frac{1}{2}(\text{id} - g) \right) = \text{Ker}(f)$$

and

$$U^- = \text{Im} \left(\frac{1}{2}(\text{id} - g) \right) = \text{Im}(f),$$

which proves the proposition. □

A projection such that $f = f^*$ is called an *orthogonal projection*.

If (a_1, \dots, a_k) are k linearly independent vectors in \mathbb{R}^n , let us determine the matrix P of the orthogonal projection onto the subspace of \mathbb{R}^n spanned by (a_1, \dots, a_k) . Let A be the $n \times k$ matrix whose j th column consists of the coordinates of the vector a_j over the canonical basis (e_1, \dots, e_n) .

Any vector in the subspace (a_1, \dots, a_k) is a linear combination of the form Ax , for some $x \in \mathbb{R}^k$. Given any $y \in \mathbb{R}^n$, the orthogonal projection $Py = Ax$ of y onto the subspace spanned by (a_1, \dots, a_k) is the vector Ax such that $y - Ax$ is orthogonal to the subspace spanned by (a_1, \dots, a_k) (prove it). This means that $y - Ax$ is orthogonal to every a_j , which is expressed by

$$A^\top(y - Ax) = 0;$$

that is,

$$A^\top Ax = A^\top y.$$

The matrix $A^\top A$ is invertible because A has full rank k , thus we get

$$x = (A^\top A)^{-1} A^\top y,$$

and so

$$Py = Ax = A(A^\top A)^{-1} A^\top y.$$

Therefore, the matrix P of the projection onto the subspace spanned by (a_1, \dots, a_k) is given by

$$P = A(A^\top A)^{-1} A^\top.$$

The reader should check that $P^2 = P$ and $P^\top = P$.

12.6 Dual Norms

In the remark following the proof of Proposition 7.8, we explained that if $(E, \|\cdot\|)$ and $(F, \|\cdot\|)$ are two normed vector spaces and if we let $\mathcal{L}(E; F)$ denote the set of all continuous (equivalently, bounded) linear maps from E to F , then, we can define the *operator norm* (or *subordinate norm*) $\|\cdot\|$ on $\mathcal{L}(E; F)$ as follows: for every $f \in \mathcal{L}(E; F)$,

$$\|f\| = \sup_{\substack{x \in E \\ x \neq 0}} \frac{\|f(x)\|}{\|x\|} = \sup_{\substack{x \in E \\ \|x\|=1}} \|f(x)\|.$$

In particular, if $F = \mathbb{C}$, then $\mathcal{L}(E; F) = E'$ is the *dual space* of E , and we get the operator norm denoted by $\|\cdot\|_*$ given by

$$\|f\|_* = \sup_{\substack{x \in E \\ \|x\|=1}} |f(x)|.$$

The norm $\|\cdot\|_*$ is called the *dual norm* of $\|\cdot\|$ on E' .

Let us now assume that E is a finite-dimensional Hermitian space, in which case $E' = E^*$. Theorem 12.5 implies that for every linear form $f \in E^*$, there is a unique vector $y \in E$ so that

$$f(x) = \langle x, y \rangle,$$

for all $x \in E$, and so we can write

$$\|f\|_* = \sup_{\substack{x \in E \\ \|x\|=1}} |\langle x, y \rangle|.$$

The above suggests defining a norm $\|\cdot\|^D$ on E .

Definition 12.8. If E is a finite-dimensional Hermitian space and $\|\cdot\|$ is any norm on E , for any $y \in E$ we let

$$\|y\|^D = \sup_{\substack{x \in E \\ \|x\|=1}} |\langle x, y \rangle|,$$

be the *dual norm* of $\|\cdot\|$ (on E). If E is a real Euclidean space, then the dual norm is defined by

$$\|y\|^D = \sup_{\substack{x \in E \\ \|x\|=1}} \langle x, y \rangle$$

for all $y \in E$.

Beware that $\|\cdot\|$ is generally *not* the Hermitian norm associated with the Hermitian inner product. The dual norm shows up in convex programming; see Boyd and Vandenberghe [17], Chapters 2, 3, 6, 9.

The fact that $\|\cdot\|^D$ is a norm follows from the fact that $\|\cdot\|_*$ is a norm and can also be checked directly. It is worth noting that the triangle inequality for $\|\cdot\|^D$ comes “for free,” in the sense that it holds for any function $p: E \rightarrow \mathbb{R}$. Indeed, we have

$$\begin{aligned} p^D(x+y) &= \sup_{p(z)=1} |\langle z, x+y \rangle| \\ &= \sup_{p(z)=1} (|\langle z, x \rangle + \langle z, y \rangle|) \\ &\leq \sup_{p(z)=1} (|\langle z, x \rangle| + |\langle z, y \rangle|) \\ &\leq \sup_{p(z)=1} |\langle z, x \rangle| + \sup_{p(z)=1} |\langle z, y \rangle| \\ &= p^D(x) + p^D(y). \end{aligned}$$

If $p: E \rightarrow \mathbb{R}$ is a function such that

- (1) $p(x) \geq 0$ for all $x \in E$, and $p(x) = 0$ iff $x = 0$;
- (2) $p(\lambda x) = |\lambda|p(x)$, for all $x \in E$ and all $\lambda \in \mathbb{C}$;
- (3) p is continuous, in the sense that for some basis (e_1, \dots, e_n) of E , the function

$$(x_1, \dots, x_n) \mapsto p(x_1 e_1 + \dots + x_n e_n)$$

from \mathbb{C}^n to \mathbb{R} is continuous;

then we say that p is a *pre-norm*. Obviously, every norm is a pre-norm, but a pre-norm may not satisfy the triangle inequality. However, we just showed that the dual norm of any pre-norm is actually a norm.

Since E is finite dimensional, the unit sphere $S^{n-1} = \{x \in E \mid \|x\| = 1\}$ is compact, so there is some $x_0 \in S^{n-1}$ such that

$$\|y\|^D = |\langle x_0, y \rangle|.$$

If $\langle x_0, y \rangle = \rho e^{i\theta}$, with $\rho \geq 0$, then

$$|\langle e^{-i\theta} x_0, y \rangle| = |e^{-i\theta} \langle x_0, y \rangle| = |e^{-i\theta} \rho e^{i\theta}| = \rho,$$

so

$$\|y\|^D = \rho = |\langle e^{-i\theta} x_0, y \rangle|,$$

with $\|e^{-i\theta} x_0\| = \|x_0\| = 1$. On the other hand,

$$\Re \langle x, y \rangle \leq |\langle x, y \rangle|,$$

so we get

$$\|y\|^D = \sup_{\substack{x \in E \\ \|x\|=1}} |\langle x, y \rangle| = \sup_{\substack{x \in E \\ \|x\|=1}} \Re \langle x, y \rangle.$$

Proposition 12.19. *For all $x, y \in E$, we have*

$$\begin{aligned} |\langle x, y \rangle| &\leq \|x\| \|y\|^D \\ |\langle x, y \rangle| &\leq \|x\|^D \|y\|. \end{aligned}$$

Proof. If $x = 0$, then $\langle x, y \rangle = 0$ and these inequalities are trivial. If $x \neq 0$, since $\|x/\|x\|\| = 1$, by definition of $\|y\|^D$, we have

$$|\langle x/\|x\|, y \rangle| \leq \sup_{\|z\|=1} |\langle z, y \rangle| = \|y\|^D,$$

which yields

$$|\langle x, y \rangle| \leq \|x\| \|y\|^D.$$

The second inequality holds because $|\langle x, y \rangle| = |\langle y, x \rangle|$. □

It is not hard to show that

$$\begin{aligned} \|y\|_1^D &= \|y\|_\infty \\ \|y\|_\infty^D &= \|y\|_1 \\ \|y\|_2^D &= \|y\|_2. \end{aligned}$$

Thus, the Euclidean norm is autodual. More generally, if $p, q \geq 1$ and $1/p + 1/q = 1$, we have

$$\|y\|_p^D = \|y\|_q.$$

It can also be shown that the dual of the spectral norm is the trace norm (or nuclear norm) from Section 16.3. We close this section by stating the following duality theorem.

Theorem 12.20. *If E is a finite-dimensional Hermitian space, then for any norm $\|\cdot\|$ on E , we have*

$$\|y\|^{DD} = \|y\|$$

for all $y \in E$.

Proof. By Proposition 12.19, we have

$$|\langle x, y \rangle| \leq \|x\|^D \|y\|,$$

so we get

$$\|y\|^{DD} = \sup_{\|x\|^D=1} |\langle x, y \rangle| \leq \|y\|, \quad \text{for all } y \in E.$$

It remains to prove that

$$\|y\| \leq \|y\|^{DD}, \quad \text{for all } y \in E.$$

Proofs of this fact can be found in Horn and Johnson [57] (Section 5.5), and in Serre [96] (Chapter 7). The proof makes use of the fact that a nonempty, closed, convex set has a

supporting hyperplane through each of its boundary points, a result known as *Minkowski's lemma*. This result is a consequence of the *Hahn–Banach theorem*; see Gallier [44]. We give the proof in the case where E is a real Euclidean space. Some minor modifications have to be made when dealing with complex vector spaces and are left as an exercise.

Since the unit ball $B = \{z \in E \mid \|z\| \leq 1\}$ is closed and convex, the Minkowski lemma says for every x such that $\|x\| = 1$, there is an affine map g , of the form

$$g(z) = \langle z, w \rangle - \langle x, w \rangle$$

with $\|w\| = 1$, such that $g(x) = 0$ and $g(z) \leq 0$ for all z such that $\|z\| \leq 1$. Then, it is clear that

$$\sup_{\|z\|=1} \langle z, w \rangle = \langle x, w \rangle,$$

and so

$$\|w\|^D = \langle x, w \rangle.$$

It follows that

$$\|x\|^{DD} \geq \langle w / \|w\|^D, x \rangle = \frac{\langle x, w \rangle}{\|w\|^D} = 1 = \|x\|$$

for all x such that $\|x\| = 1$. By homogeneity, this is true for all $y \in E$, which completes the proof in the real case. When E is a complex vector space, we have to view the unit ball B as a closed convex set in \mathbb{R}^{2n} and we use the fact that there is real affine map of the form

$$g(z) = \Re \langle z, w \rangle - \Re \langle x, w \rangle$$

such that $g(x) = 0$ and $g(z) \leq 0$ for all z with $\|z\| = 1$, so that $\|w\|^D = \Re \langle x, w \rangle$. \square

More details on dual norms and unitarily invariant norms can be found in Horn and Johnson [57] (Chapters 5 and 7).

12.7 Summary

The main concepts and results of this chapter are listed below:

- *Semilinear maps*.
- *Sesquilinear forms; Hermitian forms*.
- *Quadratic form* associated with a sesquilinear form.
- *Polarization identities*.
- *Positive and positive definite Hermitian forms; pre-Hilbert spaces, Hermitian spaces*.
- *Gram matrix* associated with a Hermitian product.

- The *Cauchy–Schwarz inequality* and the *Minkowski inequality*.
- *Hermitian inner product*, *Hermitian norm*.
- The *parallelogram law*.
- The musical isomorphisms $\flat: \overline{E} \rightarrow E^*$ and $\sharp: E^* \rightarrow \overline{E}$; Theorem 12.5 (E is finite-dimensional).
- The *adjoint* of a linear map (with respect to a Hermitian inner product).
- Existence of orthonormal bases in a Hermitian space (Proposition 12.8).
- *Gram–Schmidt orthonormalization procedure*.
- *Linear isometries (unitary transformations)*.
- The *unitary group*, *unitary matrices*.
- The *unitary group* $\mathbf{U}(n)$;
- The *special unitary group* $\mathbf{SU}(n)$.
- *QR-Decomposition* for invertible matrices.
- The *Hadamard inequality* for complex matrices.
- The *Hadamard inequality* for Hermitian positive semidefinite matrices.
- Orthogonal projections and involutions; orthogonal reflections.
- Dual norms.

Chapter 13

Spectral Theorems in Euclidean and Hermitian Spaces

13.1 Introduction

The goal of this chapter is to show that there are nice normal forms for symmetric matrices, skew-symmetric matrices, orthogonal matrices, and normal matrices. The spectral theorem for symmetric matrices states that symmetric matrices have real eigenvalues and that they can be diagonalized over an orthonormal basis. The spectral theorem for Hermitian matrices states that Hermitian matrices also have real eigenvalues and that they can be diagonalized over a complex orthonormal basis. Normal real matrices can be block diagonalized over an orthonormal basis with blocks having size at most two, and there are refinements of this normal form for skew-symmetric and orthogonal matrices.

13.2 Normal Linear Maps

We begin by studying normal maps, to understand the structure of their eigenvalues and eigenvectors. This section and the next two were inspired by Lang [67], Artin [4], Mac Lane and Birkhoff [73], Berger [8], and Bertin [12].

Definition 13.1. Given a Euclidean space E , a linear map $f: E \rightarrow E$ is *normal* if

$$f \circ f^* = f^* \circ f.$$

A linear map $f: E \rightarrow E$ is *self-adjoint* if $f = f^*$, *skew-self-adjoint* if $f = -f^*$, and *orthogonal* if $f \circ f^* = f^* \circ f = \text{id}$.

Obviously, a self-adjoint, skew-self-adjoint, or orthogonal linear map is a normal linear map. Our first goal is to show that for every normal linear map $f: E \rightarrow E$, there is an orthonormal basis (w.r.t. $\langle -, - \rangle$) such that the matrix of f over this basis has an especially

nice form: It is a block diagonal matrix in which the blocks are either one-dimensional matrices (i.e., single entries) or two-dimensional matrices of the form

$$\begin{pmatrix} \lambda & \mu \\ -\mu & \lambda \end{pmatrix}.$$

This normal form can be further refined if f is self-adjoint, skew-self-adjoint, or orthogonal. As a first step, we show that f and f^* have the same kernel when f is normal.

Proposition 13.1. *Given a Euclidean space E , if $f: E \rightarrow E$ is a normal linear map, then $\text{Ker } f = \text{Ker } f^*$.*

Proof. First, let us prove that

$$\langle f(u), f(v) \rangle = \langle f^*(u), f^*(v) \rangle$$

for all $u, v \in E$. Since f^* is the adjoint of f and $f \circ f^* = f^* \circ f$, we have

$$\begin{aligned} \langle f(u), f(u) \rangle &= \langle u, (f^* \circ f)(u) \rangle, \\ &= \langle u, (f \circ f^*)(u) \rangle, \\ &= \langle f^*(u), f^*(u) \rangle. \end{aligned}$$

Since $\langle -, - \rangle$ is positive definite,

$$\begin{aligned} \langle f(u), f(u) \rangle = 0 &\quad \text{iff} \quad f(u) = 0, \\ \langle f^*(u), f^*(u) \rangle = 0 &\quad \text{iff} \quad f^*(u) = 0, \end{aligned}$$

and since

$$\langle f(u), f(u) \rangle = \langle f^*(u), f^*(u) \rangle,$$

we have

$$f(u) = 0 \quad \text{iff} \quad f^*(u) = 0.$$

Consequently, $\text{Ker } f = \text{Ker } f^*$. □

The next step is to show that for every linear map $f: E \rightarrow E$ there is some subspace W of dimension 1 or 2 such that $f(W) \subseteq W$. When $\dim(W) = 1$, the subspace W is actually an eigenspace for some real eigenvalue of f . Furthermore, when f is normal, there is a subspace W of dimension 1 or 2 such that $f(W) \subseteq W$ and $f^*(W) \subseteq W$. The difficulty is that the eigenvalues of f are not necessarily real. One way to get around this problem is to complexify both the vector space E and the inner product $\langle -, - \rangle$.

Every real vector space E can be embedded into a complex vector space $E_{\mathbb{C}}$, and every linear map $f: E \rightarrow E$ can be extended to a linear map $f_{\mathbb{C}}: E_{\mathbb{C}} \rightarrow E_{\mathbb{C}}$.

Definition 13.2. Given a real vector space E , let $E_{\mathbb{C}}$ be the structure $E \times E$ under the addition operation

$$(u_1, u_2) + (v_1, v_2) = (u_1 + v_1, u_2 + v_2),$$

and let multiplication by a complex scalar $z = x + iy$ be defined such that

$$(x + iy) \cdot (u, v) = (xu - yv, yu + xv).$$

The space $E_{\mathbb{C}}$ is called the *complexification* of E .

It is easily shown that the structure $E_{\mathbb{C}}$ is a complex vector space. It is also immediate that

$$(0, v) = i(v, 0),$$

and thus, identifying E with the subspace of $E_{\mathbb{C}}$ consisting of all vectors of the form $(u, 0)$, we can write

$$(u, v) = u + iv.$$

Observe that if (e_1, \dots, e_n) is a basis of E (a real vector space), then (e_1, \dots, e_n) is also a basis of $E_{\mathbb{C}}$ (recall that e_i is an abbreviation for $(e_i, 0)$).

A linear map $f: E \rightarrow E$ is extended to the linear map $f_{\mathbb{C}}: E_{\mathbb{C}} \rightarrow E_{\mathbb{C}}$ defined such that

$$f_{\mathbb{C}}(u + iv) = f(u) + if(v).$$

For any basis (e_1, \dots, e_n) of E , the matrix $M(f)$ representing f over (e_1, \dots, e_n) is identical to the matrix $M(f_{\mathbb{C}})$ representing $f_{\mathbb{C}}$ over (e_1, \dots, e_n) , where we view (e_1, \dots, e_n) as a basis of $E_{\mathbb{C}}$. As a consequence, $\det(zI - M(f)) = \det(zI - M(f_{\mathbb{C}}))$, which means that f and $f_{\mathbb{C}}$ have the same characteristic polynomial (which has real coefficients). We know that every polynomial of degree n with real (or complex) coefficients always has n complex roots (counted with their multiplicity), and the roots of $\det(zI - M(f_{\mathbb{C}}))$ that are real (if any) are the eigenvalues of f .

Next, we need to extend the inner product on E to an inner product on $E_{\mathbb{C}}$.

The inner product $\langle -, - \rangle$ on a Euclidean space E is extended to the Hermitian positive definite form $\langle -, - \rangle_{\mathbb{C}}$ on $E_{\mathbb{C}}$ as follows:

$$\langle u_1 + iv_1, u_2 + iv_2 \rangle_{\mathbb{C}} = \langle u_1, u_2 \rangle + \langle v_1, v_2 \rangle + i(\langle u_2, v_1 \rangle - \langle u_1, v_2 \rangle).$$

It is easily verified that $\langle -, - \rangle_{\mathbb{C}}$ is indeed a Hermitian form that is positive definite, and it is clear that $\langle -, - \rangle_{\mathbb{C}}$ agrees with $\langle -, - \rangle$ on real vectors. Then, given any linear map $f: E \rightarrow E$, it is easily verified that the map $f_{\mathbb{C}}^*$ defined such that

$$f_{\mathbb{C}}^*(u + iv) = f^*(u) + if^*(v)$$

for all $u, v \in E$ is the adjoint of $f_{\mathbb{C}}$ w.r.t. $\langle -, - \rangle_{\mathbb{C}}$.

Assuming again that E is a Hermitian space, observe that Proposition 13.1 also holds. We deduce the following corollary.

Proposition 13.2. *Given a Hermitian space E , for any normal linear map $f: E \rightarrow E$, we have $\text{Ker}(f) \cap \text{Im}(f) = (0)$.*

Proof. Assume $v \in \text{Ker}(f) \cap \text{Im}(f) = (0)$, which means that $v = f(u)$ for some $u \in E$, and $f(v) = 0$. By Proposition 13.1, $\text{Ker}(f) = \text{Ker}(f^*)$, so $f(v) = 0$ implies that $f^*(v) = 0$. Consequently,

$$\begin{aligned} 0 &= \langle f^*(v), u \rangle \\ &= \langle v, f(u) \rangle \\ &= \langle v, v \rangle, \end{aligned}$$

and thus, $v = 0$. □

We also have the following crucial proposition relating the eigenvalues of f and f^* .

Proposition 13.3. *Given a Hermitian space E , for any normal linear map $f: E \rightarrow E$, a vector u is an eigenvector of f for the eigenvalue λ (in \mathbb{C}) iff u is an eigenvector of f^* for the eigenvalue $\bar{\lambda}$.*

Proof. First, it is immediately verified that the adjoint of $f - \lambda \text{id}$ is $f^* - \bar{\lambda} \text{id}$. Furthermore, $f - \lambda \text{id}$ is normal. Indeed,

$$\begin{aligned} (f - \lambda \text{id}) \circ (f - \lambda \text{id})^* &= (f - \lambda \text{id}) \circ (f^* - \bar{\lambda} \text{id}), \\ &= f \circ f^* - \bar{\lambda} f - \lambda f^* + \lambda \bar{\lambda} \text{id}, \\ &= f^* \circ f - \lambda f^* - \bar{\lambda} f + \bar{\lambda} \lambda \text{id}, \\ &= (f^* - \bar{\lambda} \text{id}) \circ (f - \lambda \text{id}), \\ &= (f - \lambda \text{id})^* \circ (f - \lambda \text{id}). \end{aligned}$$

Applying Proposition 13.1 to $f - \lambda \text{id}$, for every nonnull vector u , we see that

$$(f - \lambda \text{id})(u) = 0 \quad \text{iff} \quad (f^* - \bar{\lambda} \text{id})(u) = 0,$$

which is exactly the statement of the proposition. □

The next proposition shows a very important property of normal linear maps: Eigenvectors corresponding to distinct eigenvalues are orthogonal.

Proposition 13.4. *Given a Hermitian space E , for any normal linear map $f: E \rightarrow E$, if u and v are eigenvectors of f associated with the eigenvalues λ and μ (in \mathbb{C}) where $\lambda \neq \mu$, then $\langle u, v \rangle = 0$.*

Proof. Let us compute $\langle f(u), v \rangle$ in two different ways. Since v is an eigenvector of f for μ , by Proposition 13.3, v is also an eigenvector of f^* for $\bar{\mu}$, and we have

$$\langle f(u), v \rangle = \langle \lambda u, v \rangle = \lambda \langle u, v \rangle$$

and

$$\langle f(u), v \rangle = \langle u, f^*(v) \rangle = \langle u, \bar{\mu}v \rangle = \mu \langle u, v \rangle,$$

where the last identity holds because of the semilinearity in the second argument, and thus

$$\lambda \langle u, v \rangle = \mu \langle u, v \rangle,$$

that is,

$$(\lambda - \mu) \langle u, v \rangle = 0,$$

which implies that $\langle u, v \rangle = 0$, since $\lambda \neq \mu$. \square

We can also show easily that the eigenvalues of a self-adjoint linear map are real.

Proposition 13.5. *Given a Hermitian space E , all the eigenvalues of any self-adjoint linear map $f: E \rightarrow E$ are real.*

Proof. Let z (in \mathbb{C}) be an eigenvalue of f and let u be an eigenvector for z . We compute $\langle f(u), u \rangle$ in two different ways. We have

$$\langle f(u), u \rangle = \langle zu, u \rangle = z \langle u, u \rangle,$$

and since $f = f^*$, we also have

$$\langle f(u), u \rangle = \langle u, f^*(u) \rangle = \langle u, f(u) \rangle = \langle u, zu \rangle = \bar{z} \langle u, u \rangle.$$

Thus,

$$z \langle u, u \rangle = \bar{z} \langle u, u \rangle,$$

which implies that $z = \bar{z}$, since $u \neq 0$, and z is indeed real. \square

There is also a version of Proposition 13.5 for a (real) Euclidean space E and a self-adjoint map $f: E \rightarrow E$.

Proposition 13.6. *Given a Euclidean space E , if $f: E \rightarrow E$ is any self-adjoint linear map, then every eigenvalue λ of $f_{\mathbb{C}}$ is real and is actually an eigenvalue of f (which means that there is some real eigenvector $u \in E$ such that $f(u) = \lambda u$). Therefore, all the eigenvalues of f are real.*

Proof. Let $E_{\mathbb{C}}$ be the complexification of E , $\langle -, - \rangle_{\mathbb{C}}$ the complexification of the inner product $\langle -, - \rangle$ on E , and $f_{\mathbb{C}}: E_{\mathbb{C}} \rightarrow E_{\mathbb{C}}$ the complexification of $f: E \rightarrow E$. By definition of $f_{\mathbb{C}}$ and $\langle -, - \rangle_{\mathbb{C}}$, if f is self-adjoint, we have

$$\begin{aligned} \langle f_{\mathbb{C}}(u_1 + iv_1), u_2 + iv_2 \rangle_{\mathbb{C}} &= \langle f(u_1) + if(v_1), u_2 + iv_2 \rangle_{\mathbb{C}} \\ &= \langle f(u_1), u_2 \rangle + \langle f(v_1), v_2 \rangle + i(\langle u_2, f(v_1) \rangle - \langle f(u_1), v_2 \rangle) \\ &= \langle u_1, f(u_2) \rangle + \langle v_1, f(v_2) \rangle + i(\langle f(u_2), v_1 \rangle - \langle u_1, f(v_2) \rangle) \\ &= \langle u_1 + iv_1, f(u_2) + if(v_2) \rangle_{\mathbb{C}} \\ &= \langle u_1 + iv_1, f_{\mathbb{C}}(u_2 + iv_2) \rangle_{\mathbb{C}}, \end{aligned}$$

which shows that $f_{\mathbb{C}}$ is also self-adjoint with respect to $\langle -, - \rangle_{\mathbb{C}}$.

As we pointed out earlier, f and $f_{\mathbb{C}}$ have the same characteristic polynomial $\det(zI - f_{\mathbb{C}}) = \det(zI - f)$, which is a polynomial with real coefficients. Proposition 13.5 shows that the zeros of $\det(zI - f_{\mathbb{C}}) = \det(zI - f)$ are all real, and for each real zero λ of $\det(zI - f)$, the linear map $\lambda \text{id} - f$ is singular, which means that there is some nonzero $u \in E$ such that $f(u) = \lambda u$. Therefore, all the eigenvalues of f are real. \square

Given any subspace W of a Euclidean space E , recall that the *orthogonal complement* W^{\perp} of W is the subspace defined such that

$$W^{\perp} = \{u \in E \mid \langle u, w \rangle = 0, \text{ for all } w \in W\}.$$

Recall from Proposition 10.9 that $E = W \oplus W^{\perp}$ (this can be easily shown, for example, by constructing an orthonormal basis of E using the Gram–Schmidt orthonormalization procedure). The same result also holds for Hermitian spaces; see Proposition 12.10.

As a warm up for the proof of Theorem 13.10, let us prove that every self-adjoint map on a Euclidean space can be diagonalized with respect to an orthonormal basis of eigenvectors.

Theorem 13.7. (*Spectral theorem for self-adjoint linear maps on a Euclidean space*) *Given a Euclidean space E of dimension n , for every self-adjoint linear map $f: E \rightarrow E$, there is an orthonormal basis (e_1, \dots, e_n) of eigenvectors of f such that the matrix of f w.r.t. this basis is a diagonal matrix*

$$\begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{pmatrix},$$

with $\lambda_i \in \mathbb{R}$.

Proof. We proceed by induction on the dimension n of E as follows. If $n = 1$, the result is trivial. Assume now that $n \geq 2$. From Proposition 13.6, all the eigenvalues of f are real, so pick some eigenvalue $\lambda \in \mathbb{R}$, and let w be some eigenvector for λ . By dividing w by its norm, we may assume that w is a unit vector. Let W be the subspace of dimension 1 spanned by w . Clearly, $f(W) \subseteq W$. We claim that $f(W^{\perp}) \subseteq W^{\perp}$, where W^{\perp} is the orthogonal complement of W .

Indeed, for any $v \in W^{\perp}$, that is, if $\langle v, w \rangle = 0$, because f is self-adjoint and $f(w) = \lambda w$, we have

$$\begin{aligned} \langle f(v), w \rangle &= \langle v, f(w) \rangle \\ &= \langle v, \lambda w \rangle \\ &= \lambda \langle v, w \rangle = 0 \end{aligned}$$

since $\langle v, w \rangle = 0$. Therefore,

$$f(W^\perp) \subseteq W^\perp.$$

Clearly, the restriction of f to W^\perp is self-adjoint, and we conclude by applying the induction hypothesis to W^\perp (whose dimension is $n - 1$). \square

We now come back to normal linear maps. One of the key points in the proof of Theorem 13.7 is that we found a subspace W with the property that $f(W) \subseteq W$ implies that $f(W^\perp) \subseteq W^\perp$. In general, this does not happen, but normal maps satisfy a stronger property which ensures that such a subspace exists.

The following proposition provides a condition that will allow us to show that a normal linear map can be diagonalized. It actually holds for any linear map. We found the inspiration for this proposition in Berger [8].

Proposition 13.8. *Given a Hermitian space E , for any linear map $f: E \rightarrow E$ and any subspace W of E , if $f(W) \subseteq W$, then $f^*(W^\perp) \subseteq W^\perp$. Consequently, if $f(W) \subseteq W$ and $f^*(W) \subseteq W$, then $f(W^\perp) \subseteq W^\perp$ and $f^*(W^\perp) \subseteq W^\perp$.*

Proof. If $u \in W^\perp$, then

$$\langle w, u \rangle = 0 \quad \text{for all } w \in W.$$

However,

$$\langle f(w), u \rangle = \langle w, f^*(u) \rangle,$$

and $f(W) \subseteq W$ implies that $f(w) \in W$. Since $u \in W^\perp$, we get

$$0 = \langle f(w), u \rangle = \langle w, f^*(u) \rangle,$$

which shows that $\langle w, f^*(u) \rangle = 0$ for all $w \in W$, that is, $f^*(u) \in W^\perp$. Therefore, we have $f^*(W^\perp) \subseteq W^\perp$.

We just proved that if $f(W) \subseteq W$, then $f^*(W^\perp) \subseteq W^\perp$. If we also have $f^*(W) \subseteq W$, then by applying the above fact to f^* , we get $f^{**}(W^\perp) \subseteq W^\perp$, and since $f^{**} = f$, this is just $f(W^\perp) \subseteq W^\perp$, which proves the second statement of the proposition. \square

It is clear that the above proposition also holds for Euclidean spaces.

Although we are ready to prove that for every normal linear map f (over a Hermitian space) there is an orthonormal basis of eigenvectors (see Theorem 13.11 below), we now return to real Euclidean spaces.

If $f: E \rightarrow E$ is a linear map and $w = u + iv$ is an eigenvector of $f_\mathbb{C}: E_\mathbb{C} \rightarrow E_\mathbb{C}$ for the eigenvalue $z = \lambda + i\mu$, where $u, v \in E$ and $\lambda, \mu \in \mathbb{R}$, since

$$f_\mathbb{C}(u + iv) = f(u) + if(v)$$

and

$$f_\mathbb{C}(u + iv) = (\lambda + i\mu)(u + iv) = \lambda u - \mu v + i(\mu u + \lambda v),$$

we have

$$f(u) = \lambda u - \mu v \quad \text{and} \quad f(v) = \mu u + \lambda v,$$

from which we immediately obtain

$$f_{\mathbb{C}}(u - iv) = (\lambda - i\mu)(u - iv),$$

which shows that $\bar{w} = u - iv$ is an eigenvector of $f_{\mathbb{C}}$ for $\bar{z} = \lambda - i\mu$. Using this fact, we can prove the following proposition.

Proposition 13.9. *Given a Euclidean space E , for any normal linear map $f: E \rightarrow E$, if $w = u + iv$ is an eigenvector of $f_{\mathbb{C}}$ associated with the eigenvalue $z = \lambda + i\mu$ (where $u, v \in E$ and $\lambda, \mu \in \mathbb{R}$), if $\mu \neq 0$ (i.e., z is not real) then $\langle u, v \rangle = 0$ and $\langle u, u \rangle = \langle v, v \rangle$, which implies that u and v are linearly independent, and if W is the subspace spanned by u and v , then $f(W) = W$ and $f^*(W) = W$. Furthermore, with respect to the (orthogonal) basis (u, v) , the restriction of f to W has the matrix*

$$\begin{pmatrix} \lambda & \mu \\ -\mu & \lambda \end{pmatrix}.$$

If $\mu = 0$, then λ is a real eigenvalue of f , and either u or v is an eigenvector of f for λ . If W is the subspace spanned by u if $u \neq 0$, or spanned by $v \neq 0$ if $u = 0$, then $f(W) \subseteq W$ and $f^(W) \subseteq W$.*

Proof. Since $w = u + iv$ is an eigenvector of $f_{\mathbb{C}}$, by definition it is nonnull, and either $u \neq 0$ or $v \neq 0$. From the fact stated just before Proposition 13.9, $u - iv$ is an eigenvector of $f_{\mathbb{C}}$ for $\lambda - i\mu$. It is easy to check that $f_{\mathbb{C}}$ is normal. However, if $\mu \neq 0$, then $\lambda + i\mu \neq \lambda - i\mu$, and from Proposition 13.4, the vectors $u + iv$ and $u - iv$ are orthogonal w.r.t. $\langle -, - \rangle_{\mathbb{C}}$, that is,

$$\langle u + iv, u - iv \rangle_{\mathbb{C}} = \langle u, u \rangle - \langle v, v \rangle + 2i\langle u, v \rangle = 0.$$

Thus, we get $\langle u, v \rangle = 0$ and $\langle u, u \rangle = \langle v, v \rangle$, and since $u \neq 0$ or $v \neq 0$, u and v are linearly independent. Since

$$f(u) = \lambda u - \mu v \quad \text{and} \quad f(v) = \mu u + \lambda v$$

and since by Proposition 13.3 $u + iv$ is an eigenvector of $f_{\mathbb{C}}^*$ for $\lambda - i\mu$, we have

$$f^*(u) = \lambda u + \mu v \quad \text{and} \quad f^*(v) = -\mu u + \lambda v,$$

and thus $f(W) = W$ and $f^*(W) = W$, where W is the subspace spanned by u and v .

When $\mu = 0$, we have

$$f(u) = \lambda u \quad \text{and} \quad f(v) = \lambda v,$$

and since $u \neq 0$ or $v \neq 0$, either u or v is an eigenvector of f for λ . If W is the subspace spanned by u if $u \neq 0$, or spanned by v if $u = 0$, it is obvious that $f(W) \subseteq W$ and $f^*(W) \subseteq W$. Note that $\lambda = 0$ is possible, and this is why \subseteq cannot be replaced by $=$. \square

The beginning of the proof of Proposition 13.9 actually shows that for every linear map $f: E \rightarrow E$ there is some subspace W such that $f(W) \subseteq W$, where W has dimension 1 or 2. In general, it doesn't seem possible to prove that W^\perp is invariant under f . However, this happens when f is normal.

We can finally prove our first main theorem.

Theorem 13.10. (*Main spectral theorem*) *Given a Euclidean space E of dimension n , for every normal linear map $f: E \rightarrow E$, there is an orthonormal basis (e_1, \dots, e_n) such that the matrix of f w.r.t. this basis is a block diagonal matrix of the form*

$$\begin{pmatrix} A_1 & & \dots & \\ & A_2 & & \\ \vdots & \vdots & \ddots & \vdots \\ & & \dots & A_p \end{pmatrix}$$

such that each block A_j is either a one-dimensional matrix (i.e., a real scalar) or a two-dimensional matrix of the form

$$A_j = \begin{pmatrix} \lambda_j & -\mu_j \\ \mu_j & \lambda_j \end{pmatrix},$$

where $\lambda_j, \mu_j \in \mathbb{R}$, with $\mu_j > 0$.

Proof. We proceed by induction on the dimension n of E as follows. If $n = 1$, the result is trivial. Assume now that $n \geq 2$. First, since \mathbb{C} is algebraically closed (i.e., every polynomial has a root in \mathbb{C}), the linear map $f_{\mathbb{C}}: E_{\mathbb{C}} \rightarrow E_{\mathbb{C}}$ has some eigenvalue $z = \lambda + i\mu$ (where $\lambda, \mu \in \mathbb{R}$). Let $w = u + iv$ be some eigenvector of $f_{\mathbb{C}}$ for $\lambda + i\mu$ (where $u, v \in E$). We can now apply Proposition 13.9.

If $\mu = 0$, then either u or v is an eigenvector of f for $\lambda \in \mathbb{R}$. Let W be the subspace of dimension 1 spanned by $e_1 = u/\|u\|$ if $u \neq 0$, or by $e_1 = v/\|v\|$ otherwise. It is obvious that $f(W) \subseteq W$ and $f^*(W) \subseteq W$. The orthogonal W^\perp of W has dimension $n - 1$, and by Proposition 13.8, we have $f(W^\perp) \subseteq W^\perp$. But the restriction of f to W^\perp is also normal, and we conclude by applying the induction hypothesis to W^\perp .

If $\mu \neq 0$, then $\langle u, v \rangle = 0$ and $\langle u, u \rangle = \langle v, v \rangle$, and if W is the subspace spanned by $u/\|u\|$ and $v/\|v\|$, then $f(W) = W$ and $f^*(W) = W$. We also know that the restriction of f to W has the matrix

$$\begin{pmatrix} \lambda & \mu \\ -\mu & \lambda \end{pmatrix}$$

with respect to the basis $(u/\|u\|, v/\|v\|)$. If $\mu < 0$, we let $\lambda_1 = \lambda$, $\mu_1 = -\mu$, $e_1 = u/\|u\|$, and $e_2 = v/\|v\|$. If $\mu > 0$, we let $\lambda_1 = \lambda$, $\mu_1 = \mu$, $e_1 = v/\|v\|$, and $e_2 = u/\|u\|$. In all cases, it is easily verified that the matrix of the restriction of f to W w.r.t. the orthonormal basis (e_1, e_2) is

$$A_1 = \begin{pmatrix} \lambda_1 & -\mu_1 \\ \mu_1 & \lambda_1 \end{pmatrix},$$

where $\lambda_1, \mu_1 \in \mathbb{R}$, with $\mu_1 > 0$. However, W^\perp has dimension $n - 2$, and by Proposition 13.8, $f(W^\perp) \subseteq W^\perp$. Since the restriction of f to W^\perp is also normal, we conclude by applying the induction hypothesis to W^\perp . \square

After this relatively hard work, we can easily obtain some nice normal forms for the matrices of self-adjoint, skew-self-adjoint, and orthogonal linear maps. However, for the sake of completeness (and since we have all the tools to so do), we go back to the case of a Hermitian space and show that normal linear maps can be diagonalized with respect to an orthonormal basis. The proof is a slight generalization of the proof of Theorem 13.6.

Theorem 13.11. (*Spectral theorem for normal linear maps on a Hermitian space*) *Given a Hermitian space E of dimension n , for every normal linear map $f: E \rightarrow E$ there is an orthonormal basis (e_1, \dots, e_n) of eigenvectors of f such that the matrix of f w.r.t. this basis is a diagonal matrix*

$$\begin{pmatrix} \lambda_1 & & \cdots & \\ & \lambda_2 & & \\ \vdots & \vdots & \ddots & \vdots \\ & & \cdots & \lambda_n \end{pmatrix},$$

where $\lambda_j \in \mathbb{C}$.

Proof. We proceed by induction on the dimension n of E as follows. If $n = 1$, the result is trivial. Assume now that $n \geq 2$. Since \mathbb{C} is algebraically closed (i.e., every polynomial has a root in \mathbb{C}), the linear map $f: E \rightarrow E$ has some eigenvalue $\lambda \in \mathbb{C}$, and let w be some unit eigenvector for λ . Let W be the subspace of dimension 1 spanned by w . Clearly, $f(W) \subseteq W$. By Proposition 13.3, w is an eigenvector of f^* for $\bar{\lambda}$, and thus $f^*(W) \subseteq W$. By Proposition 13.8, we also have $f(W^\perp) \subseteq W^\perp$. The restriction of f to W^\perp is still normal, and we conclude by applying the induction hypothesis to W^\perp (whose dimension is $n - 1$). \square

Thus, in particular, self-adjoint, skew-self-adjoint, and orthogonal linear maps can be diagonalized with respect to an orthonormal basis of eigenvectors. In this latter case, though, an orthogonal map is called a *unitary* map. Also, Proposition 13.5 shows that the eigenvalues of a self-adjoint linear map are real. It is easily shown that skew-self-adjoint maps have eigenvalues that are pure imaginary or null, and that unitary maps have eigenvalues of absolute value 1.

Remark: There is a converse to Theorem 13.11, namely, if there is an orthonormal basis (e_1, \dots, e_n) of eigenvectors of f , then f is normal. We leave the easy proof as an exercise.

13.3 Self-Adjoint, Skew-Self-Adjoint, and Orthogonal Linear Maps

We begin with self-adjoint maps.

Theorem 13.12. *Given a Euclidean space E of dimension n , for every self-adjoint linear map $f: E \rightarrow E$, there is an orthonormal basis (e_1, \dots, e_n) of eigenvectors of f such that the matrix of f w.r.t. this basis is a diagonal matrix*

$$\begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{pmatrix},$$

where $\lambda_i \in \mathbb{R}$.

Proof. We already proved this; see Theorem 13.6. However, it is instructive to give a more direct method not involving the complexification of $\langle -, - \rangle$ and Proposition 13.5.

Since \mathbb{C} is algebraically closed, $f_{\mathbb{C}}$ has some eigenvalue $\lambda + i\mu$, and let $u + iv$ be some eigenvector of $f_{\mathbb{C}}$ for $\lambda + i\mu$, where $\lambda, \mu \in \mathbb{R}$ and $u, v \in E$. We saw in the proof of Proposition 13.9 that

$$f(u) = \lambda u - \mu v \quad \text{and} \quad f(v) = \mu u + \lambda v.$$

Since $f = f^*$,

$$\langle f(u), v \rangle = \langle u, f(v) \rangle$$

for all $u, v \in E$. Applying this to

$$f(u) = \lambda u - \mu v \quad \text{and} \quad f(v) = \mu u + \lambda v,$$

we get

$$\langle f(u), v \rangle = \langle \lambda u - \mu v, v \rangle = \lambda \langle u, v \rangle - \mu \langle v, v \rangle$$

and

$$\langle u, f(v) \rangle = \langle u, \mu u + \lambda v \rangle = \mu \langle u, u \rangle + \lambda \langle u, v \rangle,$$

and thus we get

$$\lambda \langle u, v \rangle - \mu \langle v, v \rangle = \mu \langle u, u \rangle + \lambda \langle u, v \rangle,$$

that is,

$$\mu(\langle u, u \rangle + \langle v, v \rangle) = 0,$$

which implies $\mu = 0$, since either $u \neq 0$ or $v \neq 0$. Therefore, λ is a real eigenvalue of f .

Now, going back to the proof of Theorem 13.10, only the case where $\mu = 0$ applies, and the induction shows that all the blocks are one-dimensional. \square

Theorem 13.12 implies that if $\lambda_1, \dots, \lambda_p$ are the distinct real eigenvalues of f , and E_i is the eigenspace associated with λ_i , then

$$E = E_1 \oplus \dots \oplus E_p,$$

where E_i and E_j are orthogonal for all $i \neq j$.

Remark: Another way to prove that a self-adjoint map has a real eigenvalue is to use a little bit of calculus. We learned such a proof from Herman Gluck. The idea is to consider the real-valued function $\Phi: E \rightarrow \mathbb{R}$ defined such that

$$\Phi(u) = \langle f(u), u \rangle$$

for every $u \in E$. This function is C^∞ , and if we represent f by a matrix A over some orthonormal basis, it is easy to compute the gradient vector

$$\nabla \Phi(X) = \left(\frac{\partial \Phi}{\partial x_1}(X), \dots, \frac{\partial \Phi}{\partial x_n}(X) \right)$$

of Φ at X . Indeed, we find that

$$\nabla \Phi(X) = (A + A^\top)X,$$

where X is a column vector of size n . But since f is self-adjoint, $A = A^\top$, and thus

$$\nabla \Phi(X) = 2AX.$$

The next step is to find the maximum of the function Φ on the sphere

$$S^{n-1} = \{(x_1, \dots, x_n) \in \mathbb{R}^n \mid x_1^2 + \dots + x_n^2 = 1\}.$$

Since S^{n-1} is compact and Φ is continuous, and in fact C^∞ , Φ takes a maximum at some X on S^{n-1} . But then it is well known that at an extremum X of Φ we must have

$$d\Phi_X(Y) = \langle \nabla \Phi(X), Y \rangle = 0$$

for all tangent vectors Y to S^{n-1} at X , and so $\nabla \Phi(X)$ is orthogonal to the tangent plane at X , which means that

$$\nabla \Phi(X) = \lambda X$$

for some $\lambda \in \mathbb{R}$. Since $\nabla \Phi(X) = 2AX$, we get

$$2AX = \lambda X,$$

and thus $\lambda/2$ is a real eigenvalue of A (i.e., of f).

Next, we consider skew-self-adjoint maps.

Theorem 13.13. *Given a Euclidean space E of dimension n , for every skew-self-adjoint linear map $f: E \rightarrow E$ there is an orthonormal basis (e_1, \dots, e_n) such that the matrix of f w.r.t. this basis is a block diagonal matrix of the form*

$$\begin{pmatrix} A_1 & & \cdots & \\ & A_2 & & \\ \vdots & \vdots & \ddots & \vdots \\ & & \cdots & A_p \end{pmatrix}$$

such that each block A_j is either 0 or a two-dimensional matrix of the form

$$A_j = \begin{pmatrix} 0 & -\mu_j \\ \mu_j & 0 \end{pmatrix},$$

where $\mu_j \in \mathbb{R}$, with $\mu_j > 0$. In particular, the eigenvalues of $f_{\mathbb{C}}$ are pure imaginary of the form $\pm i\mu_j$ or 0.

Proof. The case where $n = 1$ is trivial. As in the proof of Theorem 13.10, $f_{\mathbb{C}}$ has some eigenvalue $z = \lambda + i\mu$, where $\lambda, \mu \in \mathbb{R}$. We claim that $\lambda = 0$. First, we show that

$$\langle f(w), w \rangle = 0$$

for all $w \in E$. Indeed, since $f = -f^*$, we get

$$\langle f(w), w \rangle = \langle w, f^*(w) \rangle = \langle w, -f(w) \rangle = -\langle w, f(w) \rangle = -\langle f(w), w \rangle,$$

since $\langle -, - \rangle$ is symmetric. This implies that

$$\langle f(w), w \rangle = 0.$$

Applying this to u and v and using the fact that

$$f(u) = \lambda u - \mu v \quad \text{and} \quad f(v) = \mu u + \lambda v,$$

we get

$$0 = \langle f(u), u \rangle = \langle \lambda u - \mu v, u \rangle = \lambda \langle u, u \rangle - \mu \langle u, v \rangle$$

and

$$0 = \langle f(v), v \rangle = \langle \mu u + \lambda v, v \rangle = \mu \langle u, v \rangle + \lambda \langle v, v \rangle,$$

from which, by addition, we get

$$\lambda(\langle v, v \rangle + \langle v, v \rangle) = 0.$$

Since $u \neq 0$ or $v \neq 0$, we have $\lambda = 0$.

Then, going back to the proof of Theorem 13.10, unless $\mu = 0$, the case where u and v are orthogonal and span a subspace of dimension 2 applies, and the induction shows that all the blocks are two-dimensional or reduced to 0. \square

Remark: One will note that if f is skew-self-adjoint, then $if_{\mathbb{C}}$ is self-adjoint w.r.t. $\langle -, - \rangle_{\mathbb{C}}$. By Proposition 13.5, the map $if_{\mathbb{C}}$ has real eigenvalues, which implies that the eigenvalues of $f_{\mathbb{C}}$ are pure imaginary or 0.

Finally, we consider orthogonal linear maps.

Theorem 13.14. *Given a Euclidean space E of dimension n , for every orthogonal linear map $f: E \rightarrow E$ there is an orthonormal basis (e_1, \dots, e_n) such that the matrix of f w.r.t. this basis is a block diagonal matrix of the form*

$$\begin{pmatrix} A_1 & & \dots & \\ & A_2 & & \\ \vdots & \vdots & \ddots & \vdots \\ & & \dots & A_p \end{pmatrix}$$

such that each block A_j is either 1, -1 , or a two-dimensional matrix of the form

$$A_j = \begin{pmatrix} \cos \theta_j & -\sin \theta_j \\ \sin \theta_j & \cos \theta_j \end{pmatrix}$$

where $0 < \theta_j < \pi$. In particular, the eigenvalues of $f_{\mathbb{C}}$ are of the form $\cos \theta_j \pm i \sin \theta_j$, 1, or -1 .

Proof. The case where $n = 1$ is trivial. As in the proof of Theorem 13.10, $f_{\mathbb{C}}$ has some eigenvalue $z = \lambda + i\mu$, where $\lambda, \mu \in \mathbb{R}$. It is immediately verified that $f \circ f^* = f^* \circ f = \text{id}$ implies that $f_{\mathbb{C}} \circ f_{\mathbb{C}}^* = f_{\mathbb{C}}^* \circ f_{\mathbb{C}} = \text{id}$, so the map $f_{\mathbb{C}}$ is unitary. In fact, the eigenvalues of $f_{\mathbb{C}}$ have absolute value 1. Indeed, if z (in \mathbb{C}) is an eigenvalue of $f_{\mathbb{C}}$, and u is an eigenvector for z , we have

$$\langle f_{\mathbb{C}}(u), f_{\mathbb{C}}(u) \rangle = \langle zu, zu \rangle = z\bar{z}\langle u, u \rangle$$

and

$$\langle f_{\mathbb{C}}(u), f_{\mathbb{C}}(u) \rangle = \langle u, (f_{\mathbb{C}}^* \circ f_{\mathbb{C}})(u) \rangle = \langle u, u \rangle,$$

from which we get

$$z\bar{z}\langle u, u \rangle = \langle u, u \rangle.$$

Since $u \neq 0$, we have $z\bar{z} = 1$, i.e., $|z| = 1$. As a consequence, the eigenvalues of $f_{\mathbb{C}}$ are of the form $\cos \theta \pm i \sin \theta$, 1, or -1 . The theorem then follows immediately from Theorem 13.10, where the condition $\mu > 0$ implies that $\sin \theta_j > 0$, and thus, $0 < \theta_j < \pi$. \square

It is obvious that we can reorder the orthonormal basis of eigenvectors given by Theorem 13.14, so that the matrix of f w.r.t. this basis is a block diagonal matrix of the form

$$\begin{pmatrix} A_1 & \dots & & & \\ \vdots & \ddots & \vdots & & \vdots \\ & \dots & A_r & & \\ & & & -I_q & \\ \dots & & & & I_p \end{pmatrix}$$

where each block A_j is a two-dimensional rotation matrix $A_j \neq \pm I_2$ of the form

$$A_j = \begin{pmatrix} \cos \theta_j & -\sin \theta_j \\ \sin \theta_j & \cos \theta_j \end{pmatrix}$$

with $0 < \theta_j < \pi$.

The linear map f has an eigenspace $E(1, f) = \text{Ker}(f - \text{id})$ of dimension p for the eigenvalue 1, and an eigenspace $E(-1, f) = \text{Ker}(f + \text{id})$ of dimension q for the eigenvalue -1 . If $\det(f) = +1$ (f is a rotation), the dimension q of $E(-1, f)$ must be even, and the entries in $-I_q$ can be paired to form two-dimensional blocks, if we wish. In this case, every rotation in $\mathbf{SO}(n)$ has a matrix of the form

$$\begin{pmatrix} A_1 & \cdots & & \\ \vdots & \ddots & \vdots & \\ & \cdots & A_m & \\ \cdots & & & I_{n-2m} \end{pmatrix}$$

where the first m blocks A_j are of the form

$$A_j = \begin{pmatrix} \cos \theta_j & -\sin \theta_j \\ \sin \theta_j & \cos \theta_j \end{pmatrix}$$

with $0 < \theta_j \leq \pi$.

Theorem 13.14 can be used to prove a version of the Cartan–Dieudonné theorem.

Theorem 13.15. *Let E be a Euclidean space of dimension $n \geq 2$. For every isometry $f \in \mathbf{O}(E)$, if $p = \dim(E(1, f)) = \dim(\text{Ker}(f - \text{id}))$, then f is the composition of $n - p$ reflections, and $n - p$ is minimal.*

Proof. From Theorem 13.14 there are r subspaces F_1, \dots, F_r , each of dimension 2, such that

$$E = E(1, f) \oplus E(-1, f) \oplus F_1 \oplus \cdots \oplus F_r,$$

and all the summands are pairwise orthogonal. Furthermore, the restriction r_i of f to each F_i is a rotation $r_i \neq \pm \text{id}$. Each 2D rotation r_i can be written as the composition $r_i = s'_i \circ s_i$ of two reflections s_i and s'_i about lines in F_i (forming an angle $\theta_i/2$). We can extend s_i and s'_i to hyperplane reflections in E by making them the identity on F_i^\perp . Then,

$$s'_r \circ s_r \circ \cdots \circ s'_1 \circ s_1$$

agrees with f on $F_1 \oplus \cdots \oplus F_r$ and is the identity on $E(1, f) \oplus E(-1, f)$. If $E(-1, f)$ has an orthonormal basis of eigenvectors (v_1, \dots, v_q) , letting s''_j be the reflection about the hyperplane $(v_j)^\perp$, it is clear that

$$s''_q \circ \cdots \circ s''_1$$

agrees with f on $E(-1, f)$ and is the identity on $E(1, f) \oplus F_1 \oplus \cdots \oplus F_r$. But then,

$$f = s_q'' \circ \cdots \circ s_1'' \circ s_r' \circ s_r \circ \cdots \circ s_1' \circ s_1,$$

the composition of $2r + q = n - p$ reflections.

If

$$f = s_t \circ \cdots \circ s_1,$$

for t reflections s_i , it is clear that

$$F = \bigcap_{i=1}^t E(1, s_i) \subseteq E(1, f),$$

where $E(1, s_i)$ is the hyperplane defining the reflection s_i . By the Grassmann relation, if we intersect $t \leq n$ hyperplanes, the dimension of their intersection is at least $n - t$. Thus, $n - t \leq p$, that is, $t \geq n - p$, and $n - p$ is the smallest number of reflections composing f . \square

As a corollary of Theorem 13.15, we obtain the following fact: If the dimension n of the Euclidean space E is odd, then every rotation $f \in \mathbf{SO}(E)$ admits 1 as an eigenvalue.

Proof. The characteristic polynomial $\det(XI - f)$ of f has odd degree n and has real coefficients, so it must have some real root λ . Since f is an isometry, its n eigenvalues are of the form, $+1$, -1 , and $e^{\pm i\theta}$, with $0 < \theta < \pi$, so $\lambda = \pm 1$. Now, the eigenvalues $e^{\pm i\theta}$ appear in conjugate pairs, and since n is odd, the number of real eigenvalues of f is odd. This implies that $+1$ is an eigenvalue of f , since otherwise -1 would be the only real eigenvalue of f , and since its multiplicity is odd, we would have $\det(f) = -1$, contradicting the fact that f is a rotation. \square

When $n = 3$, we obtain the result due to Euler which says that every 3D rotation R has an invariant axis D , and that restricted to the plane orthogonal to D , it is a 2D rotation. Furthermore, if (a, b, c) is a unit vector defining the axis D of the rotation R and if the angle of the rotation is θ , if B is the skew-symmetric matrix

$$B = \begin{pmatrix} 0 & -c & b \\ c & 0 & -a \\ -b & a & 0 \end{pmatrix},$$

then it can be shown that

$$R = I + \sin \theta B + (1 - \cos \theta) B^2.$$

The theorems of this section and of the previous section can be immediately applied to matrices.

13.4 Normal and Other Special Matrices

First, we consider real matrices. Recall the following definitions.

Definition 13.3. Given a real $m \times n$ matrix A , the *transpose* A^\top of A is the $n \times m$ matrix $A^\top = (a_{ij}^\top)$ defined such that

$$a_{ij}^\top = a_{ji}$$

for all i, j , $1 \leq i \leq m$, $1 \leq j \leq n$. A real $n \times n$ matrix A is

- *normal* if

$$A A^\top = A^\top A,$$

- *symmetric* if

$$A^\top = A,$$

- *skew-symmetric* if

$$A^\top = -A,$$

- *orthogonal* if

$$A A^\top = A^\top A = I_n.$$

Recall from Proposition 10.12 that when E is a Euclidean space and (e_1, \dots, e_n) is an orthonormal basis for E , if A is the matrix of a linear map $f: E \rightarrow E$ w.r.t. the basis (e_1, \dots, e_n) , then A^\top is the matrix of the adjoint f^* of f . Consequently, a normal linear map has a normal matrix, a self-adjoint linear map has a symmetric matrix, a skew-self-adjoint linear map has a skew-symmetric matrix, and an orthogonal linear map has an orthogonal matrix. Similarly, if E and F are Euclidean spaces, (u_1, \dots, u_n) is an orthonormal basis for E , and (v_1, \dots, v_m) is an orthonormal basis for F , if a linear map $f: E \rightarrow F$ has the matrix A w.r.t. the bases (u_1, \dots, u_n) and (v_1, \dots, v_m) , then its adjoint f^* has the matrix A^\top w.r.t. the bases (v_1, \dots, v_m) and (u_1, \dots, u_n) .

Furthermore, if (u_1, \dots, u_n) is another orthonormal basis for E and P is the change of basis matrix whose columns are the components of the u_i w.r.t. the basis (e_1, \dots, e_n) , then P is orthogonal, and for any linear map $f: E \rightarrow E$, if A is the matrix of f w.r.t (e_1, \dots, e_n) and B is the matrix of f w.r.t. (u_1, \dots, u_n) , then

$$B = P^\top A P.$$

As a consequence, Theorems 13.10 and 13.12–13.14 can be restated as follows.

Theorem 13.16. *For every normal matrix A there is an orthogonal matrix P and a block diagonal matrix D such that $A = P D P^\top$, where D is of the form*

$$D = \begin{pmatrix} D_1 & & \cdots & \\ & D_2 & & \\ \vdots & \vdots & \ddots & \vdots \\ & & \cdots & D_p \end{pmatrix}$$

such that each block D_j is either a one-dimensional matrix (i.e., a real scalar) or a two-dimensional matrix of the form

$$D_j = \begin{pmatrix} \lambda_j & -\mu_j \\ \mu_j & \lambda_j \end{pmatrix},$$

where $\lambda_j, \mu_j \in \mathbb{R}$, with $\mu_j > 0$.

Theorem 13.17. *For every symmetric matrix A there is an orthogonal matrix P and a diagonal matrix D such that $A = P D P^\top$, where D is of the form*

$$D = \begin{pmatrix} \lambda_1 & & \cdots & \\ & \lambda_2 & & \\ \vdots & \vdots & \ddots & \vdots \\ & & \cdots & \lambda_n \end{pmatrix},$$

where $\lambda_i \in \mathbb{R}$.

Theorem 13.18. *For every skew-symmetric matrix A there is an orthogonal matrix P and a block diagonal matrix D such that $A = P D P^\top$, where D is of the form*

$$D = \begin{pmatrix} D_1 & & \cdots & \\ & D_2 & & \\ \vdots & \vdots & \ddots & \vdots \\ & & \cdots & D_p \end{pmatrix}$$

such that each block D_j is either 0 or a two-dimensional matrix of the form

$$D_j = \begin{pmatrix} 0 & -\mu_j \\ \mu_j & 0 \end{pmatrix},$$

where $\mu_j \in \mathbb{R}$, with $\mu_j > 0$. In particular, the eigenvalues of A are pure imaginary of the form $\pm i\mu_j$, or 0.

Theorem 13.19. *For every orthogonal matrix A there is an orthogonal matrix P and a block diagonal matrix D such that $A = P D P^\top$, where D is of the form*

$$D = \begin{pmatrix} D_1 & & \cdots & \\ & D_2 & & \\ \vdots & \vdots & \ddots & \vdots \\ & & \cdots & D_p \end{pmatrix}$$

such that each block D_j is either 1, -1 , or a two-dimensional matrix of the form

$$D_j = \begin{pmatrix} \cos \theta_j & -\sin \theta_j \\ \sin \theta_j & \cos \theta_j \end{pmatrix}$$

where $0 < \theta_j < \pi$. In particular, the eigenvalues of A are of the form $\cos \theta_j \pm i \sin \theta_j$, 1, or -1 .

We now consider complex matrices.

Definition 13.4. Given a complex $m \times n$ matrix A , the *transpose* A^\top of A is the $n \times m$ matrix $A^\top = (a_{ij}^\top)$ defined such that

$$a_{ij}^\top = a_{ji}$$

for all i, j , $1 \leq i \leq m$, $1 \leq j \leq n$. The *conjugate* \bar{A} of A is the $m \times n$ matrix $\bar{A} = (b_{ij})$ defined such that

$$b_{ij} = \bar{a}_{ij}$$

for all i, j , $1 \leq i \leq m$, $1 \leq j \leq n$. Given an $m \times n$ complex matrix A , the *adjoint* A^* of A is the matrix defined such that

$$A^* = \overline{(A^\top)} = (\bar{A})^\top.$$

A complex $n \times n$ matrix A is

- *normal* if

$$AA^* = A^*A,$$

- *Hermitian* if

$$A^* = A,$$

- *skew-Hermitian* if

$$A^* = -A,$$

- *unitary* if

$$AA^* = A^*A = I_n.$$

Recall from Proposition 12.12 that when E is a Hermitian space and (e_1, \dots, e_n) is an orthonormal basis for E , if A is the matrix of a linear map $f: E \rightarrow E$ w.r.t. the basis (e_1, \dots, e_n) , then A^* is the matrix of the adjoint f^* of f . Consequently, a normal linear map has a normal matrix, a self-adjoint linear map has a Hermitian matrix, a skew-self-adjoint linear map has a skew-Hermitian matrix, and a unitary linear map has a unitary matrix.

Similarly, if E and F are Hermitian spaces, (u_1, \dots, u_n) is an orthonormal basis for E , and (v_1, \dots, v_m) is an orthonormal basis for F , if a linear map $f: E \rightarrow F$ has the matrix A w.r.t. the bases (u_1, \dots, u_n) and (v_1, \dots, v_m) , then its adjoint f^* has the matrix A^* w.r.t. the bases (v_1, \dots, v_m) and (u_1, \dots, u_n) .

Furthermore, if (u_1, \dots, u_n) is another orthonormal basis for E and P is the change of basis matrix whose columns are the components of the u_i w.r.t. the basis (e_1, \dots, e_n) , then P is unitary, and for any linear map $f: E \rightarrow E$, if A is the matrix of f w.r.t (e_1, \dots, e_n) and B is the matrix of f w.r.t. (u_1, \dots, u_n) , then

$$B = P^*AP.$$

Theorem 13.11 can be restated in terms of matrices as follows. We can also say a little more about eigenvalues (easy exercise left to the reader).

Theorem 13.20. *For every complex normal matrix A there is a unitary matrix U and a diagonal matrix D such that $A = UDU^*$. Furthermore, if A is Hermitian, then D is a real matrix; if A is skew-Hermitian, then the entries in D are pure imaginary or null; and if A is unitary, then the entries in D have absolute value 1.*

13.5 Conditioning of Eigenvalue Problems

The following $n \times n$ matrix

$$A = \begin{pmatrix} 0 & & & & \\ 1 & 0 & & & \\ & 1 & 0 & & \\ & & \ddots & \ddots & \\ & & & 1 & 0 \\ & & & & 1 & 0 \end{pmatrix}$$

has the eigenvalue 0 with multiplicity n . However, if we perturb the top rightmost entry of A by ϵ , it is easy to see that the characteristic polynomial of the matrix

$$A(\epsilon) = \begin{pmatrix} 0 & & & & \epsilon \\ 1 & 0 & & & \\ & 1 & 0 & & \\ & & \ddots & \ddots & \\ & & & 1 & 0 \\ & & & & 1 & 0 \end{pmatrix}$$

is $X^n - \epsilon$. It follows that if $n = 40$ and $\epsilon = 10^{-40}$, $A(10^{-40})$ has the eigenvalues $e^{k2\pi i/40}10^{-1}$ with $k = 1, \dots, 40$. Thus, we see that a very small change ($\epsilon = 10^{-40}$) to the matrix A causes

a significant change to the eigenvalues of A (from 0 to $e^{k2\pi i/40}10^{-1}$). Indeed, the relative error is 10^{-39} . Worse, due to machine precision, since very small numbers are treated as 0, the error on the computation of eigenvalues (for example, of the matrix $A(10^{-40})$) can be very large.

This phenomenon is similar to the phenomenon discussed in Section 7.3 where we studied the effect of a small perturbation of the coefficients of a linear system $Ax = b$ on its solution. In Section 7.3, we saw that the behavior of a linear system under small perturbations is governed by the condition number $\text{cond}(A)$ of the matrix A . In the case of the eigenvalue problem (finding the eigenvalues of a matrix), we will see that the conditioning of the problem depends on the condition number of the change of basis matrix P used in reducing the matrix A to its diagonal form $D = P^{-1}AP$, rather than on the condition number of A itself. The following proposition in which we assume that A is diagonalizable and that the matrix norm $\|\cdot\|$ satisfies a special condition (satisfied by the operator norms $\|\cdot\|_p$ for $p = 1, 2, \infty$), is due to Bauer and Fike (1960).

Proposition 13.21. *Let $A \in M_n(\mathbb{C})$ be a diagonalizable matrix, P be an invertible matrix and, D be a diagonal matrix $D = \text{diag}(\lambda_1, \dots, \lambda_n)$ such that*

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

and let $\|\cdot\|$ be a matrix norm such that

$$\|\text{diag}(\alpha_1, \dots, \alpha_n)\| = \max_{1 \leq i \leq n} |\alpha_i|,$$

for every diagonal matrix. Then, for every perturbation matrix δA , if we write

$$B_i = \{z \in \mathbb{C} \mid |z - \lambda_i| \leq \text{cond}(P) \|\delta A\|\},$$

for every eigenvalue λ of $A + \delta A$, we have

$$\lambda \in \bigcup_{k=1}^n B_k.$$

Proof. Let λ be any eigenvalue of the matrix $A + \delta A$. If $\lambda = \lambda_j$ for some j , then the result is trivial. Thus, assume that $\lambda \neq \lambda_j$ for $j = 1, \dots, n$. In this case, the matrix $D - \lambda I$ is invertible (since its eigenvalues are $\lambda - \lambda_j$ for $j = 1, \dots, n$), and we have

$$\begin{aligned} P^{-1}(A + \delta A - \lambda I)P &= D - \lambda I + P^{-1}(\delta A)P \\ &= (D - \lambda I)(I + (D - \lambda I)^{-1}P^{-1}(\delta A)P). \end{aligned}$$

Since λ is an eigenvalue of $A + \delta A$, the matrix $A + \delta A - \lambda I$ is singular, so the matrix

$$I + (D - \lambda I)^{-1}P^{-1}(\delta A)P$$

must also be singular. By Proposition 7.9(2), we have

$$1 \leq \|(D - \lambda I)^{-1} P^{-1} (\delta A) P\|,$$

and since $\|\cdot\|$ is a matrix norm,

$$\|(D - \lambda I)^{-1} P^{-1} (\delta A) P\| \leq \|(D - \lambda I)^{-1}\| \|P^{-1}\| \|\delta A\| \|P\|,$$

so we have

$$1 \leq \|(D - \lambda I)^{-1}\| \|P^{-1}\| \|\delta A\| \|P\|.$$

Now, $(D - \lambda I)^{-1}$ is a diagonal matrix with entries $1/(\lambda_i - \lambda)$, so by our assumption on the norm,

$$\|(D - \lambda I)^{-1}\| = \frac{1}{\min_i (|\lambda_i - \lambda|)}.$$

As a consequence, since there is some index k for which $\min_i (|\lambda_i - \lambda|) = |\lambda_k - \lambda|$, we have

$$\|(D - \lambda I)^{-1}\| = \frac{1}{|\lambda_k - \lambda|},$$

and we obtain

$$|\lambda - \lambda_k| \leq \|P^{-1}\| \|\delta A\| \|P\| = \text{cond}(P) \|\delta A\|,$$

which proves our result. \square

Proposition 13.21 implies that for any diagonalizable matrix A , if we define $\Gamma(A)$ by

$$\Gamma(A) = \inf\{\text{cond}(P) \mid P^{-1}AP = D\},$$

then for every eigenvalue λ of $A + \delta A$, we have

$$\lambda \in \bigcup_{k=1}^n \{z \in \mathbb{C}^n \mid |z - \lambda_k| \leq \Gamma(A) \|\delta A\|\}.$$

The number $\Gamma(A)$ is called the *conditioning of A relative to the eigenvalue problem*. If A is a normal matrix, since by Theorem 13.20, A can be diagonalized with respect to a unitary matrix U , and since for the spectral norm $\|U\|_2 = 1$, we see that $\Gamma(A) = 1$. Therefore, normal matrices are very well conditioned w.r.t. the eigenvalue problem. In fact, for every eigenvalue λ of $A + \delta A$ (with A normal), we have

$$\lambda \in \bigcup_{k=1}^n \{z \in \mathbb{C}^n \mid |z - \lambda_k| \leq \|\delta A\|_2\}.$$

If A and $A + \delta A$ are both symmetric (or Hermitian), there are sharper results; see Proposition 13.27.

Note that the matrix $A(\epsilon)$ from the beginning of the section is not normal.

13.6 Rayleigh Ratios and the Courant-Fischer Theorem

A fact that is used frequently in optimization problem is that the eigenvalues of a symmetric matrix are characterized in terms of what is known as the *Rayleigh ratio*, defined by

$$R(A)(x) = \frac{x^\top Ax}{x^\top x}, \quad x \in \mathbb{R}^n, x \neq 0.$$

The following proposition is often used to prove the correctness of various optimization or approximation problems (for example PCA).

Proposition 13.22. (*Rayleigh–Ritz*) *If A is a symmetric $n \times n$ matrix with eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ and if (u_1, \dots, u_n) is any orthonormal basis of eigenvectors of A , where u_i is a unit eigenvector associated with λ_i , then*

$$\max_{x \neq 0} \frac{x^\top Ax}{x^\top x} = \lambda_n$$

(with the maximum attained for $x = u_n$), and

$$\max_{x \neq 0, x \in \{u_{n-k+1}, \dots, u_n\}^\perp} \frac{x^\top Ax}{x^\top x} = \lambda_{n-k}$$

(with the maximum attained for $x = u_{n-k}$), where $1 \leq k \leq n-1$. Equivalently, if V_k is the subspace spanned by (u_1, \dots, u_k) , then

$$\lambda_k = \max_{x \neq 0, x \in V_k} \frac{x^\top Ax}{x^\top x}, \quad k = 1, \dots, n.$$

Proof. First, observe that

$$\max_{x \neq 0} \frac{x^\top Ax}{x^\top x} = \max_x \{x^\top Ax \mid x^\top x = 1\},$$

and similarly,

$$\max_{x \neq 0, x \in \{u_{n-k+1}, \dots, u_n\}^\perp} \frac{x^\top Ax}{x^\top x} = \max_x \{x^\top Ax \mid (x \in \{u_{n-k+1}, \dots, u_n\}^\perp) \wedge (x^\top x = 1)\}.$$

Since A is a symmetric matrix, its eigenvalues are real and it can be diagonalized with respect to an orthonormal basis of eigenvectors, so let (u_1, \dots, u_n) be such a basis. If we write

$$x = \sum_{i=1}^n x_i u_i,$$

a simple computation shows that

$$x^\top Ax = \sum_{i=1}^n \lambda_i x_i^2.$$

If $x^\top x = 1$, then $\sum_{i=1}^n x_i^2 = 1$, and since we assumed that $\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_n$, we get

$$x^\top Ax = \sum_{i=1}^n \lambda_i x_i^2 \leq \lambda_n \left(\sum_{i=1}^n x_i^2 \right) = \lambda_n.$$

Thus,

$$\max_x \{x^\top Ax \mid x^\top x = 1\} \leq \lambda_n,$$

and since this maximum is achieved for $e_n = (0, 0, \dots, 1)$, we conclude that

$$\max_x \{x^\top Ax \mid x^\top x = 1\} = \lambda_n.$$

Next, observe that $x \in \{u_{n-k+1}, \dots, u_n\}^\perp$ and $x^\top x = 1$ iff $x_{n-k+1} = \cdots = x_n = 0$ and $\sum_{i=1}^{n-k} x_i^2 = 1$. Consequently, for such an x , we have

$$x^\top Ax = \sum_{i=1}^{n-k} \lambda_i x_i^2 \leq \lambda_{n-k} \left(\sum_{i=1}^{n-k} x_i^2 \right) = \lambda_{n-k}.$$

Thus,

$$\max_x \{x^\top Ax \mid (x \in \{u_{n-k+1}, \dots, u_n\}^\perp) \wedge (x^\top x = 1)\} \leq \lambda_{n-k},$$

and since this maximum is achieved for $e_{n-k} = (0, \dots, 0, 1, 0, \dots, 0)$ with a 1 in position $n-k$, we conclude that

$$\max_x \{x^\top Ax \mid (x \in \{u_{n-k+1}, \dots, u_n\}^\perp) \wedge (x^\top x = 1)\} = \lambda_{n-k},$$

as claimed. □

For our purposes, we need the version of Proposition 13.22 applying to min instead of max, whose proof is obtained by a trivial modification of the proof of Proposition 13.22.

Proposition 13.23. (*Rayleigh–Ritz*) *If A is a symmetric $n \times n$ matrix with eigenvalues $\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_n$ and if (u_1, \dots, u_n) is any orthonormal basis of eigenvectors of A , where u_i is a unit eigenvector associated with λ_i , then*

$$\min_{x \neq 0} \frac{x^\top Ax}{x^\top x} = \lambda_1$$

(with the minimum attained for $x = u_1$), and

$$\min_{x \neq 0, x \in \{u_1, \dots, u_{i-1}\}^\perp} \frac{x^\top Ax}{x^\top x} = \lambda_i$$

(with the minimum attained for $x = u_i$), where $2 \leq i \leq n$. Equivalently, if $W_k = V_{k-1}^\perp$ denotes the subspace spanned by (u_k, \dots, u_n) (with $V_0 = (0)$), then

$$\lambda_k = \min_{x \neq 0, x \in W_k} \frac{x^\top A x}{x^\top x} = \min_{x \neq 0, x \in V_{k-1}^\perp} \frac{x^\top A x}{x^\top x}, \quad k = 1, \dots, n.$$

Propositions 13.22 and 13.23 together are known the *Rayleigh–Ritz theorem*.

As an application of Propositions 13.22 and 13.23, we prove a proposition which allows us to compare the eigenvalues of two symmetric matrices A and $B = R^\top A R$, where R is a rectangular matrix satisfying the equation $R^\top R = I$.

First, we need a definition. Given an $n \times n$ symmetric matrix A and an $m \times m$ symmetric B , with $m \leq n$, if $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ are the eigenvalues of A and $\mu_1 \leq \mu_2 \leq \dots \leq \mu_m$ are the eigenvalues of B , then we say that the eigenvalues of B *interlace* the eigenvalues of A if

$$\lambda_i \leq \mu_i \leq \lambda_{n-m+i}, \quad i = 1, \dots, m.$$

Proposition 13.24. *Let A be an $n \times n$ symmetric matrix, R be an $n \times m$ matrix such that $R^\top R = I$ (with $m \leq n$), and let $B = R^\top A R$ (an $m \times m$ matrix). The following properties hold:*

- (a) *The eigenvalues of B interlace the eigenvalues of A .*
- (b) *If $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ are the eigenvalues of A and $\mu_1 \leq \mu_2 \leq \dots \leq \mu_m$ are the eigenvalues of B , and if $\lambda_i = \mu_i$, then there is an eigenvector v of B with eigenvalue μ_i such that Rv is an eigenvector of A with eigenvalue λ_i .*

Proof. (a) Let (u_1, \dots, u_n) be an orthonormal basis of eigenvectors for A , and let (v_1, \dots, v_m) be an orthonormal basis of eigenvectors for B . Let U_j be the subspace spanned by (u_1, \dots, u_j) and let V_j be the subspace spanned by (v_1, \dots, v_j) . For any i , the subspace V_i has dimension i and the subspace $R^\top U_{i-1}$ has dimension at most $i - 1$. Therefore, there is some nonzero vector $v \in V_i \cap (R^\top U_{i-1})^\perp$, and since

$$v^\top R^\top u_j = (Rv)^\top u_j = 0, \quad j = 1, \dots, i - 1,$$

we have $Rv \in (U_{i-1})^\perp$. By Proposition 13.23 and using the fact that $R^\top R = I$, we have

$$\lambda_i \leq \frac{(Rv)^\top A Rv}{(Rv)^\top Rv} = \frac{v^\top B v}{v^\top v}.$$

On the other hand, by Proposition 13.22,

$$\mu_i = \max_{x \neq 0, x \in \{v_{i+1}, \dots, v_m\}^\perp} \frac{x^\top B x}{x^\top x} = \max_{x \neq 0, x \in \{v_1, \dots, v_i\}} \frac{x^\top B x}{x^\top x},$$

so

$$\frac{w^\top Bw}{w^\top w} \leq \mu_i \quad \text{for all } w \in V_i,$$

and since $v \in V_i$, we have

$$\lambda_i \leq \frac{v^\top Bv}{v^\top v} \leq \mu_i, \quad i = 1, \dots, m.$$

We can apply the same argument to the symmetric matrices $-A$ and $-B$, to conclude that

$$-\lambda_{n-m+i} \leq -\mu_i,$$

that is,

$$\mu_i \leq \lambda_{n-m+i}, \quad i = 1, \dots, m.$$

Therefore,

$$\lambda_i \leq \mu_i \leq \lambda_{n-m+i}, \quad i = 1, \dots, m,$$

as desired.

(b) If $\lambda_i = \mu_i$, then

$$\lambda_i = \frac{(Rv)^\top ARv}{(Rv)^\top Rv} = \frac{v^\top Bv}{v^\top v} = \mu_i,$$

so v must be an eigenvector for B and Rv must be an eigenvector for A , both for the eigenvalue $\lambda_i = \mu_i$. \square

Proposition 13.24 immediately implies the *Poincaré separation theorem*. It can be used in situations, such as in quantum mechanics, where one has information about the inner products $u_i^\top Au_j$.

Proposition 13.25. (*Poincaré separation theorem*) *Let A be a $n \times n$ symmetric (or Hermitian) matrix, let r be some integer with $1 \leq r \leq n$, and let (u_1, \dots, u_r) be r orthonormal vectors. Let $B = (u_i^\top Au_j)$ (an $r \times r$ matrix), let $\lambda_1(A) \leq \dots \leq \lambda_n(A)$ be the eigenvalues of A and $\lambda_1(B) \leq \dots \leq \lambda_r(B)$ be the eigenvalues of B ; then we have*

$$\lambda_k(A) \leq \lambda_k(B) \leq \lambda_{k+n-r}(A), \quad k = 1, \dots, r.$$

Observe that Proposition 13.24 implies that

$$\lambda_1 + \dots + \lambda_m \leq \text{tr}(R^\top AR) \leq \lambda_{n-m+1} + \dots + \lambda_n.$$

If P_1 is the $n \times (n-1)$ matrix obtained from the identity matrix by dropping its last column, we have $P_1^\top P_1 = I$, and the matrix $B = P_1^\top AP_1$ is the matrix obtained from A by deleting its last row and its last column. In this case, the interlacing result is

$$\lambda_1 \leq \mu_1 \leq \lambda_2 \leq \mu_2 \leq \dots \leq \mu_{n-2} \leq \lambda_{n-1} \leq \mu_{n-1} \leq \lambda_n,$$

a genuine interlacing. We obtain similar results with the matrix P_{n-r} obtained by dropping the last $n - r$ columns of the identity matrix and setting $B = P_{n-r}^\top A P_{n-r}$ (B is the $r \times r$ matrix obtained from A by deleting its last $n - r$ rows and columns). In this case, we have the following interlacing inequalities known as *Cauchy interlacing theorem*:

$$\lambda_k \leq \mu_k \leq \lambda_{k+n-r}, \quad k = 1, \dots, r. \quad (*)$$

Another useful tool to prove eigenvalue equalities is the Courant–Fischer characterization of the eigenvalues of a symmetric matrix, also known as the Min-max (and Max-min) theorem.

Theorem 13.26. (*Courant–Fischer*) *Let A be a symmetric $n \times n$ matrix with eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ and let (u_1, \dots, u_n) be any orthonormal basis of eigenvectors of A , where u_i is a unit eigenvector associated with λ_i . If \mathcal{V}_k denotes the set of subspaces of \mathbb{R}^n of dimension k , then*

$$\begin{aligned} \lambda_k &= \max_{W \in \mathcal{V}_{n-k+1}} \min_{x \in W, x \neq 0} \frac{x^\top A x}{x^\top x} \\ \lambda_k &= \min_{W \in \mathcal{V}_k} \max_{x \in W, x \neq 0} \frac{x^\top A x}{x^\top x}. \end{aligned}$$

Proof. Let us consider the second equality, the proof of the first equality being similar. Observe that the space V_k spanned by (u_1, \dots, u_k) has dimension k , and by Proposition 13.22, we have

$$\lambda_k = \max_{x \neq 0, x \in V_k} \frac{x^\top A x}{x^\top x} \geq \min_{W \in \mathcal{V}_k} \max_{x \in W, x \neq 0} \frac{x^\top A x}{x^\top x}.$$

Therefore, we need to prove the reverse inequality; that is, we have to show that

$$\lambda_k \leq \max_{x \neq 0, x \in W} \frac{x^\top A x}{x^\top x}, \quad \text{for all } W \in \mathcal{V}_k.$$

Now, for any $W \in \mathcal{V}_k$, if we can prove that $W \cap V_{k-1}^\perp \neq (0)$, then for any nonzero $v \in W \cap V_{k-1}^\perp$, by Proposition 13.23, we have

$$\lambda_k = \min_{x \neq 0, x \in V_{k-1}^\perp} \frac{x^\top A x}{x^\top x} \leq \frac{v^\top A v}{v^\top v} \leq \max_{x \in W, x \neq 0} \frac{x^\top A x}{x^\top x}.$$

It remains to prove that $\dim(W \cap V_{k-1}^\perp) \geq 1$. However, $\dim(V_{k-1}) = k - 1$, so $\dim(V_{k-1}^\perp) = n - k + 1$, and by hypothesis $\dim(W) = k$. By the Grassmann relation,

$$\dim(W) + \dim(V_{k-1}^\perp) = \dim(W \cap V_{k-1}^\perp) + \dim(W + V_{k-1}^\perp),$$

and since $\dim(W + V_{k-1}^\perp) \leq \dim(\mathbb{R}^n) = n$, we get

$$k + n - k + 1 \leq \dim(W \cap V_{k-1}^\perp) + n;$$

that is, $1 \leq \dim(W \cap V_{k-1}^\perp)$, as claimed. \square

The Courant–Fischer theorem yields the following useful result about perturbing the eigenvalues of a symmetric matrix due to Hermann Weyl.

Proposition 13.27. *Given two $n \times n$ symmetric matrices A and $B = A + \delta A$, if $\alpha_1 \leq \alpha_2 \leq \dots \leq \alpha_n$ are the eigenvalues of A and $\beta_1 \leq \beta_2 \leq \dots \leq \beta_n$ are the eigenvalues of B , then*

$$|\alpha_k - \beta_k| \leq \rho(\delta A) \leq \|\delta A\|_2, \quad k = 1, \dots, n.$$

Proof. Let \mathcal{V}_k be defined as in the Courant–Fischer theorem and let V_k be the subspace spanned by the k eigenvectors associated with $\lambda_1, \dots, \lambda_k$. By the Courant–Fischer theorem applied to B , we have

$$\begin{aligned} \beta_k &= \min_{W \in \mathcal{V}_k} \max_{x \in W, x \neq 0} \frac{x^\top Bx}{x^\top x} \\ &\leq \max_{x \in V_k} \frac{x^\top Bx}{x^\top x} \\ &= \max_{x \in V_k} \left(\frac{x^\top Ax}{x^\top x} + \frac{x^\top \delta Ax}{x^\top x} \right) \\ &\leq \max_{x \in V_k} \frac{x^\top Ax}{x^\top x} + \max_{x \in V_k} \frac{x^\top \delta Ax}{x^\top x}. \end{aligned}$$

By Proposition 13.22, we have

$$\alpha_k = \max_{x \in V_k} \frac{x^\top Ax}{x^\top x},$$

so we obtain

$$\begin{aligned} \beta_k &\leq \max_{x \in V_k} \frac{x^\top Ax}{x^\top x} + \max_{x \in V_k} \frac{x^\top \delta Ax}{x^\top x} \\ &= \alpha_k + \max_{x \in V_k} \frac{x^\top \delta Ax}{x^\top x} \\ &\leq \alpha_k + \max_{x \in \mathbb{R}^n} \frac{x^\top \delta Ax}{x^\top x}. \end{aligned}$$

Now, by Proposition 13.22 and Proposition 7.7, we have

$$\max_{x \in \mathbb{R}^n} \frac{x^\top \delta Ax}{x^\top x} = \max_i \lambda_i(\delta A) \leq \rho(\delta A) \leq \|\delta A\|_2,$$

where $\lambda_i(\delta A)$ denotes the i th eigenvalue of δA , which implies that

$$\beta_k \leq \alpha_k + \rho(\delta A) \leq \alpha_k + \|\delta A\|_2.$$

By exchanging the roles of A and B , we also have

$$\alpha_k \leq \beta_k + \rho(\delta A) \leq \beta_k + \|\delta A\|_2,$$

and thus,

$$|\alpha_k - \beta_k| \leq \rho(\delta A) \leq \|\delta A\|_2, \quad k = 1, \dots, n,$$

as claimed. □

Proposition 13.27 also holds for Hermitian matrices.

A pretty result of Wielandt and Hoffman asserts that

$$\sum_{k=1}^n (\alpha_k - \beta_k)^2 \leq \|\delta A\|_F^2,$$

where $\|\cdot\|_F$ is the Frobenius norm. However, the proof is significantly harder than the above proof; see Lax [71].

The Courant–Fischer theorem can also be used to prove some famous inequalities due to Hermann Weyl. Given two symmetric (or Hermitian) matrices A and B , let $\lambda_i(A)$, $\lambda_i(B)$, and $\lambda_i(A+B)$ denote the i th eigenvalue of A , B , and $A+B$, respectively, arranged in nondecreasing order.

Proposition 13.28. (*Weyl*) *Given two symmetric (or Hermitian) $n \times n$ matrices A and B , the following inequalities hold: For all i, j, k with $1 \leq i, j, k \leq n$:*

1. *If $i + j = k + 1$, then*

$$\lambda_i(A) + \lambda_j(B) \leq \lambda_k(A+B).$$

2. *If $i + j = k + n$, then*

$$\lambda_k(A+B) \leq \lambda_i(A) + \lambda_j(B).$$

Proof. Observe that the first set of inequalities is obtained from the second set by replacing A by $-A$ and B by $-B$, so it is enough to prove the second set of inequalities. By the Courant–Fischer theorem, there is a subspace H of dimension $n - k + 1$ such that

$$\lambda_k(A+B) = \min_{x \in H, x \neq 0} \frac{x^\top (A+B)x}{x^\top x}.$$

Similarly, there exist a subspace F of dimension i and a subspace G of dimension j such that

$$\lambda_i(A) = \max_{x \in F, x \neq 0} \frac{x^\top Ax}{x^\top x}, \quad \lambda_j(B) = \max_{x \in G, x \neq 0} \frac{x^\top Bx}{x^\top x}.$$

We claim that $F \cap G \cap H \neq (0)$. To prove this, we use the Grassmann relation twice. First, $\dim(F \cap G \cap H) = \dim(F) + \dim(G \cap H) - \dim(F + (G \cap H)) \geq \dim(F) + \dim(G \cap H) - n$, and second,

$$\dim(G \cap H) = \dim(G) + \dim(H) - \dim(G + H) \geq \dim(G) + \dim(H) - n,$$

so

$$\dim(F \cap G \cap H) \geq \dim(F) + \dim(G) + \dim(H) - 2n.$$

However,

$$\dim(F) + \dim(G) + \dim(H) = i + j + n - k + 1$$

and $i + j = k + n$, so we have

$$\dim(F \cap G \cap H) \geq i + j + n - k + 1 - 2n = k + n + n - k + 1 - 2n = 1,$$

which shows that $F \cap G \cap H \neq (0)$. Then, for any unit vector $z \in F \cap G \cap H \neq (0)$, we have

$$\lambda_k(A + B) \leq z^\top (A + B)z, \quad \lambda_i(A) \geq z^\top Az, \quad \lambda_j(B) \geq z^\top Bz,$$

establishing the desired inequality $\lambda_k(A + B) \leq \lambda_i(A) + \lambda_j(B)$. \square

In the special case $i = j = k$, we obtain

$$\lambda_1(A) + \lambda_1(B) \leq \lambda_1(A + B), \quad \lambda_n(A + B) \leq \lambda_n(A) + \lambda_n(B).$$

It follows that λ_1 is concave, while λ_n is convex.

If $i = 1$ and $j = k$, we obtain

$$\lambda_1(A) + \lambda_k(B) \leq \lambda_k(A + B),$$

and if $i = k$ and $j = n$, we obtain

$$\lambda_k(A + B) \leq \lambda_k(A) + \lambda_n(B),$$

and combining them, we get

$$\lambda_1(A) + \lambda_k(B) \leq \lambda_k(A + B) \leq \lambda_k(A) + \lambda_n(B).$$

In particular, if B is positive semidefinite, since its eigenvalues are nonnegative, we obtain the following inequality known as the *monotonicity theorem* for symmetric (or Hermitian) matrices: if A and B are symmetric (or Hermitian) and B is positive semidefinite, then

$$\lambda_k(A) \leq \lambda_k(A + B) \quad k = 1, \dots, n.$$

The reader is referred to Horn and Johnson [57] (Chapters 4 and 7) for a very complete treatment of matrix inequalities and interlacing results, and also to Lax [71] and Serre [96].

We now have all the tools to present the important *singular value decomposition* (SVD) and the *polar form* of a matrix. However, we prefer to first illustrate how the material of this section can be used to discretize boundary value problems, and we give a brief introduction to the finite elements method.

13.7 Summary

The main concepts and results of this chapter are listed below:

- *Normal* linear maps, *self-adjoint* linear maps, *skew-self-adjoint* linear maps, and *orthogonal* linear maps.
- Properties of the eigenvalues and eigenvectors of a normal linear map.
- The *complexification* of a real vector space, of a linear map, and of a Euclidean inner product.
- The eigenvalues of a self-adjoint map in a Hermitian space are *real*.
- The eigenvalues of a self-adjoint map in a Euclidean space are *real*.
- Every self-adjoint linear map on a Euclidean space has an orthonormal basis of eigenvectors.
- Every normal linear map on a Euclidean space can be block diagonalized (blocks of size at most 2×2) with respect to an orthonormal basis of eigenvectors.
- Every normal linear map on a Hermitian space can be diagonalized with respect to an orthonormal basis of eigenvectors.
- The spectral theorems for self-adjoint, skew-self-adjoint, and orthogonal linear maps (on a Euclidean space).
- The spectral theorems for normal, symmetric, skew-symmetric, and orthogonal (real) matrices.
- The spectral theorems for normal, Hermitian, skew-Hermitian, and unitary (complex) matrices.
- The conditioning of eigenvalue problems.
- The *Rayleigh ratio* and the *Rayleigh–Ritz theorem*.
- *Interlacing inequalities* and the *Cauchy interlacing theorem*.
- The *Poincaré separation theorem*.
- The *Courant–Fischer theorem*.
- Inequalities involving perturbations of the eigenvalues of a symmetric matrix.
- The *Weyl inequalities*.

Chapter 14

Bilinear Forms and Their Geometries

14.1 Bilinear Forms

In this chapter, we study the structure of a K -vector space E endowed with a nondegenerate bilinear form $\varphi: E \times E \rightarrow K$ (for any field K), which can be viewed as a kind of generalized inner product. Unlike the case of an inner product, there may be nonzero vectors $u \in E$ such that $\varphi(u, u) = 0$, so the map $u \mapsto \varphi(u, u)$ can no longer be interpreted as a notion of square length (also, $\varphi(u, u)$ may not be real and positive!). However, the notion of orthogonality survives: we say that $u, v \in E$ are orthogonal iff $\varphi(u, v) = 0$. Under some additional conditions on φ , it is then possible to split E into orthogonal subspaces having some special properties. It turns out that the special cases where φ is symmetric (or Hermitian) or skew-symmetric (or skew-Hermitian) can be handled uniformly using a deep theorem due to Witt (the Witt decomposition theorem (1936)).

We begin with the very general situation of a bilinear form $\varphi: E \times F \rightarrow K$, where K is an arbitrary field, possibly of characteristic 2. Actually, even though at first glance this may appear to be an unnecessary abstraction, it turns out that this situation arises in attempting to prove properties of a bilinear map $\varphi: E \times E \rightarrow K$, because it may be necessary to restrict φ to different subspaces U and V of E . This general approach was pioneered by Chevalley [22], E. Artin [3], and Bourbaki [13]. The third source was a major source of inspiration, and many proofs are taken from it. Other useful references include Snapper and Troyer [99], Berger [9], Jacobson [59], Grove [52], Taylor [108], and Berndt [11].

Definition 14.1. Given two vector spaces E and F over a field K , a map $\varphi: E \times F \rightarrow K$ is a *bilinear form* iff the following conditions hold: For all $u, u_1, u_2 \in E$, all $v, v_1, v_2 \in F$, for all $\lambda, \mu \in K$, we have

$$\begin{aligned}\varphi(u_1 + u_2, v) &= \varphi(u_1, v) + \varphi(u_2, v) \\ \varphi(u, v_1 + v_2) &= \varphi(u, v_1) + \varphi(u, v_2) \\ \varphi(\lambda u, v) &= \lambda \varphi(u, v) \\ \varphi(u, \mu v) &= \mu \varphi(u, v).\end{aligned}$$

A bilinear form as in Definition 14.1 is sometimes called a *pairing*. The first two conditions imply that $\varphi(0, v) = \varphi(u, 0) = 0$ for all $u \in E$ and all $v \in F$.

If $E = F$, observe that

$$\begin{aligned}\varphi(\lambda u + \mu v, \lambda u + \mu v) &= \lambda\varphi(u, \lambda u + \mu v) + \mu\varphi(v, \lambda u + \mu v) \\ &= \lambda^2\varphi(u, u) + \lambda\mu\varphi(u, v) + \lambda\mu\varphi(v, u) + \mu^2\varphi(v, v).\end{aligned}$$

If we let $\lambda = \mu = 1$, we get

$$\varphi(u + v, u + v) = \varphi(u, u) + \varphi(u, v) + \varphi(v, u) + \varphi(v, v).$$

If φ is *symmetric*, which means that

$$\varphi(u, v) = \varphi(v, u) \quad \text{for all } u, v \in E,$$

then

$$2\varphi(u, v) = \varphi(u + v, u + v) - \varphi(u, u) - \varphi(v, v).$$

The function Φ defined such that

$$\Phi(u) = \varphi(u, u) \quad u \in E,$$

is called the *quadratic form* associated with φ . If the field K is not of characteristic 2, then φ is completely determined by its quadratic form Φ . The symmetric bilinear form φ is called the *polar form* of Φ . This suggests the following definition.

Definition 14.2. A function $\Phi: E \rightarrow K$ is a *quadratic form* on E if the following conditions hold:

- (1) We have $\Phi(\lambda u) = \lambda^2\Phi(u)$, for all $u \in E$ and all $\lambda \in E$.
- (2) The map φ' given by $\varphi'(u, v) = \Phi(u + v) - \Phi(u) - \Phi(v)$ is bilinear. Obviously, the map φ' is symmetric.

Since $\Phi(x + x) = \Phi(2x) = 4\Phi(x)$, we have

$$\varphi'(u, u) = 2\Phi(u) \quad u \in E.$$

If the field K is not of characteristic 2, then $\varphi = \frac{1}{2}\varphi'$ is the unique symmetric bilinear form such that $\varphi(u, u) = \Phi(u)$ for all $u \in E$. The bilinear form $\varphi = \frac{1}{2}\varphi'$ is called the *polar form* of Φ . In this case, there is a bijection between the set of bilinear forms on E and the set of quadratic forms on E .

If K is a field of characteristic 2, then φ' is *alternating*, which means that

$$\varphi'(u, u) = 0 \quad \text{for all } u \in E.$$

Thus, Φ cannot be recovered from the symmetric bilinear form φ' . However, there is some (nonsymmetric) bilinear form ψ such that $\Phi(u) = \psi(u, u)$ for all $u \in E$. Thus, quadratic forms are more general than symmetric bilinear forms (except in characteristic $\neq 2$).

In general, if K is a field of any characteristic, the identity

$$\varphi(u + v, u + v) = \varphi(u, u) + \varphi(u, v) + \varphi(v, u) + \varphi(v, v)$$

shows that if φ is alternating (that is, $\varphi(u, u) = 0$ for all $u \in E$), then,

$$\varphi(v, u) = -\varphi(u, v) \quad \text{for all } u, v \in E;$$

we say that φ is *skew-symmetric*. Conversely, if the field K is not of characteristic 2, then a skew-symmetric bilinear map is alternating, since $\varphi(u, u) = -\varphi(u, u)$ implies $\varphi(u, u) = 0$.

An important consequence of bilinearity is that a pairing yields a linear map from E into F^* and a linear map from F into E^* (where $E^* = \text{Hom}_K(E, K)$, the *dual* of E , is the set of linear maps from E to K , called *linear forms*).

Definition 14.3. Given a bilinear map $\varphi: E \times F \rightarrow K$, for every $u \in E$, let $l_\varphi(u)$ be the linear form in F^* given by

$$l_\varphi(u)(y) = \varphi(u, y) \quad \text{for all } y \in F,$$

and for every $v \in F$, let $r_\varphi(v)$ be the linear form in E^* given by

$$r_\varphi(v)(x) = \varphi(x, v) \quad \text{for all } x \in E.$$

Because φ is bilinear, the maps $l_\varphi: E \rightarrow F^*$ and $r_\varphi: F \rightarrow E^*$ are linear.

Definition 14.4. A bilinear map $\varphi: E \times F \rightarrow K$ is said to be *nondegenerate* iff the following conditions hold:

- (1) For every $u \in E$, if $\varphi(u, v) = 0$ for all $v \in F$, then $u = 0$, and
- (2) For every $v \in F$, if $\varphi(u, v) = 0$ for all $u \in E$, then $v = 0$.

The following proposition shows the importance of l_φ and r_φ .

Proposition 14.1. *Given a bilinear map $\varphi: E \times F \rightarrow K$, the following properties hold:*

- (a) *The map l_φ is injective iff property (1) of Definition 14.4 holds.*
- (b) *The map r_φ is injective iff property (2) of Definition 14.4 holds.*
- (c) *The bilinear form φ is nondegenerate and iff l_φ and r_φ are injective.*
- (d) *If the bilinear form φ is nondegenerate and if E and F have finite dimensions, then $\dim(E) = \dim(F)$, and $l_\varphi: E \rightarrow F^*$ and $r_\varphi: F \rightarrow E^*$ are linear isomorphisms.*

Proof. (a) Assume that (1) of Definition 14.4 holds. If $l_\varphi(u) = 0$, then $l_\varphi(u)$ is the linear form whose value is 0 for all y ; that is,

$$l_\varphi(u)(y) = \varphi(u, y) = 0 \quad \text{for all } y \in F,$$

and by (1) of Definition 14.4, we must have $u = 0$. Therefore, l_φ is injective. Conversely, if l_φ is injective, and if

$$l_\varphi(u)(y) = \varphi(u, y) = 0 \quad \text{for all } y \in F,$$

then $l_\varphi(u)$ is the zero form, and by injectivity of l_φ , we get $u = 0$; that is, (1) of Definition 14.4 holds.

(b) The proof is obtained by swapping the arguments of φ .

(c) This follows from (a) and (b).

(d) If E and F are finite dimensional, then $\dim(E) = \dim(E^*)$ and $\dim(F) = \dim(F^*)$. Since φ is nondegenerate, $l_\varphi: E \rightarrow F^*$ and $r_\varphi: F \rightarrow E^*$ are injective, so $\dim(E) \leq \dim(F^*) = \dim(F)$ and $\dim(F) \leq \dim(E^*) = \dim(E)$, which implies that

$$\dim(E) = \dim(F),$$

and thus, $l_\varphi: E \rightarrow F^*$ and $r_\varphi: F \rightarrow E^*$ are bijective. □

As a corollary of Proposition 14.1, we have the following characterization of a nondegenerate bilinear map. The proof is left as an exercise.

Proposition 14.2. *Given a bilinear map $\varphi: E \times F \rightarrow K$, if E and F have the same finite dimension, then the following properties are equivalent:*

- (1) *The map l_φ is injective.*
- (2) *The map l_φ is surjective.*
- (3) *The map r_φ is injective.*
- (4) *The map r_φ is surjective.*
- (5) *The bilinear form φ is nondegenerate.*

Observe that in terms of the canonical pairing between E^* and E given by

$$\langle f, u \rangle = f(u), \quad f \in E^*, u \in E,$$

(and the canonical pairing between F^* and F), we have

$$\varphi(u, v) = \langle l_\varphi(u), v \rangle = \langle r_\varphi(v), u \rangle.$$

Proposition 14.3. *Given a bilinear map $\varphi: E \times F \rightarrow K$, if φ is nondegenerate and E and F are finite-dimensional, then $\dim(E) = \dim(F) = n$, and for every basis (e_1, \dots, e_n) of E , there is a basis (f_1, \dots, f_n) of F such that $\varphi(e_i, f_j) = \delta_{ij}$, for all $i, j = 1, \dots, n$.*

Proof. Since φ is nondegenerate, by Proposition 14.1 we have $\dim(E) = \dim(F) = n$, and by Proposition 14.2, the linear map r_φ is bijective. Then, if (e_1^*, \dots, e_n^*) is the dual basis (in E^*) of the basis (e_1, \dots, e_n) , the vectors (f_1, \dots, f_n) given by $f_i = r_\varphi^{-1}(e_i^*)$ form a basis of F , and we have

$$\varphi(e_i, f_j) = \langle r_\varphi(f_j), e_i \rangle = \langle e_i^*, e_j \rangle = \delta_{ij},$$

as claimed. \square

If $E = F$ and φ is symmetric, then we have the following interesting result.

Theorem 14.4. *Given any bilinear form $\varphi: E \times E \rightarrow K$ with $\dim(E) = n$, if φ is symmetric and K does not have characteristic 2, then there is a basis (e_1, \dots, e_n) of E such that $\varphi(e_i, e_j) = 0$, for all $i \neq j$.*

Proof. We proceed by induction on $n \geq 0$, following a proof due to Chevalley. The base case $n = 0$ is trivial. For the induction step, assume that $n \geq 1$ and that the induction hypothesis holds for all vector spaces of dimension $n - 1$. If $\varphi(u, v) = 0$ for all $u, v \in E$, then the statement holds trivially. Otherwise, since K does not have characteristic 2, by a previous remark, there is some nonzero vector $e_1 \in E$ such that $\varphi(e_1, e_1) \neq 0$. We claim that the set

$$H = \{v \in E \mid \varphi(e_1, v) = 0\}$$

has dimension $n - 1$, and that $e_1 \notin H$.

This is because

$$H = \text{Ker}(l_\varphi(e_1)),$$

where $l_\varphi(e_1)$ is the linear form in E^* determined by e_1 . Since $\varphi(e_1, e_1) \neq 0$, we have $e_1 \notin H$, the linear form $l_\varphi(e_1)$ is not the zero form, and thus its kernel is a hyperplane H (a subspace of dimension $n - 1$). Since $\dim(H) = n - 1$ and $e_1 \notin H$, we have the direct sum

$$E = H \oplus Ke_1.$$

By the induction hypothesis applied to H , we get a basis (e_2, \dots, e_n) of vectors in H such that $\varphi(e_i, e_j) = 0$, for all $i \neq j$ with $2 \leq i, j \leq n$. Since $\varphi(e_1, v) = 0$ for all $v \in H$ and since φ is symmetric, we also have $\varphi(v, e_1) = 0$ for all $v \in H$, so we obtain a basis (e_1, \dots, e_n) of E such that $\varphi(e_i, e_j) = 0$, for all $i \neq j$. \square

If E and F are finite-dimensional vector spaces and if (e_1, \dots, e_m) is a basis of E and (f_1, \dots, f_n) is a basis of F then the bilinearity of φ yields

$$\varphi\left(\sum_{i=1}^m x_i e_i, \sum_{j=1}^n y_j f_j\right) = \sum_{i=1}^m \sum_{j=1}^n x_i \varphi(e_i, f_j) y_j.$$

This shows that φ is completely determined by the $m \times n$ matrix $M = (\varphi(e_i, e_j))$, and in matrix form, we have

$$\varphi(x, y) = x^\top M y = y^\top M^\top x,$$

where x and y are the column vectors associated with $(x_1, \dots, x_m) \in K^m$ and $(y_1, \dots, y_n) \in K^n$. As in Section 10.1, we are committing the slight abuse of notation of letting x denote both the vector $x = \sum_{i=1}^n x_i e_i$ and the column vector associated with (x_1, \dots, x_n) (and similarly for y). We call M the *matrix of φ with respect to the bases (e_1, \dots, e_m) and (f_1, \dots, f_n)* .

If $m = \dim(E) = \dim(F) = n$, then it is easy to check that φ is nondegenerate iff M is invertible iff $\det(M) \neq 0$.

As we will see later, most bilinear forms that we will encounter are equivalent to one whose matrix is of the following form:

1. $I_n, -I_n$.

2. If $p + q = n$, with $p, q \geq 1$,

$$I_{p,q} = \begin{pmatrix} I_p & 0 \\ 0 & -I_q \end{pmatrix}$$

3. If $n = 2m$,

$$J_{m,m} = \begin{pmatrix} 0 & I_m \\ -I_m & 0 \end{pmatrix}$$

4. If $n = 2m$,

$$A_{m,m} = I_{m,m} J_{m,m} = \begin{pmatrix} 0 & I_m \\ I_m & 0 \end{pmatrix}.$$

If we make changes of bases given by matrices P and Q , so that $x = Px'$ and $y = Qy'$, then the new matrix expressing φ is $P^\top M Q$. In particular, if $E = F$ and the same basis is used, then the new matrix is $P^\top M P$. This shows that if φ is nondegenerate, then the determinant of φ is determined up to a square element.

Observe that if φ is a symmetric bilinear form ($E = F$) and if K does not have characteristic 2, then by Theorem 14.4, there is a basis of E with respect to which the matrix M representing φ is a diagonal matrix. If $K = \mathbb{R}$ or $K = \mathbb{C}$, this allows us to classify completely the symmetric bilinear forms. Recall that $\Phi(u) = \varphi(u, u)$ for all $u \in E$.

Proposition 14.5. *Given any bilinear form $\varphi: E \times E \rightarrow K$ with $\dim(E) = n$, if φ is symmetric and K does not have characteristic 2, then there is a basis (e_1, \dots, e_n) of E such that*

$$\Phi\left(\sum_{i=1}^n x_i e_i\right) = \sum_{i=1}^r \lambda_i x_i^2,$$

for some $\lambda_i \in K - \{0\}$ and with $r \leq n$. Furthermore, if $K = \mathbb{C}$, then there is a basis (e_1, \dots, e_n) of E such that

$$\Phi\left(\sum_{i=1}^n x_i e_i\right) = \sum_{i=1}^r x_i^2,$$

and if $K = \mathbb{R}$, then there is a basis (e_1, \dots, e_n) of E such that

$$\Phi\left(\sum_{i=1}^n x_i e_i\right) = \sum_{i=1}^p x_i^2 - \sum_{i=p+1}^{p+q} x_i^2,$$

with $0 \leq p, q$ and $p + q \leq n$.

Proof. The first statement is a direct consequence of Theorem 14.4. If $K = \mathbb{C}$, then every λ_i has a square root μ_i , and if replace e_i by e_i/μ_i , we obtained the desired form.

If $K = \mathbb{R}$, then there are two cases:

1. If $\lambda_i > 0$, let μ_i be a positive square root of λ_i and replace e_i by e_i/μ_i .
2. If $\lambda_i < 0$, let μ_i be a positive square root of $-\lambda_i$ and replace e_i by e_i/μ_i .

□

In the nondegenerate case, the matrices corresponding to the complex and the real case are, $I_n, -I_n$, and $I_{p,q}$. Observe that the second statement of Proposition 14.4 holds in any field in which every element has a square root. In the case $K = \mathbb{R}$, we can show that (p, q) only depends on φ .

For any subspace U of E , we say that φ is *positive definite on U* iff $\varphi(u, u) > 0$ for all nonzero $u \in U$, and we say that φ is *negative definite on U* iff $\varphi(u, u) < 0$ for all nonzero $u \in U$. Then, let

$$r = \max\{\dim(U) \mid U \subseteq E, \varphi \text{ is positive definite on } U\}$$

and let

$$s = \max\{\dim(U) \mid U \subseteq E, \varphi \text{ is negative definite on } U\}$$

Proposition 14.6. (*Sylvester's inertia law*) Given any symmetric bilinear form $\varphi: E \times E \rightarrow \mathbb{R}$ with $\dim(E) = n$, for any basis (e_1, \dots, e_n) of E such that

$$\Phi\left(\sum_{i=1}^n x_i e_i\right) = \sum_{i=1}^p x_i^2 - \sum_{i=p+1}^{p+q} x_i^2,$$

with $0 \leq p, q$ and $p + q \leq n$, the integers p, q depend only on φ ; in fact, $p = r$ and $q = s$, with r and s as defined above.

Proof. If we let U be the subspace spanned by (e_1, \dots, e_p) , then φ is positive definite on U , so $r \geq p$. Similarly, if we let V be the subspace spanned by $(e_{p+1}, \dots, e_{p+q})$, then φ is negative definite on V , so $s \geq q$.

Next, if W_1 is any subspace of maximum dimension such that φ is positive definite on W_1 , and if we let V' be the subspace spanned by (e_{p+1}, \dots, e_n) , then $\varphi(u, u) \leq 0$ on V' , so $W_1 \cap V' = (0)$, which implies that $\dim(W_1) + \dim(V') \leq n$, and thus, $r + n - p \leq n$; that is, $r \leq p$. Similarly, if W_2 is any subspace of maximum dimension such that φ is negative definite on W_2 , and if we let U' be the subspace spanned by $(e_1, \dots, e_p, e_{p+q+1}, \dots, e_n)$, then $\varphi(u, u) \geq 0$ on U' , so $W_2 \cap U' = (0)$, which implies that $s + n - q \leq n$; that is, $s \leq q$. Therefore, $p = r$ and $q = s$, as claimed \square

These last two results can be generalized to ordered fields. For example, see Snapper and Troyer [99], Artin [3], and Bourbaki [13].

14.2 Sesquilinear Forms

In order to accomodate Hermitian forms, we assume that some involutive automorphism, $\lambda \mapsto \bar{\lambda}$, of the field K is given. This automorphism of K satisfies the following properties:

$$\begin{aligned} \overline{(\lambda + \mu)} &= \bar{\lambda} + \bar{\mu} \\ \overline{(\lambda\mu)} &= \bar{\lambda}\bar{\mu} \\ \overline{\bar{\lambda}} &= \lambda. \end{aligned}$$

If the automorphism $\lambda \mapsto \bar{\lambda}$ is the identity, then we are in the standard situation of a bilinear form. When $K = \mathbb{C}$ (the complex numbers), then we usually pick the automorphism of \mathbb{C} to be *conjugation*; namely, the map

$$a + ib \mapsto a - ib.$$

Definition 14.5. Given two vector spaces E and F over a field K with an involutive automorphism $\lambda \mapsto \bar{\lambda}$, a map $\varphi: E \times F \rightarrow K$ is a (right) *sesquilinear form* iff the following conditions hold: For all $u, u_1, u_2 \in E$, all $v, v_1, v_2 \in F$, for all $\lambda, \mu \in K$, we have

$$\begin{aligned} \varphi(u_1 + u_2, v) &= \varphi(u_1, v) + \varphi(u_2, v) \\ \varphi(u, v_1 + v_2) &= \varphi(u, v_1) + \varphi(u, v_2) \\ \varphi(\lambda u, v) &= \lambda \varphi(u, v) \\ \varphi(u, \mu v) &= \bar{\mu} \varphi(u, v). \end{aligned}$$

Again, $\varphi(0, v) = \varphi(u, 0) = 0$. If $E = F$, then we have

$$\begin{aligned} \varphi(\lambda u + \mu v, \lambda u + \mu v) &= \lambda \varphi(u, \lambda u + \mu v) + \mu \varphi(v, \lambda u + \mu v) \\ &= \lambda \bar{\lambda} \varphi(u, u) + \lambda \bar{\mu} \varphi(u, v) + \bar{\lambda} \mu \varphi(v, u) + \mu \bar{\mu} \varphi(v, v). \end{aligned}$$

If we let $\lambda = \mu = 1$ and then $\lambda = 1, \mu = -1$, we get

$$\begin{aligned}\varphi(u + v, u + v) &= \varphi(u, u) + \varphi(u, v) + \varphi(v, u) + \varphi(v, v) \\ \varphi(u - v, u - v) &= \varphi(u, u) - \varphi(u, v) - \varphi(v, u) + \varphi(v, v),\end{aligned}$$

so by subtraction, we get

$$2(\varphi(u, v) + \varphi(v, u)) = \varphi(u + v, u + v) - \varphi(u - v, u - v) \quad \text{for } u, v \in E.$$

If we replace v by λv (with $\lambda \neq 0$), we get

$$2(\bar{\lambda}\varphi(u, v) + \lambda\varphi(v, u)) = \varphi(u + \lambda v, u + \lambda v) - \varphi(u - \lambda v, u - \lambda v),$$

and by combining the above two equations, we get

$$2(\lambda - \bar{\lambda})\varphi(u, v) = \lambda\varphi(u + v, u + v) - \lambda\varphi(u - v, u - v) - \varphi(u + \lambda v, u + \lambda v) + \varphi(u - \lambda v, u - \lambda v).$$

If the automorphism $\lambda \mapsto \bar{\lambda}$ is not the identity, then there is some $\lambda \in K$ such that $\lambda - \bar{\lambda} \neq 0$, and if K is not of characteristic 2, then we see that the sesquilinear form φ is completely determined by its restriction to the diagonal (that is, the set of values $\{\varphi(u, u) \mid u \in E\}$). In the special case where $K = \mathbb{C}$, we can pick $\lambda = i$, and we get

$$4\varphi(u, v) = \varphi(u + v, u + v) - \varphi(u - v, u - v) + i\varphi(u + \lambda v, u + \lambda v) - i\varphi(u - \lambda v, u - \lambda v).$$

Remark: If the automorphism $\lambda \mapsto \bar{\lambda}$ is the identity, then in general φ is not determined by its value on the diagonal, unless φ is symmetric.

In the sesquilinear setting, it turns out that the following two cases are of interest:

1. We have

$$\varphi(v, u) = \overline{\varphi(u, v)}, \quad \text{for all } u, v \in E,$$

in which case we say that φ is *Hermitian*. In the special case where $K = \mathbb{C}$ and the involutive automorphism is conjugation, we see that $\varphi(u, u) \in \mathbb{R}$, for $u \in E$.

2. We have

$$\varphi(v, u) = -\overline{\varphi(u, v)}, \quad \text{for all } u, v \in E,$$

in which case we say that φ is *skew-Hermitian*.

We observed that in characteristic different from 2, a sesquilinear form is determined by its restriction to the diagonal. For Hermitian and skew-Hermitian forms, we have the following kind of converse.

Proposition 14.7. *If φ is a nonzero Hermitian or skew-Hermitian form and if $\varphi(u, u) = 0$ for all $u \in E$, then K is of characteristic 2 and the automorphism $\lambda \mapsto \bar{\lambda}$ is the identity.*

Proof. We give the proof in the Hermitian case, the skew-Hermitian case being left as an exercise. Assume that φ is alternating. From the identity

$$\varphi(u + v, u + v) = \varphi(u, u) + \varphi(u, v) + \overline{\varphi(u, v)} + \varphi(v, v),$$

we get

$$\varphi(u, v) = -\overline{\varphi(u, v)} \quad \text{for all } u, v \in E.$$

Since φ is not the zero form, there exist some nonzero vectors $u, v \in E$ such that $\varphi(u, v) = 1$. For any $\lambda \in K$, we have

$$\lambda\varphi(u, v) = \varphi(\lambda u, v) = -\overline{\varphi(\lambda u, v)} = -\overline{\lambda}\overline{\varphi(u, v)},$$

and since $\varphi(u, v) = 1$, we get

$$\lambda = -\overline{\lambda} \quad \text{for all } \lambda \in K.$$

For $\lambda = 1$, we get $1 = -1$, which means that K has characteristic 2. But then

$$\lambda = -\overline{\lambda} = \overline{\lambda} \quad \text{for all } \lambda \in K,$$

so the automorphism $\lambda \mapsto \overline{\lambda}$ is the identity. □

The definition of the linear maps l_φ and r_φ requires a small twist due to the automorphism $\lambda \mapsto \overline{\lambda}$.

Definition 14.6. Given a vector space E over a field K with an involutive automorphism $\lambda \mapsto \overline{\lambda}$, we define the K -vector space \overline{E} as E with its abelian group structure, but with scalar multiplication given by

$$(\lambda, u) \mapsto \overline{\lambda}u.$$

Given two K -vector spaces E and F , a *semilinear map* $f: E \rightarrow F$ is a function, such that for all $u, v \in E$, for all $\lambda \in K$, we have

$$\begin{aligned} f(u + v) &= f(u) + f(v) \\ f(\lambda u) &= \overline{\lambda}f(u). \end{aligned}$$

Because $\overline{\overline{\lambda}} = \lambda$, observe that a function $f: E \rightarrow F$ is semilinear iff it is a linear map $f: \overline{E} \rightarrow F$. The K -vector spaces E and \overline{E} are isomorphic, since any basis $(e_i)_{i \in I}$ of E is also a basis of \overline{E} .

The maps l_φ and r_φ are defined as follows:

For every $u \in E$, let $l_\varphi(u)$ be the linear form in F^* defined so that

$$l_\varphi(u)(y) = \overline{\varphi(u, y)} \quad \text{for all } y \in F,$$

and for every $v \in F$, let $r_\varphi(v)$ be the linear form in E^* defined so that

$$r_\varphi(v)(x) = \varphi(x, v) \quad \text{for all } x \in E.$$

The reader should check that because we used $\overline{\varphi(u, y)}$ in the definition of $l_\varphi(u)(y)$, the function $l_\varphi(u)$ is indeed a linear form in F^* . It is also easy to check that l_φ is a linear map $l_\varphi: \overline{E} \rightarrow F^*$, and that r_φ is a linear map $r_\varphi: \overline{F} \rightarrow E^*$ (equivalently, $l_\varphi: E \rightarrow F^*$ and $r_\varphi: F \rightarrow E^*$ are semilinear).

The notion of a nondegenerate sesquilinear form is identical to the notion for bilinear forms. For the convenience of the reader, we repeat the definition.

Definition 14.7. A sesquilinear map $\varphi: E \times F \rightarrow K$ is said to be *nondegenerate* iff the following conditions hold:

- (1) For every $u \in E$, if $\varphi(u, v) = 0$ for all $v \in F$, then $u = 0$, and
- (2) For every $v \in F$, if $\varphi(u, v) = 0$ for all $u \in E$, then $v = 0$.

Proposition 14.1 translates into the following proposition. The proof is left as an exercise.

Proposition 14.8. *Given a sesquilinear map $\varphi: E \times F \rightarrow K$, the following properties hold:*

- (a) *The map l_φ is injective iff property (1) of Definition 14.7 holds.*
- (b) *The map r_φ is injective iff property (2) of Definition 14.7 holds.*
- (c) *The sesquilinear form φ is nondegenerate and iff l_φ and r_φ are injective.*
- (d) *If the sesquilinear form φ is nondegenerate and if E and F have finite dimensions, then $\dim(E) = \dim(F)$, and $l_\varphi: \overline{E} \rightarrow F^*$ and $r_\varphi: \overline{F} \rightarrow E^*$ are linear isomorphisms.*

Propositions 14.2 and 14.3 also generalize to sesquilinear forms. We also have the following version of Theorem 14.4, whose proof is left as an exercise.

Theorem 14.9. *Given any sesquilinear form $\varphi: E \times E \rightarrow K$ with $\dim(E) = n$, if φ is Hermitian and K does not have characteristic 2, then there is a basis (e_1, \dots, e_n) of E such that $\varphi(e_i, e_j) = 0$, for all $i \neq j$.*

As in Section 14.1, if E and F are finite-dimensional vector spaces and if (e_1, \dots, e_m) is a basis of E and (f_1, \dots, f_n) is a basis of F then the sesquilinearity of φ yields

$$\varphi\left(\sum_{i=1}^m x_i e_i, \sum_{j=1}^n y_j f_j\right) = \sum_{i=1}^m \sum_{j=1}^n x_i \varphi(e_i, f_j) \overline{y_j}.$$

This shows that φ is completely determined by the $m \times n$ matrix $M = (\varphi(e_i, e_j))$, and in matrix form, we have

$$\varphi(x, y) = x^\top M \overline{y} = y^* M^\top x,$$

where x and \overline{y} are the column vectors associated with $(x_1, \dots, x_m) \in K^m$ and $(\overline{y}_1, \dots, \overline{y}_n) \in K^n$, and $y^* = \overline{y}^\top$. As earlier, we are committing the slight abuse of notation of letting x

denote both the vector $x = \sum_{i=1}^n x_i e_i$ and the column vector associated with (x_1, \dots, x_n) (and similarly for y). We call M the *matrix of φ with respect to the bases (e_1, \dots, e_m) and (f_1, \dots, f_n)* .

If $m = \dim(E) = \dim(F) = n$, then φ is nondegenerate iff M is invertible iff $\det(M) \neq 0$.

Observe that if φ is a Hermitian form ($E = F$) and if K does not have characteristic 2, then by Theorem 14.9, there is a basis of E with respect to which the matrix M representing φ is a diagonal matrix. If $K = \mathbb{C}$, then these entries are real, and this allows us to classify completely the Hermitian forms.

Proposition 14.10. *Given any Hermitian form $\varphi: E \times E \rightarrow \mathbb{C}$ with $\dim(E) = n$, there is a basis (e_1, \dots, e_n) of E such that*

$$\Phi\left(\sum_{i=1}^n x_i e_i\right) = \sum_{i=1}^p x_i^2 - \sum_{i=p+1}^{p+q} x_i^2,$$

with $0 \leq p, q$ and $p + q \leq n$.

The proof of Proposition 14.10 is the same as the real case of Proposition 14.5. Sylvester's inertia law (Proposition 14.6) also holds for Hermitian forms: p and q only depend on φ .

14.3 Orthogonality

In this section, we assume that we are dealing with a sesquilinear form $\varphi: E \times F \rightarrow K$. We allow the automorphism $\lambda \mapsto \bar{\lambda}$ to be the identity, in which case φ is a bilinear form. This way, we can deal with properties shared by bilinear forms and sesquilinear forms in a uniform fashion. Orthogonality is such a property.

Definition 14.8. Given a sesquilinear form $\varphi: E \times F \rightarrow K$, we say that two vectors $u \in E$ and $v \in F$ are *orthogonal* (or *conjugate*) if $\varphi(u, v) = 0$. Two subsets $E' \subseteq E$ and $F' \subseteq F$ are *orthogonal* if $\varphi(u, v) = 0$ for all $u \in E'$ and all $v \in F'$. Given a subspace U of E , the *right orthogonal space* of U , denoted U^\perp , is the subspace of F given by

$$U^\perp = \{v \in F \mid \varphi(u, v) = 0 \text{ for all } u \in U\},$$

and given a subspace V of F , the *left orthogonal space* of V , denoted V^\perp , is the subspace of E given by

$$V^\perp = \{u \in E \mid \varphi(u, v) = 0 \text{ for all } v \in V\}.$$

When E and F are distinct, there is little chance of confusing the right orthogonal subspace U^\perp of a subspace U of E and the left orthogonal subspace V^\perp of a subspace V of F . However, if $E = F$, then $\varphi(u, v) = 0$ *does not necessarily imply* that $\varphi(v, u) = 0$, that is, orthogonality is not necessarily symmetric. Thus, if both U and V are subsets of E , there

is a notational ambiguity if $U = V$. In this case, we may write U^{\perp_r} for the right orthogonal and U^{\perp_l} for the left orthogonal.

The above discussion brings up the following point: When is orthogonality symmetric?

If φ is bilinear, it is shown in E. Artin [3] (and in Jacobson [59]) that orthogonality is symmetric iff either φ is symmetric or φ is alternating ($\varphi(u, u) = 0$ for all $u \in E$).

If φ is sesquilinear, the answer is more complicated. In addition to the previous two cases, there is a third possibility:

$$\varphi(u, v) = \overline{\epsilon \varphi(v, u)} \quad \text{for all } u, v \in E,$$

where ϵ is some nonzero element in K . We say that φ is ϵ -Hermitian. Observe that

$$\varphi(u, u) = \epsilon \bar{\epsilon} \varphi(u, u),$$

so if φ is not alternating, then $\varphi(u, u) \neq 0$ for some u , and we must have $\epsilon \bar{\epsilon} = 1$. The most common cases are

1. $\epsilon = 1$, in which case φ is *Hermitian*, and
2. $\epsilon = -1$, in which case φ is *skew-Hermitian*.

If φ is alternating and K is not of characteristic 2, then the automorphism $\lambda \mapsto \bar{\lambda}$ must be the identity if φ is nonzero. If so, φ is skew-symmetric, so $\epsilon = -1$.

In summary, if φ is either symmetric, alternating, or ϵ -Hermitian, then orthogonality is symmetric, and it makes sense to talk about *the* orthogonal subspace U^\perp of U .

Observe that if φ is ϵ -Hermitian, then

$$r_\varphi = \epsilon l_\varphi.$$

This is because

$$\begin{aligned} l_\varphi(u)(y) &= \overline{\varphi(u, y)} \\ r_\varphi(u)(y) &= \varphi(y, u) \\ &= \overline{\epsilon \varphi(u, y)}, \end{aligned}$$

so $r_\varphi = \epsilon l_\varphi$.

If E and F are finite-dimensional with bases (e_1, \dots, e_m) and (f_1, \dots, f_n) , and if φ is represented by the $m \times n$ matrix M , then φ is ϵ -Hermitian iff

$$M = \epsilon M^*,$$

where $M^* = (\overline{M})^\top$ (as usual). This captures the following kinds of familiar matrices:

1. Symmetric matrices ($\epsilon = 1$)
2. Skew-symmetric matrices ($\epsilon = -1$)
3. Hermitian matrices ($\epsilon = 1$)
4. Skew-Hermitian matrices ($\epsilon = -1$).

Going back to a sesquilinear form $\varphi: E \times F \rightarrow K$, for any subspace U of E , it is easy to check that

$$U \subseteq (U^\perp)^\perp,$$

and that for any subspace V of F , we have

$$V \subseteq (V^\perp)^\perp.$$

For simplicity of notation, we write $U^{\perp\perp}$ instead of $(U^\perp)^\perp$ (and $V^{\perp\perp}$ instead of $(V^\perp)^\perp$). Given any two subspaces U_1 and U_2 of E , if $U_1 \subseteq U_2$, then $U_2^\perp \subseteq U_1^\perp$ (and similarly for any two subspaces $V_1 \subseteq V_2$ of F). As a consequence, it is easy to show that

$$U^\perp = U^{\perp\perp\perp}, \quad V^\perp = V^{\perp\perp\perp}.$$

Observe that φ is nondegenerate iff $E^\perp = \{0\}$ and $F^\perp = \{0\}$. Furthermore, since

$$\begin{aligned} \varphi(u+x, v) &= \varphi(u, v) \\ \varphi(u, v+y) &= \varphi(u, v) \end{aligned}$$

for any $x \in F^\perp$ and any $y \in E^\perp$, we see that we obtain by passing to the quotient a sesquilinear form

$$[\varphi]: (E/F^\perp) \times (F/E^\perp) \rightarrow K$$

which is nondegenerate.

Proposition 14.11. *For any sesquilinear form $\varphi: E \times F \rightarrow K$, the space E/F^\perp is finite-dimensional iff the space F/E^\perp is finite-dimensional; if so, $\dim(E/F^\perp) = \dim(F/E^\perp)$.*

Proof. Since the sesquilinear form $[\varphi]: (E/F^\perp) \times (F/E^\perp) \rightarrow K$ is nondegenerate, the maps $l_{[\varphi]}: (E/F^\perp) \rightarrow (F/E^\perp)^*$ and $r_{[\varphi]}: (F/E^\perp) \rightarrow (E/F^\perp)^*$ are injective. If $\dim(E/F^\perp) = m$, then $\dim(E/F^\perp) = \dim((E/F^\perp)^*)$, so by injectivity of $r_{[\varphi]}$, we have $\dim(F/E^\perp) = \dim((F/E^\perp)^*) \leq m$. A similar reasoning using the injectivity of $l_{[\varphi]}$ applies if $\dim(F/E^\perp) = n$, and we get $\dim(E/F^\perp) = \dim((E/F^\perp)^*) \leq n$. Therefore, $\dim(E/F^\perp) = m$ is finite iff $\dim(F/E^\perp) = n$ is finite, in which case $m = n$. \square

If U is a subspace of a space E , recall that the *codimension* of U is the dimension of E/U , which is also equal to the dimension of any subspace V such that E is a direct sum of U and V ($E = U \oplus V$).

Proposition 14.11 implies the following useful fact.

Proposition 14.12. *Let $\varphi: E \times F \rightarrow K$ be any nondegenerate sesquilinear form. A subspace U of E has finite dimension iff U^\perp has finite codimension in F . If $\dim(U)$ is finite, then $\text{codim}(U^\perp) = \dim(U)$, and $U^{\perp\perp} = U$.*

Proof. Since φ is nondegenerate $E^\perp = \{0\}$ and $F^\perp = \{0\}$, so the first two statements follow from proposition 14.11 applied to the restriction of φ to $U \times F$. Since U^\perp and $U^{\perp\perp}$ are orthogonal, and since $\text{codim}(U^\perp)$ is finite, $\dim(U^{\perp\perp})$ is finite and we have $\dim(U^{\perp\perp}) = \text{codim}(U^\perp) = \dim(U)$. Since $U \subseteq U^{\perp\perp}$, we must have $U = U^{\perp\perp}$. \square

Proposition 14.13. *Let $\varphi: E \times F \rightarrow K$ be any sesquilinear form. Given any two subspaces U and V of E , we have*

$$(U + V)^\perp = U^\perp \cap V^\perp.$$

Furthermore, if φ is nondegenerate and if U and V are finite-dimensional, then

$$(U \cap V)^\perp = U^\perp + V^\perp.$$

Proof. If $w \in (U + V)^\perp$, then $\varphi(u + v, w) = 0$ for all $u \in U$ and all $v \in V$. In particular, with $v = 0$, we have $\varphi(u, w) = 0$ for all $u \in U$, and with $u = 0$, we have $\varphi(v, w) = 0$ for all $v \in V$, so $w \in U^\perp \cap V^\perp$. Conversely, if $w \in U^\perp \cap V^\perp$, then $\varphi(u, w) = 0$ for all $u \in U$ and $\varphi(v, w) = 0$ for all $v \in V$. By bilinearity, $\varphi(u + v, w) = \varphi(u, w) + \varphi(v, w) = 0$, which shows that $w \in (U + V)^\perp$. Therefore, the first identity holds.

Now, assume that φ is nondegenerate and that U and V are finite-dimensional, and let $W = U^\perp + V^\perp$. Using the equation that we just established and the fact that U and V are finite-dimensional, by Proposition 14.12, we get

$$W^\perp = U^{\perp\perp} \cap V^{\perp\perp} = U \cap V.$$

We can apply Proposition 14.11 to the restriction of φ to $U \times W$ (since $U^\perp \subseteq W$ and $W^\perp \subseteq U$), and we get

$$\dim(U/W^\perp) = \dim(U/(U \cap V)) = \dim(W/U^\perp) = \text{codim}(U^\perp) - \text{codim}(W),$$

and since $\text{codim}(U^\perp) = \dim(U)$, we deduce that

$$\dim(U \cap V) = \text{codim}(W).$$

However, by Proposition 14.12, we have $\dim(U \cap V) = \text{codim}((U \cap V)^\perp)$, so $\text{codim}(W) = \text{codim}((U \cap V)^\perp)$, and since $W \subseteq W^{\perp\perp} = (U \cap V)^\perp$, we must have $W = (U \cap V)^\perp$, as claimed. \square

In view of Proposition 14.11, we can make the following definition.

Definition 14.9. Let $\varphi: E \times F \rightarrow K$ be any sesquilinear form. If E/F^\perp and F/E^\perp are finite-dimensional, then their common dimension is called the *rank* of the form φ . If E/F^\perp and F/E^\perp have infinite dimension, we say that φ has infinite rank.

Not surprisingly, the rank of φ is related to the ranks of l_φ and r_φ .

Proposition 14.14. *Let $\varphi: E \times F \rightarrow K$ be any sesquilinear form. If φ has finite rank r , then l_φ and r_φ have the same rank, which is equal to r .*

Proof. Because for every $u \in E$,

$$l_\varphi(u)(y) = \overline{\varphi(u, y)} \quad \text{for all } y \in F,$$

and for every $v \in F$,

$$r_\varphi(v)(x) = \varphi(x, v) \quad \text{for all } x \in E,$$

it is clear that the kernel of $l_\varphi: \overline{E} \rightarrow F^*$ is equal to F^\perp and that, the kernel of $r_\varphi: \overline{F} \rightarrow E^*$ is equal to E^\perp . Therefore, $\text{rank}(l_\varphi) = \dim(\text{Im } l_\varphi) = \dim(E/F^\perp) = r$, and similarly $\text{rank}(r_\varphi) = \dim(\overline{F}/E^\perp) = r$. \square

Remark: If the sesquilinear form φ is represented by the matrix $m \times n$ matrix M with respect to the bases (e_1, \dots, e_m) in E and (f_1, \dots, f_n) in F , it can be shown that the matrix representing l_φ with respect to the bases (e_1, \dots, e_m) and (f_1^*, \dots, f_n^*) is M^* , and that the matrix representing r_φ with respect to the bases (f_1, \dots, f_n) and (e_1^*, \dots, e_m^*) is M . It follows that the rank of φ is equal to the rank of M .

14.4 Adjoint of a Linear Map

Let E_1 and E_2 be two K -vector spaces, and let $\varphi_1: E_1 \times E_1 \rightarrow K$ be a sesquilinear form on E_1 and $\varphi_2: E_2 \times E_2 \rightarrow K$ be a sesquilinear form on E_2 . It is also possible to deal with the more general situation where we have four vector spaces E_1, F_1, E_2, F_2 and two sesquilinear forms $\varphi_1: E_1 \times F_1 \rightarrow K$ and $\varphi_2: E_2 \times F_2 \rightarrow K$, but we will leave this generalization as an exercise. We also assume that l_{φ_1} and r_{φ_1} are bijective, which implies that φ_1 is nondegenerate. This is automatic if the space E_1 is finite dimensional and φ_1 is nondegenerate.

Given any linear map $f: E_1 \rightarrow E_2$, for any fixed $u \in E_2$, we can consider the linear form in E_1^* given by

$$x \mapsto \varphi_2(f(x), u), \quad x \in E_1.$$

Since $r_{\varphi_1}: \overline{E_1} \rightarrow E_1^*$ is bijective, there is a unique $y \in E_1$ (because the vector spaces E_1 and $\overline{E_1}$ only differ by scalar multiplication), so that

$$\varphi_2(f(x), u) = \varphi_1(x, y), \quad \text{for all } x \in E_1.$$

If we denote this unique $y \in E_1$ by $f^{*!}(u)$, then we have

$$\varphi_2(f(x), u) = \varphi_1(x, f^{*!}(u)), \quad \text{for all } x \in E_1, \text{ and all } u \in E_2.$$

Thus, we get a function $f^{*l}: E_2 \rightarrow E_1$. We claim that this function is a linear map. For any $v_1, v_2 \in E_2$, we have

$$\begin{aligned}\varphi_2(f(x), v_1 + v_2) &= \varphi_2(f(x), v_1) + \varphi_2(f(x), v_2) \\ &= \varphi_1(x, f^{*l}(v_1)) + \varphi_1(x, f^{*l}(v_2)) \\ &= \varphi_1(x, f^{*l}(v_1) + f^{*l}(v_2)) \\ &= \varphi_1(x, f^{*l}(v_1 + v_2)),\end{aligned}$$

for all $x \in E_1$. Since r_{φ_1} is injective, we conclude that

$$f^{*l}(v_1 + v_2) = f^{*l}(v_1) + f^{*l}(v_2).$$

For any $\lambda \in K$, we have

$$\begin{aligned}\varphi_2(f(x), \lambda v) &= \overline{\lambda} \varphi_2(f(x), v) \\ &= \overline{\lambda} \varphi_1(x, f^{*l}(v)) \\ &= \varphi_1(x, \lambda f^{*l}(v)) \\ &= \varphi_1(x, f^{*l}(\lambda v)),\end{aligned}$$

for all $x \in E_1$. Since r_{φ_1} is injective, we conclude that

$$f^{*l}(\lambda v) = \lambda f^{*l}(v).$$

Therefore, f^{*l} is linear. We call it the *left adjoint* of f .

Now, for any fixed $u \in E_2$, we can consider the linear form in E_1^* given by

$$x \mapsto \overline{\varphi_2(u, f(x))} \quad x \in E_1.$$

Since $l_{\varphi_1}: \overline{E_1} \rightarrow E_1^*$ is bijective, there is a unique $y \in E_1$ so that

$$\overline{\varphi_2(u, f(x))} = \overline{\varphi_1(y, x)}, \quad \text{for all } x \in E_1.$$

If we denote this unique $y \in E_1$ by $f^{*r}(u)$, then we have

$$\varphi_2(u, f(x)) = \varphi_1(f^{*r}(u), x), \quad \text{for all } x \in E_1, \text{ and all } u \in E_2.$$

Thus, we get a function $f^{*r}: E_2 \rightarrow E_1$. As in the previous situation, it is easy to check that f^{*r} is linear. We call it the *right adjoint* of f . In summary, we make the following definition.

Definition 14.10. Let E_1 and E_2 be two K -vector spaces, and let $\varphi_1: E_1 \times E_1 \rightarrow K$ and $\varphi_2: E_2 \times E_2 \rightarrow K$ be two sesquilinear forms. Assume that l_{φ_1} and r_{φ_1} are bijective, so that φ_1 is nondegenerate. For every linear map $f: E_1 \rightarrow E_2$, there exist unique linear maps $f^{*l}: E_2 \rightarrow E_1$ and $f^{*r}: E_2 \rightarrow E_1$, such that

$$\begin{aligned}\varphi_2(f(x), u) &= \varphi_1(x, f^{*l}(u)), \quad \text{for all } x \in E_1, \text{ and all } u \in E_2 \\ \varphi_2(u, f(x)) &= \varphi_1(f^{*r}(u), x), \quad \text{for all } x \in E_1, \text{ and all } u \in E_2.\end{aligned}$$

The map f^{*l} is called the *left adjoint* of f , and the map f^{*r} is called the *right adjoint* of f .

If E_1 and E_2 are finite-dimensional with bases (e_1, \dots, e_m) and (f_1, \dots, f_n) , then we can work out the matrices A^{*l} and A^{*r} corresponding to the left adjoint f^{*l} and the right adjoint f^{*r} of f . Assuming that f is represented by the $n \times m$ matrix A , φ_1 is represented by the $m \times m$ matrix M_1 , and φ_2 is represented by the $n \times n$ matrix M_2 , we find that

$$\begin{aligned} A^{*l} &= (\overline{M_1})^{-1} A^* \overline{M_2} \\ A^{*r} &= (M_1^\top)^{-1} A^* M_2^\top. \end{aligned}$$

If φ_1 and φ_2 are symmetric bilinear forms, then $f^{*l} = f^{*r}$. This also holds if φ is ϵ -Hermitian. Indeed, since

$$\varphi_2(u, f(x)) = \varphi_1(f^{*r}(u), x),$$

we get

$$\overline{\epsilon \varphi_2(f(x), u)} = \overline{\epsilon \varphi_1(x, f^{*r}(u))},$$

and since $\lambda \mapsto \bar{\lambda}$ is an involution, we get

$$\varphi_2(f(x), u) = \varphi_1(x, f^{*r}(u)).$$

Since we also have

$$\varphi_2(f(x), u) = \varphi_1(x, f^{*l}(u)),$$

we obtain

$$\varphi_1(x, f^{*r}(u)) = \varphi_1(x, f^{*l}(u)) \quad \text{for all } x \in E_1, \text{ and all } u \in E_2,$$

and since φ_1 is nondegenerate, we conclude that $f^{*l} = f^{*r}$. Whenever $f^{*l} = f^{*r}$, we use the simpler notation f^* .

If $f: E_1 \rightarrow E_2$ and $g: E_1 \rightarrow E_2$ are two linear maps, we have the following properties:

$$\begin{aligned} (f + g)^{*l} &= f^{*l} + g^{*l} \\ \text{id}^{*l} &= \text{id} \\ (\lambda f)^{*l} &= \bar{\lambda} f^{*l}, \end{aligned}$$

and similarly for right adjoints. If E_3 is another space, φ_3 is a sesquilinear form on E_3 , and if l_{φ_2} and r_{φ_2} are bijective, then for any linear maps $f: E_1 \rightarrow E_2$ and $g: E_2 \rightarrow E_3$, we have

$$(g \circ f)^{*l} = f^{*l} \circ g^{*l},$$

and similarly for right adjoints. Furthermore, if $E_1 = E_2$ and $\varphi: E \times E \rightarrow K$ is ϵ -Hermitian, for any linear map $f: E \rightarrow E$ (recall that in this case $f^{*l} = f^{*r} = f^*$), we have

$$f^{**} = \epsilon \bar{\epsilon} f.$$

14.5 Isometries Associated with Sesquilinear Forms

The notion of adjoint is a good tool to investigate the notion of isometry between spaces equipped with sesquilinear forms. First, we define metric maps and isometries.

Definition 14.11. If (E_1, φ_1) and (E_2, φ_2) are two pairs of spaces and sesquilinear maps $\varphi_1: E_1 \times E_1 \rightarrow K$ and $\varphi_2: E_2 \times E_2 \rightarrow K$, a *metric map* from (E_1, φ_1) to (E_2, φ_2) is a linear map $f: E_1 \rightarrow E_2$ such that

$$\varphi_1(u, v) = \varphi_2(f(u), f(v)) \quad \text{for all } u, v \in E_1.$$

We say that φ_1 and φ_2 are *equivalent* iff there is a metric map $f: E_1 \rightarrow E_2$ which is bijective. Such a metric map is called an *isometry*.

The problem of classifying sesquilinear forms up to equivalence is an important but very difficult problem. Solving this problem depends intimately on properties of the field K , and a complete answer is only known in a few cases. The problem is easily solved for $K = \mathbb{R}$, $K = \mathbb{C}$. It is also solved for finite fields and for $K = \mathbb{Q}$ (the rationals), but the solution is surprisingly involved!

It is hard to say anything interesting if φ_1 is degenerate and if the linear map f does not have adjoints. The next few propositions make use of natural conditions on φ_1 that yield a useful criterion for being a metric map.

Proposition 14.15. *With the same assumptions as in Definition 14.10, if $f: E_1 \rightarrow E_2$ is a bijective linear map, then we have*

$$\begin{aligned} \varphi_1(x, y) &= \varphi_2(f(x), f(y)) \quad \text{for all } x, y \in E_1 \text{ iff} \\ f^{-1} &= f^{*l} = f^{*r}. \end{aligned}$$

Proof. We have

$$\varphi_1(x, y) = \varphi_2(f(x), f(y))$$

iff

$$\varphi_1(x, y) = \varphi_2(f(x), f(y)) = \varphi_1(x, f^{*l}(f(y)))$$

iff

$$\varphi_1(x, (\text{id} - f^{*l} \circ f)(y)) = 0 \quad \text{for all } x \in E_1 \text{ and all } y \in E_2.$$

Since φ_1 is nondegenerate, we must have

$$f^{*l} \circ f = \text{id},$$

which implies that $f^{-1} = f^{*l}$. similarly,

$$\varphi_1(x, y) = \varphi_2(f(x), f(y))$$

iff

$$\varphi_1(x, y) = \varphi_2(f(x), f(y)) = \varphi_1(f^{*r}(f(x)), y)$$

iff

$$\varphi_1((\text{id} - f^{*r} \circ f)(x), y) = 0 \quad \text{for all } x \in E_1 \text{ and all } y \in E_2.$$

Since φ_1 is nondegenerate, we must have

$$f^{*r} \circ f = \text{id},$$

which implies that $f^{-1} = f^{*r}$. Therefore, $f^{-1} = f^{*l} = f^{*r}$. For the converse, do the computations in reverse. \square

As a corollary, we get the following important proposition.

Proposition 14.16. *If $\varphi: E \times E \rightarrow K$ is a sesquilinear map, and if l_φ and r_φ are bijective, for every bijective linear map $f: E \rightarrow E$, then we have*

$$\begin{aligned} \varphi(f(x), f(y)) &= \varphi(x, y) \quad \text{for all } x, y \in E \text{ iff} \\ f^{-1} &= f^{*l} = f^{*r}. \end{aligned}$$

We also have the following facts.

Proposition 14.17. *(1) If $\varphi: E \times E \rightarrow K$ is a sesquilinear map and if l_φ is injective, then for every linear map $f: E \rightarrow E$, if*

$$\varphi(f(x), f(y)) = \varphi(x, y) \quad \text{for all } x, y \in E, \tag{*}$$

then f is injective.

(2) If E is finite-dimensional and if φ is nondegenerate, then the linear maps $f: E \rightarrow E$ satisfying () form a group. The inverse of f is given by $f^{-1} = f^*$.*

Proof. (1) If $f(x) = 0$, then

$$\varphi(x, y) = \varphi(f(x), f(y)) = \varphi(0, f(y)) = 0 \quad \text{for all } y \in E.$$

Since l_φ is injective, we must have $x = 0$, and thus f is injective.

(2) If E is finite-dimensional, since a linear map satisfying (*) is injective, it is a bijection. By Proposition 14.16, we have $f^{-1} = f^*$. We also have

$$\varphi(f(x), f(y)) = \varphi((f^* \circ f)(x), y) = \varphi(x, y) = \varphi((f \circ f^*)(x), y) = \varphi(f^*(x), f^*(y)),$$

which shows that f^* satisfies (*). If $\varphi(f(x), f(y)) = \varphi(x, y)$ for all $x, y \in E$ and $\varphi(g(x), g(y)) = \varphi(x, y)$ for all $x, y \in E$, then we have

$$\varphi((g \circ f)(x), (g \circ f)(y)) = \varphi(f(x), f(y)) = \varphi(x, y) \quad \text{for all } x, y \in E.$$

Obviously, the identity map id_E satisfies (*). Therefore, the set of linear maps satisfying (*) is a group. \square

The above considerations motivate the following definition.

Definition 14.12. Let $\varphi: E \times E \rightarrow K$ be a sesquilinear map, and assume that E is finite-dimensional and that φ is nondegenerate. A linear map $f: E \rightarrow E$ is an *isometry* of E (with respect to φ) iff

$$\varphi(f(x), f(y)) = \varphi(x, y) \quad \text{for all } x, y \in E.$$

The set of all isometries of E is a group denoted by $\mathbf{Isom}(\varphi)$.

If φ is symmetric, then the group $\mathbf{Isom}(\varphi)$ is denoted $\mathbf{O}(\varphi)$ and called the *orthogonal group* of φ . If φ is alternating, then the group $\mathbf{Isom}(\varphi)$ is denoted $\mathbf{Sp}(\varphi)$ and called the *symplectic group* of φ . If φ is ϵ -Hermitian, then the group $\mathbf{Isom}(\varphi)$ is denoted $\mathbf{U}_\epsilon(\varphi)$ and called the ϵ -*unitary group* of φ . When $\epsilon = 1$, we drop ϵ and just say *unitary group*.

If (e_1, \dots, e_n) is a basis of E , φ is represented by the $n \times n$ matrix M , and f is represented by the $n \times n$ matrix A , then we find that $f \in \mathbf{Isom}(\varphi)$ iff

$$A^* M^\top A = M^\top \quad \text{iff} \quad A^\top M \bar{A} = M,$$

and A^{-1} is given by $A^{-1} = (M^\top)^{-1} A^* M^\top = (\bar{M})^{-1} A^* \bar{M}$.

More specifically, we define the following groups, using the matrices $I_{p,q}$, $J_{m,m}$ and $A_{m,m}$ defined at the end of Section 14.1.

(1) $K = \mathbb{R}$. We have

$$\begin{aligned} \mathbf{O}(n) &= \{A \in \text{Mat}_n(\mathbb{R}) \mid A^\top A = I_n\} \\ \mathbf{O}(p, q) &= \{A \in \text{Mat}_{p+q}(\mathbb{R}) \mid A^\top I_{p,q} A = I_{p,q}\} \\ \mathbf{Sp}(2n, \mathbb{R}) &= \{A \in \text{Mat}_{2n}(\mathbb{R}) \mid A^\top J_{n,n} A = J_{n,n}\} \\ \mathbf{SO}(n) &= \{A \in \text{Mat}_n(\mathbb{R}) \mid A^\top A = I_n, \det(A) = 1\} \\ \mathbf{SO}(p, q) &= \{A \in \text{Mat}_{p+q}(\mathbb{R}) \mid A^\top I_{p,q} A = I_{p,q}, \det(A) = 1\}. \end{aligned}$$

The group $\mathbf{O}(n)$ is the *orthogonal group*, $\mathbf{Sp}(2n, \mathbb{R})$ is the *real symplectic group*, and $\mathbf{SO}(n)$ is the *special orthogonal group*. We can define the group

$$\{A \in \text{Mat}_{2n}(\mathbb{R}) \mid A^\top A_{n,n} A = A_{n,n}\},$$

but it is isomorphic to $\mathbf{O}(n, n)$.

(2) $K = \mathbb{C}$. We have

$$\begin{aligned} \mathbf{U}(n) &= \{A \in \text{Mat}_n(\mathbb{C}) \mid A^* A = I_n\} \\ \mathbf{U}(p, q) &= \{A \in \text{Mat}_{p+q}(\mathbb{C}) \mid A^* I_{p,q} A = I_{p,q}\} \\ \mathbf{Sp}(2n, \mathbb{C}) &= \{A \in \text{Mat}_{2n}(\mathbb{C}) \mid A^* J_{n,n} A = J_{n,n}\} \\ \mathbf{SU}(n) &= \{A \in \text{Mat}_n(\mathbb{C}) \mid A^* A = I_n, \det(A) = 1\} \\ \mathbf{SU}(p, q) &= \{A \in \text{Mat}_{p+q}(\mathbb{C}) \mid A^* I_{p,q} A = I_{p,q}, \det(A) = 1\}. \end{aligned}$$

The group $\mathbf{U}(n)$ is the *unitary group*, $\mathbf{Sp}(2n, \mathbb{C})$ is the *complex symplectic group*, and $\mathbf{SU}(n)$ is the *special unitary group*.

It can be shown that if $A \in \mathbf{Sp}(2n, \mathbb{R})$ or if $A \in \mathbf{Sp}(2n, \mathbb{C})$, then $\det(A) = 1$.

14.6 Totally Isotropic Subspaces. Witt Decomposition

In this section, we deal with ϵ -Hermitian forms, $\varphi: E \times E \rightarrow K$. In general, E may have subspaces U such that $U \cap U^\perp \neq (0)$, or worse, such that $U \subseteq U^\perp$ (that is, φ is zero on U). We will see that such subspaces play a crucial role in the decomposition of E into orthogonal subspaces.

Definition 14.13. Given an ϵ -Hermitian form $\varphi: E \times E \rightarrow K$, a nonzero vector $u \in E$ is said to be *isotropic* if $\varphi(u, u) = 0$. It is convenient to consider 0 to be isotropic. Given any subspace U of E , the subspace $\text{rad}(U) = U \cap U^\perp$ is called the *radical* of U . We say that

- (i) U is *degenerate* if $\text{rad}(U) \neq (0)$ (equivalently if there is some nonzero vector $u \in U$ such that $x \in U^\perp$). Otherwise, we say that U is *nondegenerate*.
- (ii) U is *totally isotropic* if $U \subseteq U^\perp$ (equivalently if the restriction of φ to U is zero).

By definition, the trivial subspace $U = (0)$ ($= \{0\}$) is nondegenerate. Observe that a subspace U is nondegenerate iff the restriction of φ to U is nondegenerate. A degenerate subspace is sometimes called an *isotropic* subspace. Other authors say that a subspace U is *isotropic* if it contains some (nonzero) isotropic vector. A subspace which has no nonzero isotropic vector is often called *anisotropic*. The space of all isotropic vectors is a cone often called the *light cone* (a terminology coming from the theory of relativity). This is not to be confused with the cone of silence (from Get Smart)! It should also be noted that some authors (such as Serre) use the term *isotropic* instead of *totally isotropic*. The apparent lack of standard terminology is almost as bad as in graph theory!

It is clear that any direct sum of pairwise orthogonal totally isotropic subspaces is totally isotropic. Thus, every totally isotropic subspace is contained in some maximal totally isotropic subspace.

First, let us show that in order to study an ϵ -Hermitian form on a space E , it suffices to restrict our attention to nondegenerate forms.

Proposition 14.18. *Given an ϵ -Hermitian form $\varphi: E \times E \rightarrow K$ on E , we have:*

- (a) *If U and V are any two orthogonal subspaces of E , then*

$$\text{rad}(U + V) = \text{rad}(U) + \text{rad}(V).$$

- (b) *$\text{rad}(\text{rad}(E)) = \text{rad}(E)$.*

- (c) *If U is any subspace supplementary to $\text{rad}(E)$, so that*

$$E = \text{rad}(E) \oplus U,$$

then U is nondegenerate, and $\text{rad}(E)$ and U are orthogonal.

Proof. (a) If U and V are orthogonal, then $U \subseteq V^\perp$ and $V \subseteq U^\perp$. We get

$$\begin{aligned} \text{rad}(U + V) &= (U + V) \cap (U + V)^\perp \\ &= (U + V) \cap U^\perp \cap V^\perp \\ &= U \cap U^\perp \cap V^\perp + V \cap U^\perp \cap V^\perp \\ &= U \cap U^\perp + V \cap V^\perp \\ &= \text{rad}(U) + \text{rad}(V). \end{aligned}$$

(b) By definition, $\text{rad}(E) = E^\perp$, and obviously $E = E^{\perp\perp}$, so we get

$$\text{rad}(\text{rad}(E)) = E^\perp \cap E^{\perp\perp} = E^\perp \cap E = E^\perp = \text{rad}(E).$$

(c) If $E = \text{rad}(E) \oplus U$, by definition of $\text{rad}(E)$, the subspaces $\text{rad}(E)$ and U are orthogonal. From (a) and (b), we get

$$\text{rad}(E) = \text{rad}(E) + \text{rad}(U).$$

Since $\text{rad}(U) = U \cap U^\perp \subseteq U$ and since $\text{rad}(E) \oplus U$ is a direct sum, we have a direct sum

$$\text{rad}(E) = \text{rad}(E) \oplus \text{rad}(U),$$

which implies that $\text{rad}(U) = (0)$; that is, U is nondegenerate. \square

Proposition 14.18(c) shows that the restriction of φ to any supplement U of $\text{rad}(E)$ is nondegenerate and φ is zero on $\text{rad}(U)$, so we may restrict our attention to nondegenerate forms.

The following is also a key result.

Proposition 14.19. *Given an ϵ -Hermitian form $\varphi: E \times E \rightarrow K$ on E , if U is a finite-dimensional nondegenerate subspace of E , then $E = U \oplus U^\perp$.*

Proof. By hypothesis, the restriction φ_U of φ to U is nondegenerate, so the semilinear map $r_{\varphi_U}: U \rightarrow U^*$ is injective. Since U is finite-dimensional, r_{φ_U} is actually bijective, so for every $v \in E$, if we consider the linear form in U^* given by $u \mapsto \varphi(u, v)$ ($u \in U$), there is a unique $v_0 \in U$ such that

$$\varphi(u, v_0) = \varphi(u, v) \quad \text{for all } u \in U;$$

that is, $\varphi(u, v - v_0) = 0$ for all $u \in U$, so $v - v_0 \in U^\perp$. It follows that $v = v_0 + v - v_0$, with $v_0 \in U$ and $v - v_0 \in U^\perp$, and since U is nondegenerate $U \cap U^\perp = (0)$, and $E = U \oplus U^\perp$. \square

As a corollary of Proposition 14.19, we get the following result.

Proposition 14.20. *Given an ϵ -Hermitian form $\varphi: E \times E \rightarrow K$ on E , if φ is nondegenerate and if U is a finite-dimensional subspace of E , then $\text{rad}(U) = \text{rad}(U^\perp)$, and the following conditions are equivalent:*

(i) U is nondegenerate.

(ii) U^\perp is nondegenerate.

(iii) $E = U \oplus U^\perp$.

Proof. By definition, $\text{rad}(U^\perp) = U^\perp \cap U^{\perp\perp}$, and since φ is nondegenerate and U is finite-dimensional, $U^{\perp\perp} = U$, so $\text{rad}(U^\perp) = U^\perp \cap U^{\perp\perp} = U^\perp \cap U = (0)$.

By Proposition 14.19, (i) implies (iii). If $E = U \oplus U^\perp$, then $\text{rad}(U) = U \cap U^\perp = (0)$, so U is nondegenerate and (iii) implies (i). Since $\text{rad}(U^\perp) = \text{rad}(U)$, (iii) also implies (ii). Now, if U^\perp is nondegenerate, we have $U^\perp \cap U^{\perp\perp} = (0)$, and since $U \subseteq U^{\perp\perp}$, we get

$$U \cap U^\perp \subseteq U^{\perp\perp} \cap U^\perp = (0),$$

which shows that U is nondegenerate, proving the implication (ii) \implies (i). \square

If E is finite-dimensional, we have the following results.

Proposition 14.21. *Given an ϵ -Hermitian form $\varphi: E \times E \rightarrow K$ on a finite-dimensional space E , if φ is nondegenerate, then for every subspace U of E we have*

(i) $\dim(U) + \dim(U^\perp) = \dim(E)$.

(ii) $U^{\perp\perp} = U$.

Proof. (i) Since φ is nondegenerate and E is finite-dimensional, the semilinear map $l_\varphi: E \rightarrow E^*$ is bijective. By transposition, the inclusion $i: U \rightarrow E$ yields a surjection $r: E^* \rightarrow U^*$ (with $r(f) = f \circ i$ for every $f \in E^*$; the map $f \circ i$ is the restriction of the linear form f to U). It follows that the semilinear map $r \circ l_\varphi: E \rightarrow U^*$ given by

$$(r \circ l_\varphi)(x)(u) = \overline{\varphi(x, u)} \quad x \in E, u \in U$$

is surjective, and its kernel is U^\perp . Thus, we have

$$\dim(U^*) + \dim(U^\perp) = \dim(E),$$

and since $\dim(U) = \dim(U^*)$ because U is finite-dimensional, we get

$$\dim(U) + \dim(U^\perp) = \dim(U^*) + \dim(U^\perp) = \dim(E).$$

(ii) Applying the above formula to U^\perp , we deduce that $\dim(U) = \dim(U^{\perp\perp})$. Since $U \subseteq U^{\perp\perp}$, we must have $U^{\perp\perp} = U$. \square

Remark: We already proved in Proposition 14.12 that if U is finite-dimensional, then $\text{codim}(U^\perp) = \dim(U)$ and $U^{\perp\perp} = U$, but it doesn't hurt to give another proof. Observe that (i) implies that

$$\dim(U) + \dim(\text{rad}(U)) \leq \dim(E).$$

We can now proceed with the Witt decomposition, but before that, we quickly take care of the structure theorem for alternating bilinear forms (the case where $\varphi(u, u) = 0$ for all $u \in E$). For an alternating bilinear form, the space E is totally isotropic. For example in dimension 2, the matrix

$$B = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}$$

defines the alternating form given by

$$\varphi((x_1, y_1), (x_2, y_2)) = x_1y_2 - x_2y_1.$$

This case is surprisingly general.

Proposition 14.22. *Let $\varphi: E \times E \rightarrow K$ be an alternating bilinear form on E . If $u, v \in E$ are two (nonzero) vectors such that $\varphi(u, v) = \lambda \neq 0$, then u and v are linearly independent. If we let $u_1 = \lambda^{-1}u$ and $v_1 = v$, then $\varphi(u_1, v_1) = 1$, and the restriction of φ to the plane spanned by u_1 and v_1 is represented by the matrix*

$$\begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}.$$

Proof. If u and v were linearly dependent, as $u, v \neq 0$, we could write $v = \mu u$ for some $\mu \neq 0$, but then, since φ is alternating, we would have

$$\lambda = \varphi(u, v) = \varphi(u, \mu u) = \mu \varphi(u, u) = 0,$$

contradicting the fact that $\lambda \neq 0$. The rest is obvious. \square

Proposition 14.22 yields a plane spanned by two vectors u_1, v_1 such that $\varphi(u_1, u_1) = \varphi(v_1, v_1) = 0$ and $\varphi(u_1, v_1) = 1$. Such a plane is called a *hyperbolic plane*. If E is finite-dimensional, we obtain the following theorem.

Theorem 14.23. *Let $\varphi: E \times E \rightarrow K$ be an alternating bilinear form on a space E of finite dimension n . Then, there is a direct sum decomposition of E into pairwise orthogonal subspaces*

$$E = W_1 \oplus \cdots \oplus W_r \oplus \text{rad}(E),$$

where each W_i is a hyperbolic plane and $\text{rad}(E) = E^\perp$. Therefore, there is a basis of E of the form

$$(u_1, v_1, \dots, u_r, v_r, w_1, \dots, w_{n-2r}),$$

with respect to which the matrix representing φ is a block diagonal matrix M of the form

$$M = \begin{pmatrix} J & & & 0 \\ & J & & \\ & & \ddots & \\ & & & J \\ 0 & & & & 0_{n-2r} \end{pmatrix},$$

with

$$J = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}.$$

Proof. If $\varphi = 0$, then $E = E^\perp$ and we are done. Otherwise, there are two nonzero vectors $u, v \in E$ such that $\varphi(u, v) \neq 0$, so by Proposition 14.22, we obtain a hyperbolic plane W_2 spanned by two vectors u_1, v_1 such that $\varphi(u_1, v_1) = 1$. The subspace W_1 is nondegenerate (for example, $\det(J) = -1$), so by Proposition 14.20, we get a direct sum

$$E = W_1 \oplus W_1^\perp.$$

By Proposition 14.13, we also have

$$E^\perp = (W_1 \oplus W_1^\perp)^\perp = W_1^\perp \cap W_1^{\perp\perp} = \text{rad}(W_1^\perp).$$

By the induction hypothesis applied to W_1^\perp , we obtain our theorem. \square

The following corollary follows immediately.

Proposition 14.24. *Let $\varphi: E \times E \rightarrow K$ be an alternating bilinear form on a space E of finite dimension n .*

- (1) *The rank of φ is even.*
- (2) *If φ is nondegenerate, then $\dim(E) = n$ is even.*
- (3) *Two alternating bilinear forms $\varphi_1: E_1 \times E_1 \rightarrow K$ and $\varphi_2: E_2 \times E_2 \rightarrow K$ are equivalent iff $\dim(E_1) = \dim(E_2)$ and φ_1 and φ_2 have the same rank.*

The only part that requires a proof is part (3), which is left as an easy exercise.

If φ is nondegenerate, then $n = 2r$, and a basis of E as in Theorem 14.23 is called a *symplectic basis*. The space E is called a *hyperbolic space* (or *symplectic space*).

Observe that if we reorder the vectors in the basis

$$(u_1, v_1, \dots, u_r, v_r, w_1, \dots, w_{n-2r})$$

to obtain the basis

$$(u_1, \dots, u_r, v_1, \dots, v_r, w_1, \dots, w_{n-2r}),$$

then the matrix representing φ becomes

$$\begin{pmatrix} 0 & I_r & 0 \\ -I_r & 0 & 0 \\ 0 & 0 & 0_{n-2r} \end{pmatrix}.$$

This particularly simple matrix is often preferable, especially when dealing with the matrices (symplectic matrices) representing the isometries of φ (in which case $n = 2r$).

We now return to the Witt decomposition. From now on, $\varphi: E \times E \rightarrow K$ is an ϵ -Hermitian form. The following assumption will be needed:

Property (T). For every $u \in E$, there is some $\alpha \in K$ such that $\varphi(u, u) = \alpha + \epsilon\bar{\alpha}$.

Property (T) is always satisfied if φ is alternating, or if K is of characteristic $\neq 2$ and $\epsilon = \pm 1$, with $\alpha = \frac{1}{2}\varphi(u, u)$.

The following (bizarre) technical lemma will be needed.

Lemma 14.25. *Let φ be an ϵ -Hermitian form on E and assume that φ satisfies property (T). For any totally isotropic subspace $U \neq (0)$ of E , for every $x \in E$ not orthogonal to U , and for every $\alpha \in K$, there is some $y \in U$ so that*

$$\varphi(x + y, x + y) = \alpha + \epsilon\bar{\alpha}.$$

Proof. By property (T), we have $\varphi(x, x) = \beta + \epsilon\bar{\beta}$ for some $\beta \in K$. For any $y \in U$, since φ is ϵ -Hermitian, $\varphi(y, x) = \epsilon\overline{\varphi(x, y)}$, and since U is totally isotropic $\varphi(y, y) = 0$, so we have

$$\begin{aligned} \varphi(x + y, x + y) &= \varphi(x, x) + \varphi(x, y) + \varphi(y, x) + \varphi(y, y) \\ &= \beta + \epsilon\bar{\beta} + \varphi(x, y) + \overline{\epsilon\varphi(x, y)} \\ &= \beta + \varphi(x, y) + \epsilon(\beta + \varphi(x, y)). \end{aligned}$$

Since x is not orthogonal to U , the function $y \mapsto \varphi(x, y) + \beta$ is not the constant function. Consequently, this function takes the value α for some $y \in U$, which proves the lemma. \square

Definition 14.14. Let φ be an ϵ -Hermitian form on E . A *Witt decomposition* of E is a triple (U, U', W) , such that

- (i) $E = U \oplus U' \oplus W$ (a direct sum)
- (ii) U and U' are totally isotropic
- (iii) W is nondegenerate and orthogonal to $U \oplus U'$.

Furthermore, if E is finite-dimensional, then $\dim(U) = \dim(U')$ and in a suitable basis, the matrix representing φ is of the form

$$\begin{pmatrix} 0 & A & 0 \\ \epsilon \overline{A} & 0 & 0 \\ 0 & 0 & B \end{pmatrix}$$

We say that φ is a *neutral form* if it is nondegenerate, E is finite-dimensional, and if $W = (0)$.

Sometimes, we use the notation $U_1 \overset{\perp}{\oplus} U_2$ to indicate that in a direct sum $U_1 \oplus U_2$, the subspaces U_1 and U_2 are orthogonal. Then, in Definition 14.14, we can write that $E = (U \oplus U') \overset{\perp}{\oplus} W$.

As a warm up for Proposition 14.27, we prove an analog of Proposition 14.22 in the case of a symmetric bilinear form.

Proposition 14.26. *Let $\varphi: E \times E \rightarrow K$ be a nondegenerate symmetric bilinear form with K a field of characteristic different from 2. For any nonzero isotropic vector u , there is another nonzero isotropic vector v such that $\varphi(u, v) = 2$, and u and v are linearly independent. In the basis $(u, v/2)$, the restriction of φ to the plane spanned by u and $v/2$ is of the form*

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

Proof. Since φ is nondegenerate, there is some nonzero vector z such that (rescaling z if necessary) $\varphi(u, z) = 1$. If

$$v = 2z - \varphi(z, z)u,$$

then since $\varphi(u, u) = 0$ and $\varphi(u, z) = 1$, note that

$$\varphi(u, v) = \varphi(u, 2z - \varphi(z, z)u) = 2\varphi(u, z) - \varphi(z, z)\varphi(u, u) = 2,$$

and

$$\begin{aligned} \varphi(v, v) &= \varphi(2z - \varphi(z, z)u, 2z - \varphi(z, z)u) \\ &= 4\varphi(z, z) - 4\varphi(z, z)\varphi(u, z) + \varphi(z, z)^2\varphi(u, u) \\ &= 4\varphi(z, z) - 4\varphi(z, z) = 0. \end{aligned}$$

If u and z were linearly dependent, as $u, z \neq 0$, we could write $z = \mu u$ for some $\mu \neq 0$, but then, we would have

$$\varphi(u, z) = \varphi(u, \mu u) = \mu\varphi(u, u) = 0,$$

contradicting the fact that $\varphi(u, z) \neq 0$. Then u and $v = 2z - \varphi(z, z)u$ are also linearly independent, since otherwise z could be expressed as a multiple of u . The rest is obvious. \square

Proposition 14.26 yields a plane spanned by two vectors u_1, v_1 such that $\varphi(u_1, u_1) = \varphi(v_1, v_1) = 0$ and $\varphi(u_1, v_1) = 1$. Such a plane is called an *Artinian plane*. Proposition 14.26 also shows that nonzero isotropic vectors come in pair.

Remark: Some authors refer to the above plane as a *hyperbolic plane*. Berger (and others) point out that this terminology is undesirable because the notion of hyperbolic plane already exists in differential geometry and refers to a very different object.

We leave it as an exercise to figure out that the group of isometries of the Artinian plane, the set of all 2×2 matrices A such that

$$A^\top \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} A = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix},$$

consists of all matrices of the form

$$\begin{pmatrix} \lambda & 0 \\ 0 & \lambda^{-1} \end{pmatrix} \quad \text{or} \quad \begin{pmatrix} 0 & \lambda \\ \lambda^{-1} & 0 \end{pmatrix}, \quad \lambda \in K - \{0\}.$$

In particular, if $K = \mathbb{R}$, then this group denoted $\mathbf{O}(1, 1)$ has four connected components.

The first step in showing the existence of a Witt decomposition is this.

Proposition 14.27. *Let φ be an ϵ -Hermitian form on E , assume that φ is nondegenerate and satisfies property (T), and let U be any totally isotropic subspace of E of finite dimension $\dim(U) = r$.*

- (1) *If U' is any totally isotropic subspace of dimension r and if $U' \cap U^\perp = (0)$, then $U \oplus U'$ is nondegenerate, and for any basis (u_1, \dots, u_r) of U , there is a basis (u'_1, \dots, u'_r) of U' such that $\varphi(u_i, u'_j) = \delta_{ij}$, for all $i, j = 1, \dots, r$.*
- (2) *If W is any totally isotropic subspace of dimension at most r and if $W \cap U^\perp = (0)$, then there exists a totally isotropic subspace U' with $\dim(U') = r$ such that $W \subseteq U'$ and $U' \cap U^\perp = (0)$.*

Proof. (1) Let φ' be the restriction of φ to $U \times U'$. Since $U' \cap U^\perp = (0)$, for any $v \in U'$, if $\varphi(u, v) = 0$ for all $u \in U$, then $v = 0$. Thus, φ' is nondegenerate (we only have to check on the left since φ is ϵ -Hermitian). Then, the assertion about bases follows from the version of Proposition 14.3 for sesquilinear forms. Since U is totally isotropic, $U \subseteq U^\perp$, and since $U' \cap U^\perp = (0)$, we must have $U' \cap U = (0)$, which show that we have a direct sum $U \oplus U'$.

It remains to prove that $U + U'$ is nondegenerate. Observe that

$$H = (U + U') \cap (U + U')^\perp = (U + U') \cap U^\perp \cap U'^\perp.$$

Since U is totally isotropic, $U \subseteq U^\perp$, and since $U' \cap U^\perp = (0)$, we have

$$(U + U') \cap U^\perp = (U \cap U^\perp) + (U' \cap U^\perp) = U + (0) = U,$$

thus $H = U \cap U'^\perp$. Since φ' is nondegenerate, $U \cap U'^\perp = (0)$, so $H = (0)$ and $U + U'$ is nondegenerate.

(2) We proceed by descending induction on $s = \dim(W)$. The base case $s = r$ is trivial. For the induction step, it suffices to prove that if $s < r$, then there is a totally isotropic subspace W' containing W such that $\dim(W') = s + 1$ and $W' \cap U^\perp = (0)$.

Since $s = \dim(W) < \dim(U)$, the restriction of φ to $U \times W$ is degenerate. Since $W \cap U^\perp = (0)$, we must have $U \cap W^\perp \neq (0)$. We claim that

$$W^\perp \not\subseteq W + U^\perp.$$

If we had

$$W^\perp \subseteq W + U^\perp,$$

then because U and W are finite-dimensional and φ is nondegenerate, by Proposition 14.12, $U^{\perp\perp} = U$ and $W^{\perp\perp} = W$, so by taking orthogonals, $W^\perp \subseteq W + U^\perp$ would yield

$$(W + U^\perp)^\perp \subseteq W^{\perp\perp},$$

that is,

$$W^\perp \cap U \subseteq W,$$

thus $W^\perp \cap U \subseteq W \cap U$, and since U is totally isotropic, $U \subseteq U^\perp$, which yields

$$W^\perp \cap U \subseteq W \cap U \subseteq W \cap U^\perp = (0),$$

contradicting the fact that $U \cap W^\perp \neq (0)$.

Therefore, there is some $u \in W^\perp$ such that $u \notin W + U^\perp$. Since $U \subseteq U^\perp$, we can add to u any vector $z \in W^\perp \cap U \subseteq U^\perp$ so that $u + z \in W^\perp$ and $u + z \notin W + U^\perp$ (if $u + z \in W + U^\perp$, since $z \in U^\perp$, then $u \in W + U^\perp$, a contradiction). Since $W^\perp \cap U \neq (0)$ is totally isotropic and $u \notin W + U^\perp = (W^\perp \cap U)^\perp$, we can invoke Lemma 14.25 to find a $z \in W^\perp \cap U$ such that $\varphi(u + z, u + z) = 0$. If we write $x = u + z$, then $x \notin W + U^\perp$, so $W' = W + Kx$ is a totally isotropic subspace of dimension $s + 1$. Furthermore, we claim that $W' \cap U^\perp = (0)$.

Otherwise, we would have $y = w + \lambda x \in U^\perp$, for some $w \in W$ and some $\lambda \in K$, and then we would have $\lambda x = -w + y \in W + U^\perp$. If $\lambda \neq 0$, then $x \in W + U^\perp$, a contradiction. Therefore, $\lambda = 0$, $y = w$, and since $y \in U^\perp$ and $w \in W$, we have $y \in W \cap U^\perp = (0)$, which means that $y = 0$. Therefore, W' is the required subspace and this completes the proof. \square

Here are some consequences of Proposition 14.27. If we set $W = (0)$ in Proposition 14.27(2), then we get:

Proposition 14.28. *Let φ be an ϵ -Hermitian form on E which is nondegenerate and satisfies property (T). For any totally isotropic subspace U of E of finite dimension r , there exists a totally isotropic subspace U' of dimension r such that $U \cap U' = (0)$ and $U \oplus U'$ is nondegenerate.*

Proposition 14.29. *Any two ϵ -Hermitian neutral forms satisfying property (T) defined on spaces of the same dimension are equivalent.*

Note that under the conditions of Proposition 14.28, $(U, U', (U \oplus U')^\perp)$ is a Witt decomposition for E . By Proposition 14.27(1), the block A in the matrix of φ is the identity matrix.

The following proposition shows that every subspace U of E can be embedded into a nondegenerate subspace.

Proposition 14.30. *Let φ be an ϵ -Hermitian form on E which is nondegenerate and satisfies property (T). For any subspace U of E of finite dimension, if we write*

$$U = V \oplus^\perp W,$$

for some orthogonal complement W of $V = \text{rad}(U)$, and if we let $r = \dim(\text{rad}(U))$, then there exists a totally isotropic subspace V' of dimension r such that $V \cap V' = (0)$, and $(V \oplus V') \oplus^\perp W = V' \oplus U$ is nondegenerate. Furthermore, any isometry f from U into another space (E', φ') where φ' is an ϵ -Hermitian form satisfying the same assumptions as φ can be extended to an isometry on $(V \oplus V') \oplus^\perp W$.

Proof. Since W is nondegenerate, W^\perp is also nondegenerate, and $V \subseteq W^\perp$. Therefore, we can apply Proposition 14.28 to the restriction of φ to W^\perp and to V to obtain the required V' . We know that $V \oplus V'$ is nondegenerate and orthogonal to W , which is also nondegenerate, so $(V \oplus V') \oplus^\perp W = V' \oplus U$ is nondegenerate.

We leave the second statement about extending f as an exercise (use the fact that $f(U) = f(V) \oplus^\perp f(W)$, where $V_1 = f(V)$ is totally isotropic of dimension r , to find another totally isotropic subspace V'_1 of dimension r such that $V_1 \cap V'_1 = (0)$ and $V_1 \oplus V'_1$ is orthogonal to $f(W)$). \square

The subspace $(V \oplus V') \oplus^\perp W = V' \oplus U$ is often called a *nondegenerate completion* of U . The subspace $V \oplus V'$ is called an *Artinian space*. Proposition 14.27 show that $V \oplus V'$ has a basis $(u_1, v_1, \dots, u_r, v_r)$ consisting of vectors $u_i \in V$ and $v_j \in V'$ such that $\varphi(u_i, u_j) = \delta_{ij}$. The subspace spanned by (u_i, v_i) is an Artinian plane, so $V \oplus V'$ it is the orthogonal direct sum of r Artinian planes. Such a space is often denoted by Ar_{2r} .

We now sharpen Proposition 14.27.

Theorem 14.31. *Let φ be an ϵ -Hermitian form on E which is nondegenerate and satisfies property (T). Let U_1 and U_2 be two totally isotropic maximal subspaces of E , with U_1 or U_2 of finite dimension. Write $U = U_1 \cap U_2$, let S_1 be a supplement of U in U_1 and S_2 be a supplement of U in U_2 (so that $U_1 = U \oplus S_1$, $U_2 = U \oplus S_2$), and let $S = S_1 + S_2$. Then, there exist two subspaces W and D of E such that:*

- (a) The subspaces S , $U + W$, and D are nondegenerate and pairwise orthogonal.
- (b) We have a direct sum $E = S \overset{\perp}{\oplus} (U \oplus W) \overset{\perp}{\oplus} D$.
- (c) The subspace D contains no nonzero isotropic vector (D is anisotropic).
- (d) The subspace W is totally isotropic.

Furthermore, U_1 and U_2 are both finite dimensional, and we have $\dim(U_1) = \dim(U_2)$, $\dim(W) = \dim(U)$, $\dim(S_1) = \dim(S_2)$, and $\text{codim}(D) = 2 \dim(F_1)$.

Proof. First observe that if X is a totally isotropic maximal subspace of E , then any isotropic vector $x \in E$ orthogonal to X must belong to X , since otherwise, $X + Kx$ would be a totally isotropic subspace strictly containing X , contradicting the maximality of X . As a consequence, if x_i is any isotropic vector such that $x_i \in U_i^\perp$ (for $i = 1, 2$), then $x_i \in U_i$.

We claim that

$$S_1 \cap S_2^\perp = (0) \quad \text{and} \quad S_2 \cap S_1^\perp = (0).$$

Assume that $y \in S_1$ is orthogonal to S_2 . Since $U_1 = U \oplus S_1$ and U_1 is totally isotropic, y is orthogonal to U_1 , and thus orthogonal to U , so that y is orthogonal to $U_2 = U \oplus S_2$. Since $S_1 \subseteq U_1$ and U_1 is totally isotropic, y is an isotropic vector orthogonal to U_2 , which by a previous remark implies that $y \in U_2$. Then, since $S_1 \subseteq U_1$ and $U \oplus S_1$ is a direct sum, we have

$$y \in S_1 \cap U_2 = S_1 \cap U_1 \cap U_2 = S_1 \cap U = (0).$$

Therefore $S_1 \cap S_2^\perp = (0)$. A similar proof show that $S_2 \cap S_1^\perp = (0)$. If U_1 is finite-dimensional (the case where U_2 is finite-dimensional is similar), then S_1 is finite-dimensional, so by Proposition 14.12, S_1^\perp has finite codimension. Since $S_2 \cap S_1^\perp = (0)$, and since any supplement of S_1^\perp has finite dimension, we must have

$$\dim(S_2) \leq \text{codim}(S_1^\perp) = \dim(S_1).$$

By a similar argument, $\dim(S_1) \leq \dim(S_2)$, so we have

$$\dim(S_1) = \dim(S_2).$$

By Proposition 14.27(1), we conclude that $S = S_1 + S_2$ is nondegenerate.

By Proposition 14.20, the subspace $N = S^\perp = (S_1 + S_2)^\perp$ is nondegenerate. Since $U_1 = U \oplus S_1$, $U_2 = U \oplus S_2$, and U_1, U_2 are totally isotropic, U is orthogonal to S_1 and to S_2 , so $U \subseteq N$. Since U is totally isotropic, by Proposition 14.28 applied to N , there is a totally isotropic subspace W of N such that $\dim(W) = \dim(U)$, $U \cap W = (0)$, and $U + W$ is nondegenerate. Consequently, (d) is satisfied by W .

To satisfy (a) and (b), we pick D to be the orthogonal of $U \oplus W$ in N . Then, $N = (U \oplus W) \overset{\perp}{\oplus} D$ and $E = S \overset{\perp}{\oplus} N$, so $E = S \overset{\perp}{\oplus} (U \oplus W) \overset{\perp}{\oplus} D$.

As to (c), since D is orthogonal $U \oplus W$, D is orthogonal to U , and since $D \subseteq N$ and N is orthogonal to $S_1 + S_2$, D is orthogonal to S_1 , so D is orthogonal to $U_1 = U \oplus S_1$. If $y \in D$ is any isotropic vector, since $y \in U_1^\perp$, by a previous remark, $y \in U_1$, so $y \in D \cap U_1$. But, $D \subseteq N$ with $N \cap (S_1 + S_2) = (0)$, and $D \cap (U + W) = (0)$, so $D \cap (U + S_1) = D \cap U_1 = (0)$, which yields $y = 0$. The statements about dimensions are easily obtained. \square

We obtain the following corollaries.

Theorem 14.32. *Let φ be an ϵ -Hermitian form on E which is nondegenerate and satisfies property (T).*

- (1) *Any two totally isotropic maximal spaces of finite dimension have the same dimension.*
- (2) *For any totally isotropic maximal subspace U of finite dimension r , there is another totally isotropic maximal subspace U' of dimension r such that $U \cap U' = (0)$, and $U \oplus U'$ is nondegenerate. Furthermore, if $D = (U \oplus U')^\perp$, then (U, U', D) is a Witt decomposition of E , and there are no nonzero isotropic vectors in D (D is anisotropic).*
- (3) *If E has finite dimension $n \geq 1$, then E has a Witt decomposition (U, U', D) as in (2). There is a basis of E such that*

- (a) *if φ is alternating ($\epsilon = -1$ and $\lambda = \bar{\lambda}$ for all $\lambda \in K$), then $n = 2m$ and φ is represented by a matrix of the form*

$$\begin{pmatrix} 0 & I_m \\ -I_m & 0 \end{pmatrix}$$

- (b) *if φ is symmetric ($\epsilon = +1$ and $\lambda = \bar{\lambda}$ for all $\lambda \in K$), then φ is represented by a matrix of the form*

$$\begin{pmatrix} 0 & I_r & 0 \\ I_r & 0 & 0 \\ 0 & 0 & P \end{pmatrix},$$

where either $n = 2r$ and P does not occur, or $n > 2r$ and P is a definite symmetric matrix.

- (c) *if φ is ϵ -Hermitian (the involutive automorphism $\lambda \mapsto \bar{\lambda}$ is not the identity), then φ is represented by a matrix of the form*

$$\begin{pmatrix} 0 & I_r & 0 \\ \epsilon I_r & 0 & 0 \\ 0 & 0 & P \end{pmatrix},$$

where either $n = 2r$ and P does not occur, or $n > 2r$ and P is a definite matrix such that $P^ = \epsilon P$.*

Proof. Part (1) follows from Theorem 14.31. By Proposition 14.28, we obtain a totally isotropic subspace U' of dimension r such that $U \cap U' = (0)$. By applying Theorem 14.31 to $U_1 = U$ and $U_2 = U'$, we get $U = W = (0)$, which proves (2). Part (3) is an immediate consequence of (2). \square

As a consequence of Theorem 14.32, we make the following definition.

Definition 14.15. Let E be a vector space of finite dimension n , and let φ be an ϵ -Hermitian form on E which is nondegenerate and satisfies property (T). The *index* (or *Witt index*) ν of φ , is the common dimension of all totally isotropic maximal subspaces of E . We have $2\nu \leq n$.

Neutral forms only exist if n is even, in which case, $\nu = n/2$. Forms of index $\nu = 0$ have no nonzero isotropic vectors. When $K = \mathbb{R}$, this is satisfied by positive definite or negative definite symmetric forms. When $K = \mathbb{C}$, this is satisfied by positive definite or negative definite Hermitian forms. The vector space of a neutral Hermitian form ($\epsilon = +1$) is an Artinian space, and the vector space of a neutral alternating form is a hyperbolic space.

If the field K is algebraically closed, we can describe all nondegenerate quadratic forms.

Proposition 14.33. *If K is algebraically closed and E has dimension n , then for every nondegenerate quadratic form Φ , there is a basis (e_1, \dots, e_n) such that Φ is given by*

$$\Phi\left(\sum_{i=1}^n x_i e_i\right) = \begin{cases} \sum_{i=1}^m x_i x_{m+i} & \text{if } n = 2m \\ \sum_{i=1}^m x_i x_{m+i} + x_{2m+1}^2 & \text{if } n = 2m + 1. \end{cases}$$

Proof. We work with the polar form φ of Φ . Let U_1 and U_2 be some totally isotropic subspaces such that $U_1 \cap U_2 = (0)$ given by Theorem 14.32, and let q be their common dimension. Then, $W = U = (0)$. Since we can pick bases (e_1, \dots, e_q) in U_1 and (e_{q+1}, \dots, e_{2q}) in U_2 such that $\varphi(e_i, e_{i+q}) = 0$, for $i, j = 1, \dots, q$, it suffices to prove that $\dim(D) \leq 1$. If $x, y \in D$ with $x \neq 0$, from the identity

$$\Phi(y - \lambda x) = \Phi(y) - \lambda\varphi(x, y) + \lambda^2\Phi(x)$$

and the fact that $\Phi(x) \neq 0$ since $x \in D$ and $x \neq 0$, we see that the equation $\Phi(y - \lambda x) = 0$ has at least one solution. Since $\Phi(z) \neq 0$ for every nonzero $z \in D$, we get $y = \lambda x$, and thus $\dim(D) \leq 1$, as claimed. \square

We also have the following proposition which has applications in number theory.

Proposition 14.34. *If Φ is any nondegenerate quadratic form such that there is some nonzero vector $x \in E$ with $\Phi(x) = 0$, then for every $\alpha \in K$, there is some $y \in E$ such that $\Phi(y) = \alpha$.*

The proof is left as an exercise. We now turn to the Witt extension theorem.

14.7 Witt's Theorem

Witt's theorem was referred to as a “scandal” by Emil Artin. What he meant by this is that one had to wait until 1936 (Witt [115]) to formulate and prove a theorem at once so simple in its statement and underlying concepts, and so useful in various domains (geometry, arithmetic of quadratic forms).¹

Besides Witt's original proof (Witt [115]), Chevalley's proof [22] seems to be the “best” proof that applies to the symmetric as well as the skew-symmetric case. The proof in Bourbaki [13] is based on Chevalley's proof, and so are a number of other proofs. This is the one we follow (slightly reorganized). In the symmetric case, Serre's exposition is hard to beat (see Serre [97], Chapter IV).

Theorem 14.35. (Witt, 1936) *Let E and E' be two finite-dimensional spaces respectively equipped with two nondegenerate ϵ -Hermitian forms φ and φ' satisfying condition (T), and assume that there is an isometry between (E, φ) and (E', φ') . For any subspace U of E , every injective metric linear map f from U into E' extends to an isometry from E to E' .*

Proof. Since (E, φ) and (E', φ') are isometric, we may assume that $E' = E$ and $\varphi' = \varphi$ (if $h: E \rightarrow E'$ is an isometry, then $h^{-1} \circ f$ is an injective metric map from U into E . The details are left to the reader). We begin with the following observation. If U_1 and U_2 are two subspaces of E such that $U_1 \cap U_2 = (0)$ and if we have metric linear maps $f_1: U_1 \rightarrow E$ and $f_2: U_2 \rightarrow E$ such that

$$\varphi(f_1(u_1), f_2(u_2)) = \varphi(u_1, u_2) \quad \text{for } u_i \in U_i \ (i = 1, 2), \quad (*)$$

then the linear map $f: U_1 \oplus U_2 \rightarrow E$ given by $f(u_1 + u_2) = f_1(u_1) + f_2(u_2)$ extends f_1 and f_2 and is metric. Indeed, since f_1 and f_2 are metric and using $(*)$, we have

$$\begin{aligned} \varphi(f_1(u_1) + f_2(u_2), f_1(v_1) + f_2(v_2)) &= \varphi(f_1(u_1), f_1(v_1)) + \varphi(f_1(u_1), f_2(v_2)) \\ &\quad + \varphi(f_2(u_2), f_1(v_1)) + \varphi(f_2(u_2), f_2(v_2)) \\ &= \varphi(u_1, v_1) + \varphi(u_1, v_2) + \varphi(u_2, v_1) + \varphi(u_2, v_2) \\ &= \varphi(u_1 + u_2, v_2 + v_2). \end{aligned}$$

Furthermore, if f_1 and f_2 are injective, then so is f .

We now proceed by induction on the dimension r of U . The case $r = 0$ is trivial. For the induction step, $r \geq 1$ so $U \neq (0)$, and let H be any hyperplane in U . Let $f: U \rightarrow E$ be an injective metric linear map. By the induction hypothesis, the restriction f_0 of f to H extends to an isometry g_0 of E . If g_0 extends f , we are done. Otherwise, H is the subspace of elements of U left fixed by $g_0^{-1} \circ f$. If the theorem holds in this situation, namely the

¹Curiously, some references to Witt's paper claim its date of publication to be 1936, but others say 1937. The answer to this mystery is that Volume 176 of *Crelle Journal* was published in four issues. The cover page of volume 176 mentions the year 1937, but Witt's paper is dated May 1936. This is not the only paper of Witt appearing in this volume!

subspace of U left fixed by f is a hyperplane H in U , then we have an isometry g_1 of E extending $g_0^{-1} \circ f$, and $g_0 \circ g_1$ is an isometry of E extending f . Therefore, we are reduced to the following situation:

Case (H). The subspace of U left fixed by f is a hyperplane H in U .

In this case, the set $D = \{f(u) - u \mid u \in U\}$ is a line in U (a one-dimensional subspace). For all $u, v \in U$, we have

$$\varphi(f(u), f(v) - v) = \varphi(f(u), f(v)) - \varphi(f(u), v) = \varphi(u, v) - \varphi(f(u), v) = \varphi(u - f(u), v),$$

that is

$$\varphi(f(u), f(v) - v) = \varphi(u - f(u), v) \quad \text{for all } u, v \in U, \quad (**)$$

and if $u \in H$, which means that $f(u) = u$, we get $u \in D^\perp$. Therefore, $H \subseteq D^\perp$. Since φ is nondegenerate, we have $\dim(D) + \dim(D^\perp) = \dim(E)$, and since $\dim(D) = 1$, the subspace D^\perp is a hyperplane in E .

Hypothesis (V). We can find a subspace V of E orthogonal to D and such that $V \cap U = V \cap f(U) = (0)$.

Then, we have

$$\varphi(f(u), v) = \varphi(u, v) \quad \text{for all } u \in U \text{ and all } v \in V,$$

since $\varphi(f(u), v) - \varphi(u, v) = \varphi(f(u) - u, v) = 0$, with $f(u) - u \in D$ and $v \in V$ orthogonal to D . By the remark at the beginning of the proof, with $f_1 = f$ and f_2 the inclusion of V into E , we can extend f to an injective metric map on $U \oplus V$ leaving all vectors in V fixed. In this case, the set $\{f(w) - w \mid w \in U \oplus V\}$ is still the line D . We show below that the fact that f can be extended to $U \oplus V$ implies that f can be extended to the whole of E .

We are reduced to proving that a subspace V as above exists. We distinguish between two cases.

Case (a). $U \not\subseteq D^\perp$.

In this case, formula (**) show that $f(U)$ is not contained in D^\perp (check this!). Consequently,

$$U \cap D^\perp = f(U) \cap D^\perp = H.$$

We can pick V to be any supplement of H in D^\perp , and the above formula shows that $V \cap U = V \cap f(U) = (0)$. Since $U \oplus V$ contains the hyperplane D^\perp (since $D^\perp = H \oplus V$ and $H \subseteq U$), and $U \oplus V \neq D^\perp$ (since U is not contained in D^\perp and $V \subseteq D^\perp$), we must have $E = U \oplus V$, and as we showed as a consequence of hypothesis (V), f can be extended to an isometry of $U \oplus V = E$.

Case (b). $U \subseteq D^\perp$.

In this case, formula (**) shows that $f(U) \subseteq D^\perp$ so $U + f(U) \subseteq D^\perp$, and since $D = \{f(u) - u \mid u \in U\}$, we have $D \subseteq D^\perp$; that is, the line D is isotropic.

We show that case (b) can be reduced to the situation where $U = D^\perp$ and f is an isometry of U . For this, we show that there exists a subspace V of D^\perp , such that

$$D^\perp = U \oplus V = f(U) \oplus V.$$

This is obvious if $U = f(U)$. Otherwise, let $x \in U$ with $x \notin H$, and let $y \in f(U)$ with $y \notin H$. Since $f(H) = H$ (pointwise), f is injective, and H is a hyperplane in U , we have

$$U = H \oplus Kx, \quad f(U) = H \oplus Ky.$$

We claim that $x + y \notin U$. Otherwise, since $y = x + y - x$, with $x + y, x \in U$ and since $y \in f(U)$, we would have $y \in U \cap f(U) = H$, a contradiction. Similarly, $x + y \notin f(U)$. It follows that

$$U + f(U) = U \oplus K(x + y) = f(U) \oplus K(x + y).$$

Now, pick W to be any supplement of $U + f(U)$ in D^\perp so that $D^\perp = (U + f(U)) \oplus W$, and let

$$V = K(x + y) + W.$$

Then, since $x \in U, y \in f(U), W \subseteq D^\perp$, and $U + f(U) \subseteq D^\perp$, we have $V \subseteq D^\perp$. We also have

$$U \oplus V = U \oplus K(x + y) \oplus W = (U + f(U)) \oplus W = D^\perp$$

and

$$f(U) \oplus V = f(U) \oplus K(x + y) \oplus W = (U + f(U)) \oplus W = D^\perp,$$

so as we showed as a consequence of hypothesis (V), f can be extended to an isometry of the hyperplane D^\perp , and D is still the line $\{f(w) - w \mid w \in U \oplus V\}$.

The above argument shows that we are reduced to the situation where $U = D^\perp$ is a hyperplane in E and f is an isometry of U . If we pick any $v \notin U$, then $E = U \oplus Kv$, and if we can find some $v_1 \in E$ such that

$$\begin{aligned} \varphi(f(u), v_1) &= \varphi(u, v) \quad \text{for all } u \in U \\ \varphi(v_1, v_1) &= \varphi(v, v), \end{aligned}$$

then as we showed at the beginning of the proof, we can extend f to a metric map g of $U + Kv = E$ such that $g(v) = v_1$.

To find v_1 , let us prove that for every $v \in E$, there is some $v' \in E$ such that

$$\varphi(f(u), v') = \varphi(u, v) \quad \text{for all } u \in U. \quad (\dagger)$$

This is because the linear form $u \mapsto \varphi(f^{-1}(u), v)$ ($u \in U$) is the restriction of a linear form $\psi \in E^*$, and since φ is nondegenerate, there is some (unique) $v' \in E$, such that

$$\psi(x) = \varphi(x, v') \quad \text{for all } x \in E,$$

which implies that

$$\varphi(u, v') = \varphi(f^{-1}(u), v) \quad \text{for all } u \in U,$$

and since f is an automorphism of U , that (\dagger) holds. Furthermore, observe that formula (\dagger) still holds if we add to v' a vector y in D , since $f(U) = U = D^\perp$. Therefore, for any $v_1 = v' + y$ with $y \in D$, if we extend f to a linear map of E by setting $g(v) = v_1$, then by (\dagger) we have

$$\varphi(g(u), g(v)) = \varphi(u, v) \quad \text{for all } u \in U.$$

We still need to pick $y \in D$ so that $v_1 = v' + y$ satisfies $\varphi(v_1, v_1) = \varphi(v, v)$. However, since $v \notin U = D^\perp$, the vector v is not orthogonal D , and by lemma 14.25, there is some $y \in D$ such that

$$\varphi(v' + y, v' + y) = \varphi(v, v).$$

Then, if we let $v_1 = v' + y$, as we showed at the beginning of the proof, we can extend f to a metric map g of $U + Kv = E$ by setting $g(v) = v_1$. Since φ is nondegenerate, g is an isometry. \square

The first corollary of Witt's theorem is sometimes called the Witt's cancellation theorem.

Theorem 14.36. (*Witt Cancellation Theorem*) *Let (E_1, φ_1) and (E_2, φ_2) be two pairs of finite-dimensional spaces and nondegenerate ϵ -Hermitian forms satisfying condition (T), and assume that (E_1, φ_1) and (E_2, φ_2) are isometric. For any subspace U of E_1 and any subspace V of E_2 , if there is an isometry $f: U \rightarrow V$, then there is an isometry $g: U^\perp \rightarrow V^\perp$.*

Proof. If $f: U \rightarrow V$ is an isometry between U and V , by Witt's theorem (Theorem 14.36), the linear map f extends to an isometry g between E_1 and E_2 . We claim that g maps U^\perp into V^\perp . This is because if $v \in U^\perp$, we have $\varphi_1(u, v) = 0$ for all $u \in U$, so

$$\varphi_2(g(u), g(v)) = \varphi_1(u, v) = 0 \quad \text{for all } u \in U,$$

and since g is a bijection between U and V , we have $g(U) = V$, so we see that $g(v)$ is orthogonal to V for every $v \in U^\perp$; that is, $g(U^\perp) \subseteq V^\perp$. Since g is a metric map and since φ_1 is nondegenerate, the restriction of g to U^\perp is an isometry from U^\perp to V^\perp . \square

A pair (E, φ) where E is finite-dimensional and φ is a nondegenerate ϵ -Hermitian form is often called an ϵ -Hermitian space. When $\epsilon = 1$ and φ is symmetric, we use the term *Euclidean space* or *quadratic space*. When $\epsilon = -1$ and φ is alternating, we use the term *symplectic space*. When $\epsilon = 1$ and the automorphism $\lambda \mapsto \bar{\lambda}$ is not the identity we use the term *Hermitian space*, and when $\epsilon = -1$, we use the term *skew-Hermitian space*.

We also have the following result showing that the group of isometries of an ϵ -Hermitian space is transitive on totally isotropic subspaces of the same dimension.

Theorem 14.37. *Let E be a finite-dimensional vector space and let φ be a nondegenerate ϵ -Hermitian form on E satisfying condition (T). Then for any two totally isotropic subspaces U and V of the same dimension, there is an isometry $f \in \mathbf{Isom}(\varphi)$ such that $f(U) = V$. Furthermore, every linear automorphism of U is induced by an isometry of E .*

Remark: Witt's cancelation theorem can be used to define an equivalence relation on ϵ -Hermitian spaces and to define a group structure on these equivalence classes. This way, we obtain the *Witt group*, but we will not discuss it here.

14.8 Symplectic Groups

In this section, we are dealing with a nondegenerate alternating form φ on a vector space E of dimension n . As we saw earlier, n must be even, say $n = 2m$. By Theorem 14.23, there is a direct sum decomposition of E into pairwise orthogonal subspaces

$$E = W_1 \overset{\perp}{\oplus} \cdots \overset{\perp}{\oplus} W_m,$$

where each W_i is a hyperbolic plane. Each W_i has a basis (u_i, v_i) , with $\varphi(u_i, u_i) = \varphi(v_i, v_i) = 0$ and $\varphi(u_i, v_i) = 1$, for $i = 1, \dots, m$. In the basis

$$(u_1, \dots, u_m, v_1, \dots, v_m),$$

φ is represented by the matrix

$$J_{m,m} = \begin{pmatrix} 0 & I_m \\ -I_m & 0 \end{pmatrix}.$$

The symplectic group $\mathbf{Sp}(2m, K)$ is the group of isometries of φ . The maps in $\mathbf{Sp}(2m, K)$ are called *symplectic* maps. With respect to the above basis, $\mathbf{Sp}(2m, K)$ is the group of $2m \times 2m$ matrices A such that

$$A^\top J_{m,m} A = J_{m,m}.$$

Matrices satisfying the above identity are called *symplectic* matrices. In this section, we show that $\mathbf{Sp}(2m, K)$ is a subgroup of $\mathbf{SL}(2m, K)$ (that is, $\det(A) = +1$ for all $A \in \mathbf{Sp}(2m, K)$), and we show that $\mathbf{Sp}(2m, K)$ is generated by special linear maps called *symplectic transvections*.

First, we leave it as an easy exercise to show that $\mathbf{Sp}(2, K) = \mathbf{SL}(2, K)$. The reader should also prove that $\mathbf{Sp}(2m, K)$ has a subgroup isomorphic to $\mathbf{GL}(m, K)$.

Next we characterize the symplectic maps f that leave fixed every vector in some given hyperplane H , that is,

$$f(v) = v \quad \text{for all } v \in H.$$

Since φ is nondegenerate, by Proposition 14.21, the orthogonal H^\perp of H is a line (that is, $\dim(H^\perp) = 1$). For every $u \in E$ and every $v \in H$, since f is an isometry and $f(v) = v$ for all $v \in H$, we have

$$\begin{aligned} \varphi(f(u) - u, v) &= \varphi(f(u), v) - \varphi(u, v) \\ &= \varphi(f(u), v) - \varphi(f(u), f(v)) \\ &= \varphi(f(u), v - f(v)) \\ &= \varphi(f(u), 0) = 0, \end{aligned}$$

which shows that $f(u) - u \in H^\perp$ for all $u \in E$. Therefore, $f - \text{id}$ is a linear map from E into the line H^\perp whose kernel contains H , which means that there is some nonzero vector $w \in H^\perp$ and some linear form ψ such that

$$f(u) = u + \psi(u)w, \quad u \in E.$$

Since f is an isometry, we must have $\varphi(f(u), f(v)) = \varphi(u, v)$ for all $u, v \in E$, which means that

$$\begin{aligned} \varphi(u, v) &= \varphi(f(u), f(v)) \\ &= \varphi(u + \psi(u)w, v + \psi(v)w) \\ &= \varphi(u, v) + \psi(u)\varphi(w, v) + \psi(v)\varphi(u, w) + \psi(u)\psi(v)\varphi(w, w) \\ &= \varphi(u, v) + \psi(u)\varphi(w, v) - \psi(v)\varphi(w, u), \end{aligned}$$

which yields

$$\psi(u)\varphi(w, v) = \psi(v)\varphi(w, u) \quad \text{for all } u, v \in E.$$

Since φ is nondegenerate, we can pick some v_0 such that $\varphi(w, v_0) \neq 0$, and we get $\psi(u)\varphi(w, v_0) = \psi(v_0)\varphi(w, u)$ for all $u \in E$; that is,

$$\psi(u) = \lambda\varphi(w, u) \quad \text{for all } u \in E,$$

for some $\lambda \in K$. Therefore, f is of the form

$$f(u) = u + \lambda\varphi(w, u)w, \quad \text{for all } u \in E.$$

It is also clear that every f of the above form is a symplectic map. If $\lambda = 0$, then $f = \text{id}$. Otherwise, if $\lambda \neq 0$, then $f(u) = u$ iff $\varphi(w, u) = 0$ iff $u \in (Kw)^\perp = H$, where H is a hyperplane. Thus, f fixes every vector in the hyperplane H . Note that since φ is alternating, $\varphi(w, w) = 0$, which means that $w \in H$.

In summary, we have characterized all the symplectic maps that leave every vector in some hyperplane fixed, and we make the following definition.

Definition 14.16. Given a nondegenerate alternating form φ on a space E , a *symplectic transvection (of direction w)* is a linear map f of the form

$$f(u) = u + \lambda\varphi(w, u)w, \quad \text{for all } u \in E,$$

for some nonzero $w \in E$ and some $\lambda \in K$. If $\lambda \neq 0$, the subspace of vectors left fixed by f is the hyperplane $H = (Kw)^\perp$. The map f is also denoted $\tau_{u, \lambda}$.

Observe that

$$\tau_{u, \lambda} \circ \tau_{u, \mu} = \tau_{u, \lambda + \mu}$$

and $\tau_{u, \lambda} = \text{id}$ iff $\lambda = 0$. The above shows that $\det(\tau_{u, \lambda}) = 1$, since when $\lambda \neq 0$, we have $\tau_{u, \lambda} = (\tau_{u, \lambda/2})^2$.

Our next goal is to show that if u and v are any two nonzero vectors in E , then there is a simple symplectic map f such that $f(u) = v$.

Proposition 14.38. *Given any two nonzero vectors $u, v \in E$, there is a symplectic map f such that $f(u) = v$, and f is either a symplectic transvection, or the composition of two symplectic transvections.*

Proof. There are two cases.

Case 1. $\varphi(u, v) \neq 0$.

In this case, $u \neq v$, since $\varphi(u, u) = 0$. Let us look for a symplectic transvection of the form $\tau_{v-u, \lambda}$. We want

$$v = u + \lambda\varphi(v - u, u)(v - u) = u + \lambda\varphi(v, u)(v - u),$$

which yields

$$(\lambda\varphi(v, u) - 1)(v - u) = 0.$$

Since $\varphi(u, v) \neq 0$ and $\varphi(v, u) = -\varphi(u, v)$, we can pick $\lambda = \varphi(v, u)^{-1}$ and $\tau_{v-u, \lambda}$ maps u to v .

Case 2. $\varphi(u, v) = 0$.

If $u = v$, use $\tau_{u, 0} = \text{id}$. Now, assume $u \neq v$. We claim that it is possible to pick some $w \in E$ such that $\varphi(u, w) \neq 0$ and $\varphi(v, w) \neq 0$. Indeed, if $(Ku)^\perp = (Kv)^\perp$, then pick any nonzero vector w not in the hyperplane $(Ku)^\perp$. Otherwise, $(Ku)^\perp$ and $(Kv)^\perp$ are two distinct hyperplanes, so neither is contained in the other (they have the same dimension), so pick any nonzero vector w_1 such that $w_1 \in (Ku)^\perp$ and $w_1 \notin (Kv)^\perp$, and pick any nonzero vector w_2 such that $w_2 \in (Kv)^\perp$ and $w_2 \notin (Ku)^\perp$. If we let $w = w_1 + w_2$, then $\varphi(u, w) = \varphi(u, w_2) \neq 0$, and $\varphi(v, w) = \varphi(v, w_1) \neq 0$. From case 1, we have some symplectic transvection τ_{w-u, λ_1} such that $\tau_{w-u, \lambda_1}(u) = w$, and some symplectic transvection τ_{v-w, λ_2} such that $\tau_{v-w, \lambda_2}(w) = v$, so the composition $\tau_{v-w, \lambda_2} \circ \tau_{w-u, \lambda_1}$ maps u to v . \square

Next, we would like to extend Proposition 14.38 to two hyperbolic planes W_1 and W_2 .

Proposition 14.39. *Given any two hyperbolic planes W_1 and W_2 given by bases (u_1, v_1) and (u_2, v_2) (with $\varphi(u_i, u_i) = \varphi(v_i, v_i) = 0$ and $\varphi(u_i, v_i) = 1$, for $i = 1, 2$), there is a symplectic map f such that $f(u_1) = u_2$, $f(v_1) = v_2$, and f is the composition of at most four symplectic transvections.*

Proof. From Proposition 14.38, we can map u_1 to u_2 , using a map f which is the composition of at most two symplectic transvections. Say $v_3 = f(v_1)$. We claim that there is a map g such that $g(u_2) = u_2$ and $g(v_3) = v_2$, and g is the composition of at most two symplectic transvections. If so, $g \circ f$ maps the pair (u_1, v_1) to the pair (u_2, v_2) , and $g \circ f$ consists of at most four symplectic transvections. Thus, we need to prove the following claim:

Claim. If (u, v) and (u, v') are hyperbolic bases determining two hyperbolic planes, then there is a symplectic map g such that $g(u) = u$, $g(v) = v'$, and g is the composition of at most two symplectic transvections. There are two case.

Case 1. $\varphi(v, v') \neq 0$.

In this case, there is a symplectic transvection $\tau_{v'-v, \lambda}$ such that $\tau_{v'-v, \lambda}(v) = v'$. We also have

$$\varphi(u, v' - v) = \varphi(u, v') - \varphi(u, v) = 1 - 1 = 0.$$

Therefore, $\tau_{v'-v, \lambda}(u) = u$, and $g = \tau_{v'-v, \lambda}$ does the job.

Case 2. $\varphi(v, v') = 0$.

First, check that $(u, u + v)$ is also a hyperbolic basis. Furthermore,

$$\varphi(v, u + v) = \varphi(v, u) + \varphi(v, v) = \varphi(v, u) = -1 \neq 0.$$

Thus, there is a symplectic transvection τ_{u, λ_1} such that $\tau_{u, \lambda_1}(v) = u + v$ and $\tau_{u, \lambda_1}(u) = u$. We also have

$$\varphi(u + v, v') = \varphi(u, v') + \varphi(v, v') = \varphi(u, v') = 1 \neq 0,$$

so there is a symplectic transvection $\tau_{v'-u-v, \lambda_2}$ such that $\tau_{v'-u-v, \lambda_2}(u + v) = v'$. Since

$$\varphi(u, v' - u - v) = \varphi(u, v') - \varphi(u, u) - \varphi(u, v) = 1 - 0 - 1 = 0,$$

we have $\tau_{v'-u-v, \lambda_2}(u) = u$. Then, the composition $g = \tau_{v'-u-v, \lambda_2} \circ \tau_{u, \lambda_1}$ is such that $g(u) = u$ and $g(v) = v'$. \square

We will use Proposition 14.39 in an inductive argument to prove that the symplectic transvections generate the symplectic group. First, make the following observation: If U is a nondegenerate subspace of E , so that

$$E = U \oplus U^\perp,$$

and if τ is a transvection of H^\perp , then we can form the linear map $\text{id}_U \oplus \tau$ whose restriction to U is the identity and whose restriction to U^\perp is τ , and $\text{id}_U \oplus \tau$ is a transvection of E .

Theorem 14.40. *The symplectic group $\mathbf{Sp}(2m, K)$ is generated by the symplectic transvections. For every transvection $f \in \mathbf{Sp}(2m, K)$, we have $\det(f) = 1$.*

Proof. Let G be the subgroup of $\mathbf{Sp}(2m, K)$ generated by the transvections. We need to prove that $G = \mathbf{Sp}(2m, K)$. Let $(u_1, v_1, \dots, u_m, v_m)$ be a symplectic basis of E , and let $f \in \mathbf{Sp}(2m, K)$ be any symplectic map. Then, f maps $(u_1, v_1, \dots, u_m, v_m)$ to another symplectic basis $(u'_1, v'_1, \dots, u'_m, v'_m)$. If we prove that there is some $g \in G$ such that $g(u_i) = u'_i$ and $g(v_i) = v'_i$ for $i = 1, \dots, m$, then $f = g$ and $G = \mathbf{Sp}(2m, K)$.

We use induction on i to prove that there is some $g_i \in G$ so that g_i maps $(u_1, v_1, \dots, u_i, v_i)$ to $(u'_1, v'_1, \dots, u'_i, v'_i)$.

The base case $i = 1$ follows from Proposition 14.39.

For the induction step, assume that we have some $g_i \in G$ mapping $(u_1, v_1, \dots, u_i, v_i)$ to $(u'_1, v'_1, \dots, u'_i, v'_i)$, and let $(u''_{i+1}, v''_{i+1}, \dots, u''_m, v''_m)$ be the image of $(u_{i+1}, v_{i+1}, \dots, u_m, v_m)$

by g_i . If U is the subspace spanned by $(u'_1, v'_1, \dots, u'_m, v'_m)$, then each hyperbolic plane W'_{i+k} given by (u'_{i+k}, v'_{i+k}) and each hyperbolic plane W''_{i+k} given by (u''_{i+k}, v''_{i+k}) belongs to U^\perp . Using the remark before the theorem and Proposition 14.39, we can find a transvection τ mapping W''_{i+1} onto W'_{i+1} and leaving every vector in U fixed. Then, $\tau \circ g_i$ maps $(u_1, v_1, \dots, u_{i+1}, v_{i+1})$ to $(u'_1, v'_1, \dots, u'_{i+1}, v'_{i+1})$, establishing the induction step.

For the second statement, since we already proved that every transvection has a determinant equal to $+1$, this also holds for any composition of transvections in G , and since $G = \mathbf{Sp}(2m, K)$, we are done. \square

It can also be shown that the center of $\mathbf{Sp}(2m, K)$ is reduced to the subgroup $\{\text{id}, -\text{id}\}$. The *projective symplectic group* $\mathbf{PSp}(2m, K)$ is the quotient group $\mathbf{Sp}(2m, K)/\{\text{id}, -\text{id}\}$. All symplectic projective groups are simple, except $\mathbf{PSp}(2, \mathbb{F}_2)$, $\mathbf{PSp}(2, \mathbb{F}_3)$, and $\mathbf{PSp}(4, \mathbb{F}_2)$, see Grove [52].

The orders of the symplectic groups over finite fields can be determined. For details, see Artin [3], Jacobson [59] and Grove [52].

An interesting property of symplectic spaces is that the determinant of a skew-symmetric matrix B is the square of some polynomial $\text{Pf}(B)$ called the *Pfaffian*; see Jacobson [59] and Artin [3]. We leave considerations of the Pfaffian to the exercises.

We now take a look at the orthogonal groups.

14.9 Orthogonal Groups

In this section, we are dealing with a nondegenerate symmetric bilinear form φ over a finite-dimensional vector space E of dimension n over a field of characteristic not equal to 2. Recall that the orthogonal group $\mathbf{O}(\varphi)$ is the group of isometries of φ ; that is, the group of linear maps $f: E \rightarrow E$ such that

$$\varphi(f(u), f(v)) = \varphi(u, v) \quad \text{for all } u, v \in E.$$

The elements of $\mathbf{O}(\varphi)$ are also called *orthogonal transformations*. If M is the matrix of φ in any basis, then a matrix A represents an orthogonal transformation iff

$$A^\top M A = M.$$

Since φ is nondegenerate, M is invertible, so we see that $\det(A) = \pm 1$. The subgroup

$$\mathbf{SO}(\varphi) = \{f \in \mathbf{O}(\varphi) \mid \det(f) = 1\}$$

is called the *special orthogonal group* (of φ), and its members are called *rotations* (or *proper orthogonal transformations*). Isometries $f \in \mathbf{O}(\varphi)$ such that $\det(f) = -1$ are called *improper orthogonal transformations*, or sometimes *reversions*.

If H is any nondegenerate hyperplane in E , then $D = H^\perp$ is a nondegenerate line and we have

$$E = H \oplus H^\perp.$$

For any nonzero vector $u \in D = H^\perp$ Consider the map τ_u given by

$$\tau_u(v) = v - 2 \frac{\varphi(v, u)}{\varphi(u, u)} u \quad \text{for all } v \in E.$$

If we replace u by λu with $\lambda \neq 0$, we have

$$\tau_{\lambda u}(v) = v - 2 \frac{\varphi(v, \lambda u)}{\varphi(\lambda u, \lambda u)} \lambda u = v - 2 \frac{\lambda \varphi(v, u)}{\lambda^2 \varphi(u, u)} \lambda u = v - 2 \frac{\varphi(v, u)}{\varphi(u, u)} u,$$

which shows that τ_u depends only on the line D , and thus only the hyperplane H . Therefore, denote by τ_H the linear map τ_u determined as above by any nonzero vector $u \in H^\perp$. Note that if $v \in H$, then

$$\tau_H(v) = v,$$

and if $v \in D$, then

$$\tau_H(v) = -v.$$

A simple computation shows that

$$\varphi(\tau_H(u), \tau_H(v)) = \varphi(u, v) \quad \text{for all } u, v \in E,$$

so $\tau_H \in \mathbf{O}(\varphi)$, and by picking a basis consisting of u and vectors in H , that $\det(\tau_H) = -1$. It is also clear that $\tau_H^2 = \text{id}$.

Definition 14.17. If H is any nondegenerate hyperplane in E , for any nonzero vector $u \in H^\perp$, the linear map τ_H given by

$$\tau_H(v) = v - 2 \frac{\varphi(v, u)}{\varphi(u, u)} u \quad \text{for all } v \in E$$

is an involutive isometry of E called the *reflection through (or about) the hyperplane H* .

Remarks:

1. It can be shown that if $f \in \mathbf{O}(\varphi)$ leaves every vector in some hyperplane H fixed, then either $f = \text{id}$ or $f = \tau_H$; see Taylor [108] (Chapter 11). Thus, there is no analog to symplectic transvections in the orthogonal group.
2. If $K = \mathbb{R}$ and φ is the usual Euclidean inner product, the matrices corresponding to hyperplane reflections are called *Householder matrices*.

Our goal is to prove that $\mathbf{O}(\varphi)$ is generated by the hyperplane reflections. The following proposition is needed.

Proposition 14.41. *Let φ be a nondegenerate symmetric bilinear form on a vector space E . For any two nonzero vectors $u, v \in E$, if $\varphi(u, u) = \varphi(v, v)$ and $v - u$ is nonisotropic, then the hyperplane reflection $\tau_H = \tau_{v-u}$ maps u to v , with $H = (K(v - u))^\perp$.*

Proof. Since $v - u$ is not isotropic, $\varphi(v - u, v - u) \neq 0$, and we have

$$\begin{aligned} \tau_{v-u}(u) &= u - 2 \frac{\varphi(u, v - u)}{\varphi(v - u, v - u)}(v - u) \\ &= u - 2 \frac{\varphi(u, v) - \varphi(u, u)}{\varphi(v, v) - 2\varphi(u, v) + \varphi(u, u)}(v - u) \\ &= u - \frac{2(\varphi(u, v) - \varphi(u, u))}{2(\varphi(u, u) - \varphi(u, v))}(v - u) \\ &= v, \end{aligned}$$

which proves the proposition. \square

We can now obtain a cheap version of the Cartan–Dieudonné theorem.

Theorem 14.42. *(Cartan–Dieudonné, weak form) Let φ be a nondegenerate symmetric bilinear form on a K -vector space E of dimension n ($\text{char}(K) \neq 2$). Then, every isometry $f \in \mathbf{O}(\varphi)$ with $f \neq \text{id}$ is the composition of at most $2n - 1$ hyperplane reflections.*

Proof. We proceed by induction on n . For $n = 0$, this is trivial (since $\mathbf{O}(\varphi) = \{\text{id}\}$).

Next, assume that $n \geq 1$. Since φ is nondegenerate, we know that there is some nonisotropic vector $u \in E$. There are three cases.

Case 1. $f(u) = u$.

Since φ is nondegenerate and u is nonisotropic, the hyperplane $H = (Ku)^\perp$ is nondegenerate, $E = H \oplus (Ku)^\perp$, and since $f(u) = u$, we must have $f(H) = H$. The restriction f' of f to H is an isometry of H . By the induction hypothesis, we can write

$$f' = \tau'_k \circ \cdots \circ \tau'_1,$$

where τ_i is some hyperplane reflection about a hyperplane L_i in H , with $k \leq 2n - 3$. We can extend each τ'_i to a reflection τ_i about the hyperplane $L_i \oplus Ku$ so that $\tau_i(u) = u$, and clearly,

$$f = \tau_k \circ \cdots \circ \tau_1.$$

Case 2. $f(u) = -u$.

If τ is the hyperplane reflection about the hyperplane $H = (Ku)^\perp$, then $g = \tau \circ f$ is an isometry of E such that $g(u) = u$, and we are back to Case (1). Since $\tau^2 = 1$ We obtain

$$f = \tau \circ \tau_k \circ \cdots \circ \tau_1$$

where τ and the τ_i are hyperplane reflections, with $k \geq 2n - 3$, and we get a total of $2n - 2$ hyperplane reflections.

Case 3. $f(u) \neq u$ and $f(u) \neq -u$.

Note that $f(u) - u$ and $f(u) + u$ are orthogonal, since

$$\begin{aligned}\varphi(f(u) - u, f(u) + u) &= \varphi(f(u), f(u)) + \varphi(f(u), u) - \varphi(u, f(u)) - \varphi(u, u) \\ &= \varphi(u, u) - \varphi(u, u) = 0.\end{aligned}$$

We also have

$$\begin{aligned}\varphi(u, u) &= \varphi((f(u) + u - (f(u) - u))/2, (f(u) + u - (f(u) - u))/2) \\ &= \frac{1}{4}\varphi(f(u) + u, f(u) + u) + \frac{1}{4}\varphi(f(u) - u, f(u) - u),\end{aligned}$$

so $f(u) + u$ and $f(u) - u$ cannot be both isotropic, since u is not isotropic.

If $f(u) - u$ is not isotropic, then the reflection $\tau_{f(u)-u}$ is such that

$$\tau_{f(u)-u}(u) = f(u),$$

and since $\tau_{f(u)-u}^2 = \text{id}$, if $g = \tau_{f(u)-u} \circ f$, then $g(u) = u$, and we are back to case (1). We obtain

$$f = \tau_{f(u)-u} \circ \tau_k \circ \cdots \circ \tau_1$$

where $\tau_{f(u)-u}$ and the τ_i are hyperplane reflections, with $k \geq 2n - 3$, and we get a total of $2n - 2$ hyperplane reflections.

If $f(u) + u$ is not isotropic, then the reflection $\tau_{f(u)+u}$ is such that

$$\tau_{f(u)+u}(u) = -f(u),$$

and since $\tau_{f(u)+u}^2 = \text{id}$, if $g = \tau_{f(u)+u} \circ f$, then $g(u) = -u$, and we are back to case (2). We obtain

$$f = \tau_{f(u)+u} \circ \tau \circ \tau_k \circ \cdots \circ \tau_1$$

where $\tau, \tau_{f(u)+u}$ and the τ_i are hyperplane reflections, with $k \geq 2n - 3$, and we get a total of $2n - 1$ hyperplane reflections. This proves the induction step. \square

The bound $2n - 1$ is not optimal. The strong version of the Cartan–Dieudonné theorem says that at most n reflections are needed, but the proof is harder. Here is a neat proof due to E. Artin (see [3], Chapter III, Section 4).

Case 1 remains unchanged. Case 2 is slightly different: $f(u) - u$ is not isotropic. Since $\varphi(f(u) + u, f(u) - u) = 0$, as in the first subcase of Case (3), $g = \tau_{f(u)-u} \circ f$ is such that $g(u) = u$ and we are back to Case 1. This only costs one more reflection.

The new (bad) case is:

Case 3'. $f(u) - u$ is nonzero and isotropic for all nonisotropic $u \in E$. In this case, what saves us is that E must be an Artinian space of dimension $n = 2m$ and that f must be a rotation ($f \in \mathbf{SO}(\varphi)$).

If we accept this fact, then pick any hyperplane reflection τ . Then, since f is a rotation, $g = \tau \circ f$ is *not* a rotation because $\det(g) = \det(\tau)\det(f) = (-1)(+1) = -1$, so $g(u) - u$ is not isotropic for all nonisotropic $u \in E$, we are back to Case 2, and using the induction hypothesis, we get

$$\tau \circ f = \tau_k \circ \dots \circ \tau_1,$$

where each τ_i is a hyperplane reflection, and $k \leq 2m$. Since $\tau \circ f$ is not a rotation, actually $k \leq 2m - 1$, and then $f = \tau \circ \tau_k \circ \dots \circ \tau_1$, the composition of at most $k + 1 \leq 2m$ hyperplane reflections.

Therefore, except for the fact that in Case 3', E must be an Artinian space of dimension $n = 2m$ and that f must be a rotation, which has not been proven yet, we proved the following theorem.

Theorem 14.43. (*Cartan–Dieudonné, strong form*) *Let φ be a nondegenerate symmetric bilinear form on a K -vector space E of dimension n ($\text{char}(K) \neq 2$). Then, every isometry $f \in \mathbf{O}(\varphi)$ with $f \neq \text{id}$ is the composition of at most n hyperplane reflections.*

To fill in the gap, we need two propositions.

Proposition 14.44. *Let (E, φ) be an Artinian space of dimension $2m$, and let U be a totally isotropic subspace of dimension m . For any isometry $f \in \mathbf{O}(\varphi)$, we have $\det(f) = 1$ (f is a rotation).*

Proof. We know that we can find a basis $(u_1, \dots, u_m, v_1, \dots, v_m)$ of E such (u_1, \dots, u_m) is a basis of U and φ is represented by the matrix

$$\begin{pmatrix} 0 & I_m \\ I_m & 0 \end{pmatrix}.$$

Since $f(U) = U$, the matrix representing f is of the form

$$A = \begin{pmatrix} B & C \\ 0 & D \end{pmatrix}.$$

The condition $A^\top A_{m,m} A = A_{m,m}$ translates as

$$\begin{pmatrix} B^\top & 0 \\ C^\top & D^\top \end{pmatrix} \begin{pmatrix} 0 & I_m \\ I_m & 0 \end{pmatrix} \begin{pmatrix} B & C \\ 0 & D \end{pmatrix} = \begin{pmatrix} 0 & I_m \\ I_m & 0 \end{pmatrix}$$

that is,

$$\begin{pmatrix} B^\top & 0 \\ C^\top & D^\top \end{pmatrix} \begin{pmatrix} 0 & D \\ B & C \end{pmatrix} = \begin{pmatrix} 0 & B^\top D \\ D^\top B & C^\top D + D^\top C \end{pmatrix} = \begin{pmatrix} 0 & I_m \\ I_m & 0 \end{pmatrix},$$

which implies that $B^\top D = I$, and so

$$\det(A) = \det(B) \det(D) = \det(B^\top) \det(D) = \det(B^\top D) = \det(I) = 1,$$

as claimed □

Proposition 14.45. *Let φ be a nondegenerate symmetric bilinear form on a space E of dimension n , and let f be any isometry $f \in \mathbf{O}(\varphi)$ such that $f(u) - u$ is nonzero and isotropic for every nonisotropic vector $u \in E$. Then, E is an Artinian space of dimension $n = 2m$, and f is a rotation ($f \in \mathbf{SO}(\varphi)$).*

Proof. We follow E. Artin's proof (see [3], Chapter III, Section 4). First, consider the case $n = 2$. Since we are assuming that E has some nonzero isotropic vector, by Proposition 14.26, E is an Artinian plane and there is a basis in which φ is represented by the matrix

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix},$$

we have $\varphi((x_1, x_2), (x_1, x_2)) = 2x_1x_2$, and the matrices representing isometries are of the form

$$\begin{pmatrix} \lambda & 0 \\ 0 & \lambda^{-1} \end{pmatrix} \quad \text{or} \quad \begin{pmatrix} 0 & \lambda \\ \lambda^{-1} & 0 \end{pmatrix}, \quad \lambda \in K - \{0\}.$$

In the second case,

$$\begin{pmatrix} 0 & \lambda \\ \lambda^{-1} & 0 \end{pmatrix} \begin{pmatrix} \lambda \\ 1 \end{pmatrix} = \begin{pmatrix} \lambda \\ 1 \end{pmatrix},$$

but $u = (\lambda, 1)$ is a nonisotropic vector such that $f(u) - u = 0$. Therefore, we must be in the first case, and $\det(f) = +1$.

Let us now assume that $n \geq 3$. Let y be some nonzero isotropic vector. Since $n \geq 3$, the orthogonal space $(Ky)^\perp$ has dimension at least 2, and we know that $\text{rad}(Ky) = \text{rad}((Ky)^\perp)$, which implies that $(Ky)^\perp$ contains some nonisotropic vector, say x . We have $\varphi(x, y) = 0$, so $\varphi(x + \epsilon y, x + \epsilon y) = \varphi(x, x) \neq 0$, for $\epsilon = \pm 1$. Then, by hypothesis, the vectors $f(x) - x$, $f(x + y) - (x + y) = f(x) - x + (f(y) - y)$, and $f(x - y) - (x - y) = f(x) - x - (f(y) - y)$ are isotropic. The last two vectors can be written as $f(x) - x + \epsilon(f(y) - y)$ with $\epsilon = \pm 1$, so we have

$$\begin{aligned} 0 &= \varphi(f(x) - x + \epsilon(f(y) - y), f(x) - x + \epsilon(f(y) - y)) \\ &= 2\epsilon\varphi(f(x) - x, f(y) - y) + \epsilon^2\varphi(f(y) - y, f(y) - y). \end{aligned}$$

If we write the two equations corresponding to $\epsilon = \pm 1$, and then add them up, we get

$$\varphi(f(y) - y, f(y) - y) = 0.$$

Therefore, we proved that $f(u) - u$ is isotropic for every $u \in E$. If we let $W = \text{Im}(f - \text{id})$, then every vector in W is isotropic, and thus W is totally isotropic (recall that we assumed

that $\text{char}(K) \neq 2$, so φ is determined by Φ). For any $u \in E$ and any $v \in W^\perp$, since W is totally isotropic, we have

$$\varphi(f(u) - u, f(v) - v) = 0,$$

and since $f(u) - u \in W$ and $v \in W^\perp$, we have $\varphi(f(u) - u, v) = 0$, and so

$$\begin{aligned} 0 &= \varphi(f(u) - u, f(v) - v) \\ &= \varphi(f(u), f(v)) - \varphi(u, f(v)) - \varphi(f(u) - u, v) \\ &= \varphi(u, v) - \varphi(u, f(v)) \\ &= \varphi(u, v - f(v)), \end{aligned}$$

for all $u \in E$. Since φ is nonsingular, this means that $f(v) = v$, for all $v \in W^\perp$. However, by hypothesis, no nonisotropic vector is left fixed, which implies that W^\perp is also totally isotropic. In summary, we proved that $W \subseteq W^\perp$ and $W^\perp \subseteq W^{\perp\perp} = W$, that is,

$$W = W^\perp.$$

Since, $\dim(W) + \dim(W^\perp) = n$, we conclude that W is a totally isotropic subspace of E such that

$$\dim(W) = n/2.$$

By Proposition 14.27, the space E is an Artinian space of dimension $n = 2m$. Since $W = W^\perp$ and $f(W^\perp) = W^\perp$, by Proposition 14.44, the isometry f is a rotation. \square

Remarks:

1. Another way to finish the proof of Proposition 14.45 is to prove that if f is an isometry, then

$$\text{Ker}(f - \text{id}) = (\text{Im}(f - \text{id}))^\perp.$$

After having proved that $W = \text{Im}(f - \text{id})$ is totally isotropic, we get

$$\text{Ker}(f - \text{id}) = \text{Im}(f - \text{id}),$$

which implies that $(f - \text{id})^2 = 0$. From this, we deduce that $\det(f) = 1$. For details, see Jacobson [59] (Chapter 6, Section 6).

2. If $f = \tau_{H_k} \circ \cdots \circ \tau_{H_1}$, where the H_i are hyperplanes, then it can be shown that

$$\dim(H_1 \cap H_2 \cap \cdots \cap H_s) \geq n - s.$$

Now, since each H_i is left fixed by τ_{H_i} , we see that every vector in $H_1 \cap \cdots \cap H_s$ is left fixed by f . In particular, if $s < n$, then f has some nonzero fixed point. As a consequence, an isometry without fixed points requires n hyperplane reflections.

Witt's Theorem can be sharpened to isometries in $\mathbf{SO}(\varphi)$, but some condition on U is needed.

Theorem 14.46. (*Witt–Sharpened Version*) *Let E be a finite-dimensional space equipped with a nondegenerate symmetric bilinear forms φ . For any subspace U of E , every linear injective metric map f from U into E extends to an isometry g of E with a prescribed value ± 1 of $\det(g)$ iff*

$$\dim(U) + \dim(\text{rad}(U)) < \dim(E) = n.$$

If

$$\dim(U) + \dim(\text{rad}(U)) = \dim(E) = n,$$

and $\det(f) = -1$, then there is no $g \in \mathbf{SO}(\varphi)$ extending f .

Proof. If g_1 and g_2 are two extensions of f such that $\det(g_1)\det(g_2) = -1$, then $h = g_1^{-1} \circ g_2$ is an isometry such that $\det(h) = -1$, and h leaves every vector of U fixed. Conversely, if h is an isometry such that $\det(h) = -1$, and $h(u) = u$ for all $u \in U$, then for any extension g_1 of f , the map $g_2 = h \circ g_1$ is another extension of f such that $\det(g_2) = -\det(g_1)$. Therefore, we need to show that a map h as above exists.

If $\dim(U) + \dim(\text{rad}(U)) < \dim(E)$, consider the nondegenerate completion \bar{U} of U given by Proposition 14.30. We know that $\dim(\bar{U}) = \dim(U) + \dim(\text{rad}(U)) < n$, and since \bar{U} is nondegenerate, we have

$$E = \bar{U} \oplus \bar{U}^\perp,$$

with $\bar{U}^\perp \neq (0)$. Pick any isometry τ of \bar{U}^\perp such that $\det(\tau) = -1$, and extend it to an isometry h of E whose restriction to \bar{U} is the identity.

If $\dim(U) + \dim(\text{rad}(U)) = \dim(E) = n$, then $U = V \oplus W$ with $V = \text{rad}(U)$ and since $\dim(\bar{U}) = \dim(U) + \dim(\text{rad}(U)) = n$, we have

$$E = \bar{U} = (V \oplus V') \oplus W,$$

where $V \oplus V' = \text{Ar}_{2r} = W^\perp$ is an Artinian space. Any isometry h of E which is the identity on U and with $\det(h) = -1$ is the identity on W , and thus it must map $W^\perp = \text{Ar}_{2r} = V \oplus V'$ into itself, and the restriction h' of h to Ar_{2r} has $\det(h') = -1$. However, h' is the identity on $V = \text{rad}(U)$, a totally isotopic subspace of Ar_{2r} of dimension r , and by Proposition 14.44, we have $\det(h') = +1$, a contradiction. \square

It can be shown that the center of $\mathbf{O}(\varphi)$ is $\{\text{id}, -\text{id}\}$. For further properties of orthogonal groups, see Grove [52], Jacobson [59], Taylor [108], and Artin [3].

Chapter 15

Variational Approximation of Boundary-Value Problems; Introduction to the Finite Elements Method

15.1 A One-Dimensional Problem: Bending of a Beam

Consider a beam of unit length supported at its ends in 0 and 1, stretched along its axis by a force P , and subjected to a transverse load $f(x)dx$ per element dx , as illustrated in Figure 15.1.

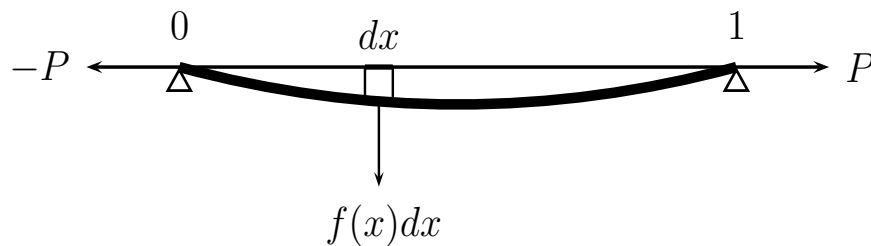


Figure 15.1: Vertical deflection of a beam

The bending moment $u(x)$ at the abscissa x is the solution of a boundary problem (BP) of the form

$$\begin{aligned} -u''(x) + c(x)u(x) &= f(x), & 0 < x < 1 \\ u(0) &= \alpha \\ u(1) &= \beta, \end{aligned}$$

where $c(x) = P/(EI(x))$, where E is the Young's modulus of the material of which the beam is made and $I(x)$ is the principal moment of inertia of the cross-section of the beam at the abscissa x , and with $\alpha = \beta = 0$. For this problem, we may assume that $c(x) \geq 0$ for all $x \in [0, 1]$.

Remark: The vertical deflection $w(x)$ of the beam and the bending moment $u(x)$ are related by the equation

$$u(x) = -EI \frac{d^2 w}{dx^2}.$$

If we seek a solution $u \in C^2([0, 1])$, that is, a function whose first and second derivatives exist and are continuous, then it can be shown that the problem has a unique solution (assuming c and f to be continuous functions on $[0, 1]$).

Except in very rare situations, this problem has no closed-form solution, so we are led to seek approximations of the solutions.

One way to proceed is to use the *finite difference method*, where we discretize the problem and replace derivatives by differences. Another way is to use a variational approach. In this approach, we follow a somewhat surprising path in which we come up with a so-called “weak formulation” of the problem, by using a trick based on integrating by parts!

First, let us observe that we can always assume that $\alpha = \beta = 0$, by looking for a solution of the form $u(x) - (\alpha(1-x) + \beta x)$. This turns out to be crucial when we integrate by parts. There are a lot of subtle mathematical details involved to make what follows rigorous, but here, we will take a “relaxed” approach.

First, we need to specify the space of “weak solutions.” This will be the vector space V of continuous functions f on $[0, 1]$, with $f(0) = f(1) = 0$, and which are piecewise continuously differentiable on $[0, 1]$. This means that there is a finite number of points x_0, \dots, x_{N+1} with $x_0 = 0$ and $x_{N+1} = 1$, such that $f'(x_i)$ is undefined for $i = 1, \dots, N$, but otherwise f' is defined and continuous on each interval (x_i, x_{i+1}) for $i = 0, \dots, N$.¹ The space V becomes a Euclidean vector space under the inner product

$$\langle f, g \rangle_V = \int_0^1 (f(x)g(x) + f'(x)g'(x))dx,$$

for all $f, g \in V$. The associated norm is

$$\|f\|_V = \left(\int_0^1 (f(x)^2 + f'(x)^2)dx \right)^{1/2}.$$

Assume that u is a solution of our original boundary problem (BP), so that

$$\begin{aligned} -u''(x) + c(x)u(x) &= f(x), & 0 < x < 1 \\ u(0) &= 0 \\ u(1) &= 0. \end{aligned}$$

¹We also assume that $f'(x)$ has a limit when x tends to a boundary of (x_i, x_{i+1}) .

Multiply the differential equation by any arbitrary *test function* $v \in V$, obtaining

$$-u''(x)v(x) + c(x)u(x)v(x) = f(x)v(x), \quad (*)$$

and integrate this equation! We get

$$-\int_0^1 u''(x)v(x)dx + \int_0^1 c(x)u(x)v(x)dx = \int_0^1 f(x)v(x)dx. \quad (\dagger)$$

Now, the trick is to use integration by parts on the first term. Recall that

$$(u'v)' = u''v + u'v',$$

and to be careful about discontinuities, write

$$\int_0^1 u''(x)v(x)dx = \sum_{i=0}^N \int_{x_i}^{x_{i+1}} u''(x)v(x)dx.$$

Using integration by parts, we have

$$\begin{aligned} \int_{x_i}^{x_{i+1}} u''(x)v(x)dx &= \int_{x_i}^{x_{i+1}} (u'(x)v(x))'dx - \int_{x_i}^{x_{i+1}} u'(x)v'(x)dx \\ &= [u'(x)v(x)]_{x=x_i}^{x=x_{i+1}} - \int_{x_i}^{x_{i+1}} u'(x)v'(x)dx \\ &= u'(x_{i+1})v(x_{i+1}) - u'(x_i)v(x_i) - \int_{x_i}^{x_{i+1}} u'(x)v'(x)dx. \end{aligned}$$

It follows that

$$\begin{aligned} \int_0^1 u''(x)v(x)dx &= \sum_{i=0}^N \int_{x_i}^{x_{i+1}} u''(x)v(x)dx \\ &= \sum_{i=0}^N \left(u'(x_{i+1})v(x_{i+1}) - u'(x_i)v(x_i) - \int_{x_i}^{x_{i+1}} u'(x)v'(x)dx \right) \\ &= u'(1)v(1) - u'(0)v(0) - \int_0^1 u'(x)v'(x)dx. \end{aligned}$$

However, the test function v satisfies the boundary conditions $v(0) = v(1) = 0$ (recall that $v \in V$), so we get

$$\int_0^1 u''(x)v(x)dx = - \int_0^1 u'(x)v'(x)dx.$$

Consequently, the equation (\dagger) becomes

$$\int_0^1 u'(x)v'(x)dx + \int_0^1 c(x)u(x)v(x)dx = \int_0^1 f(x)v(x)dx,$$

or

$$\int_0^1 (u'v' + cuv)dx = \int_0^1 fvdx, \quad \text{for all } v \in V. \quad (**)$$

Thus, it is natural to introduce the bilinear form $a: V \times V \rightarrow \mathbb{R}$ given by

$$a(u, v) = \int_0^1 (u'v' + cuv)dx, \quad \text{for all } u, v \in V,$$

and the linear form $\tilde{f}: V \rightarrow \mathbb{R}$ given by

$$\tilde{f}(v) = \int_0^1 f(x)v(x)dx, \quad \text{for all } v \in V.$$

Then, (**) becomes

$$a(u, v) = \tilde{f}(v), \quad \text{for all } v \in V.$$

We also introduce the *energy function* J given by

$$J(v) = \frac{1}{2}a(v, v) - \tilde{f}(v) \quad v \in V.$$

Then, we have the following theorem.

Theorem 15.1. *Let u be any solution of the boundary problem (BP).*

(1) *Then we have*

$$a(u, v) = \tilde{f}(v), \quad \text{for all } v \in V, \quad (\text{WF})$$

where

$$a(u, v) = \int_0^1 (u'v' + cuv)dx, \quad \text{for all } u, v \in V,$$

and

$$\tilde{f}(v) = \int_0^1 f(x)v(x)dx, \quad \text{for all } v \in V.$$

(2) *If $c(x) \geq 0$ for all $x \in [0, 1]$, then a function $u \in V$ is a solution of (WF) iff u minimizes $J(v)$, that is,*

$$J(u) = \inf_{v \in V} J(v),$$

with

$$J(v) = \frac{1}{2}a(v, v) - \tilde{f}(v) \quad v \in V.$$

Furthermore, u is unique.

Proof. We already proved (1).

To prove (2), first we show that

$$\|v\|_V^2 \leq 2a(v, v), \quad \text{for all } v \in V.$$

For this, it suffices to prove that

$$\|v\|_V^2 \leq 2 \int_0^1 (f'(x))^2 dx, \quad \text{for all } v \in V.$$

However, by Cauchy-Schwarz for functions, for every $x \in [0, 1]$, we have

$$|v(x)| = \left| \int_0^x v'(t) dt \right| \leq \int_0^1 |v'(t)| dt \leq \left(\int_0^1 |v'(t)|^2 dt \right)^{1/2},$$

and so

$$\|v\|_V^2 = \int_0^1 ((v(x))^2 + (v'(x))^2) dx \leq 2 \int_0^1 (v'(x))^2 dx \leq 2a(v, v),$$

since

$$a(v, v) = \int_0^1 ((v')^2 + cv^2) dx.$$

Next, it is easy to check that

$$J(u + v) - J(u) = a(u, v) - \tilde{f}(v) + \frac{1}{2}a(v, v), \quad \text{for all } u, v \in V.$$

Then, if u is a solution of (WF), we deduce that

$$J(u + v) - J(u) = \frac{1}{2}a(v, v) \geq \frac{1}{4}\|v\|_V^2 \geq 0 \quad \text{for all } v \in V.$$

since $a(u, v) - \tilde{f}(v) = 0$ for all $v \in V$. Therefore, J achieves a minimum for u .

We also have

$$J(u + \theta v) - J(u) = \theta(a(u, v) - \tilde{f}(v)) + \frac{\theta^2}{2}a(v, v) \quad \text{for all } \theta \in \mathbb{R},$$

and so $J(u + \theta v) - J(u) \geq 0$ for all $\theta \in \mathbb{R}$. Consequently, if J achieves a minimum for u , then $a(u, v) = \tilde{f}(v)$, which means that u is a solution of (WF).

Finally, assuming that $c(x) \geq 0$, we claim that if $v \in V$ and $v \neq 0$, then $a(v, v) > 0$. This is because if $a(v, v) = 0$, since

$$\|v\|_V^2 \leq 2a(v, v) \quad \text{for all } v \in V,$$

we would have $\|v\|_V = 0$, that is, $v = 0$. Then, if $v \neq 0$, from

$$J(u + v) - J(u) = \frac{1}{2}a(v, v) \quad \text{for all } v \in V$$

we see that $J(u + v) > J(u)$, so the minimum u is unique □

Theorem 15.1 shows that every solution u of our boundary problem (BP) is a solution (in fact, unique) of the equation (WF).

The equation (WF) is called the *weak form* or *variational equation* associated with the boundary problem. This idea to derive these equations is due to *Ritz and Galerkin*.

Now, the natural question is whether the variational equation (WF) has a solution, and whether this solution, if it exists, is also a solution of the boundary problem (it must belong to $C^2([0, 1])$, which is far from obvious). Then, (BP) and (WF) would be equivalent.

Some fancy tools of analysis can be used to prove these assertions. The first difficulty is that the vector space V is not the right space of solutions, because in order for the variational problem to have a solution, it must be complete. So, we must construct a completion of the vector space V . This can be done and we get the *Sobolev space* $H_0^1(0, 1)$. Then, the question of the regularity of the “weak solution” can also be tackled.

We will not worry about all this. Instead, let us find *approximations* of the problem (WF). Instead of using the infinite-dimensional vector space V , we consider *finite-dimensional* subspaces V_a (with $\dim(V_a) = n$) of V , and we consider the *discrete problem*:

Find a function $u^{(a)} \in V_a$, such that

$$a(u^{(a)}, v) = \tilde{f}(v), \quad \text{for all } v \in V_a. \quad (\text{DWF})$$

Since V_a is finite dimensional (of dimension n), let us pick a basis of functions (w_1, \dots, w_n) in V_a , so that every function $u \in V_a$ can be written as

$$u = u_1 w_1 + \dots + u_n w_n.$$

Then, the equation (DWF) holds iff

$$a(u, w_j) = \tilde{f}(w_j), \quad j = 1, \dots, n,$$

and by plugging $u_1 w_1 + \dots + u_n w_n$ for u , we get a system of k linear equations

$$\sum_{i=1}^n a(w_i, w_j) u_i = \tilde{f}(w_j), \quad 1 \leq j \leq n.$$

Because $a(v, v) \geq \frac{1}{2} \|v\|_{V_a}$, the bilinear form a is symmetric positive definite, and thus the matrix $(a(w_i, w_j))$ is symmetric positive definite, and thus invertible. Therefore, (DWF) has a solution given by a *linear system*!

From a practical point of view, we have to compute the integrals

$$a_{ij} = a(w_i, w_j) = \int_0^1 (w_i' w_j' + c w_i w_j) dx,$$

and

$$b_j = \tilde{f}(w_j) = \int_0^1 f(x) w_j(x) dx.$$

However, if the basis functions are simple enough, this can be done “by hand.” Otherwise, numerical integration methods must be used, but there are some good ones.

Let us also remark that the proof of Theorem 15.1 also shows that the unique solution of (DWF) is the unique minimizer of J over all functions in V_a . It is also possible to compare the approximate solution $u^{(a)} \in V_a$ with the exact solution $u \in V$.

Theorem 15.2. *Suppose $c(x) \geq 0$ for all $x \in [0, 1]$. For every finite-dimensional subspace V_a ($\dim(V_a) = n$) of V , for every basis (w_1, \dots, w_n) of V_a , the following properties hold:*

(1) *There is a unique function $u^{(a)} \in V_a$ such that*

$$a(u^{(a)}, v) = \tilde{f}(v), \quad \text{for all } v \in V_a, \quad (\text{DWF})$$

and if $u^{(a)} = u_1 w_1 + \dots + u_n w_n$, then $\mathbf{u} = (u_1, \dots, u_n)$ is the solution of the linear system

$$A\mathbf{u} = \mathbf{b}, \quad (*)$$

with $A = (a_{ij}) = (a(w_i, w_j))$ and $b_j = \tilde{f}(w_j)$, $1 \leq i, j \leq n$. Furthermore, the matrix $A = (a_{ij})$ is symmetric positive definite.

(2) *The unique solution $u^{(a)} \in V_a$ of (DWF) is the unique minimizer of J over V_a , that is,*

$$J(u^{(a)}) = \inf_{v \in V_a} J(v),$$

(3) *There is a constant C independent of V_a and of the unique solution $u \in V$ of (WF), such that*

$$\|u - u^{(a)}\|_V \leq C \inf_{v \in V_a} \|u - v\|_V.$$

We proved (1) and (2), but we will omit the proof of (3) which can be found in Ciarlet [24].

Let us now give examples of the subspaces V_a used in practice. They usually consist of piecewise polynomial functions.

Pick an integer $N \geq 1$ and subdivide $[0, 1]$ into $N + 1$ intervals $[x_i, x_{i+1}]$, where

$$x_i = hi, \quad h = \frac{1}{N+1}, \quad i = 0, \dots, N+1.$$

We will use the following fact: every polynomial $P(x)$ of degree $2m + 1$ ($m \geq 0$) is completely determined by its values as well as the values of its first m derivatives at two distinct points $\alpha, \beta \in \mathbb{R}$.

There are various ways to prove this. One way is to use the Bernstein basis, because the k th derivative of a polynomial is given by a formula in terms of its control points. For example, for $m = 1$, every degree 3 polynomial can be written as

$$P(x) = (1-x)^3 b_0 + 3(1-x)^2 x b_1 + 3(1-x)x^2 b_2 + x^3 b_3,$$

with $b_0, b_1, b_2, b_3 \in \mathbb{R}$, and we showed that

$$\begin{aligned} P'(0) &= 3(b_1 - b_0) \\ P'(1) &= 3(b_3 - b_2). \end{aligned}$$

Given $P(0)$ and $P(1)$, we determine b_0 and b_3 , and from $P'(0)$ and $P'(1)$, we determine b_1 and b_2 .

In general, for a polynomial of degree m written as

$$P(x) = \sum_{j=0}^m b_j B_j^m(x)$$

in terms of the Bernstein basis $(B_0^m(x), \dots, B_m^m(x))$ with

$$B_j^m(x) = \binom{m}{j} (1-x)^{m-j} x^j,$$

it can be shown that the k th derivative of P at zero is given by

$$P^{(k)}(0) = m(m-1) \cdots (m-k+1) \left(\sum_{i=0}^k \binom{k}{i} (-1)^{k-i} b_i \right),$$

and there is a similar formula for $P^{(k)}(1)$.

Actually, we need to use the Bernstein basis of polynomials $B_k^m[r, s]$, where

$$B_j^m[r, s](x) = \binom{m}{j} \left(\frac{s-x}{s-r} \right)^{m-j} \left(\frac{x-r}{s-r} \right)^j,$$

with $r < s$, in which case

$$P^{(k)}(0) = \frac{m(m-1) \cdots (m-k+1)}{(s-r)^k} \left(\sum_{i=0}^k \binom{k}{i} (-1)^{k-i} b_i \right),$$

with a similar formula for $P^{(k)}(1)$. In our case, we set $r = x_i, s = x_{i+1}$.

Now, if the $2m+2$ values

$$P(0), P^{(1)}(0), \dots, P^{(m)}(0), P(1), P^{(1)}(1), \dots, P^{(m)}(1)$$

are given, we obtain a triangular system that determines uniquely the $2m + 2$ control points b_0, \dots, b_{2m+1} .

Recall that $C^m([0, 1])$ denotes the set of C^m functions f on $[0, 1]$, which means that $f, f^{(1)}, \dots, f^{(m)}$ exist and are continuous on $[0, 1]$.

We define the vector space V_N^m as the subspace of $C^m([0, 1])$ consisting of all functions f such that

1. $f(0) = f(1) = 0$.
2. The restriction of f to $[x_i, x_{i+1}]$ is a polynomial of degree $2m + 1$, for $i = 0, \dots, N$.

Observe that the functions in V_N^0 are the piecewise affine functions f with $f(0) = f(1) = 0$; an example is shown in Figure 15.2.

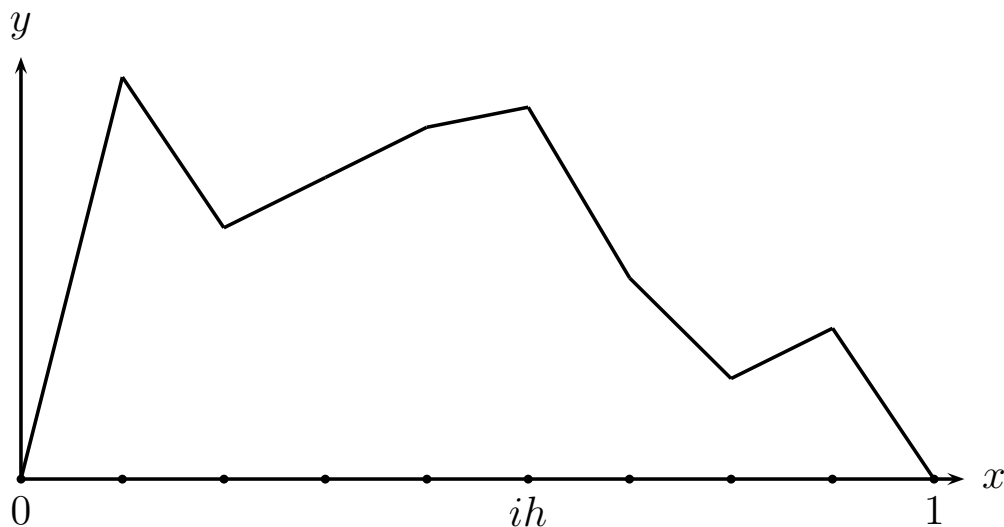


Figure 15.2: A piecewise affine function

This space has dimension N , and a basis consists of the “hat functions” w_i , where the only two nonflat parts of the graph of w_i are the line segments from $(x_{i-1}, 0)$ to $(x_i, 1)$, and from $(x_i, 1)$ to $(x_{i+1}, 0)$, for $i = 1, \dots, N$, see Figure 15.3.

The basis functions w_i have a small support, which is good because in computing the integrals giving $a(w_i, w_j)$, we find that we get a tridiagonal matrix. They also have the nice property that every function $v \in V_N^0$ has the following expression on the basis (w_i) :

$$v(x) = \sum_{i=1}^N v(ih)w_i(x), \quad x \in [0, 1].$$

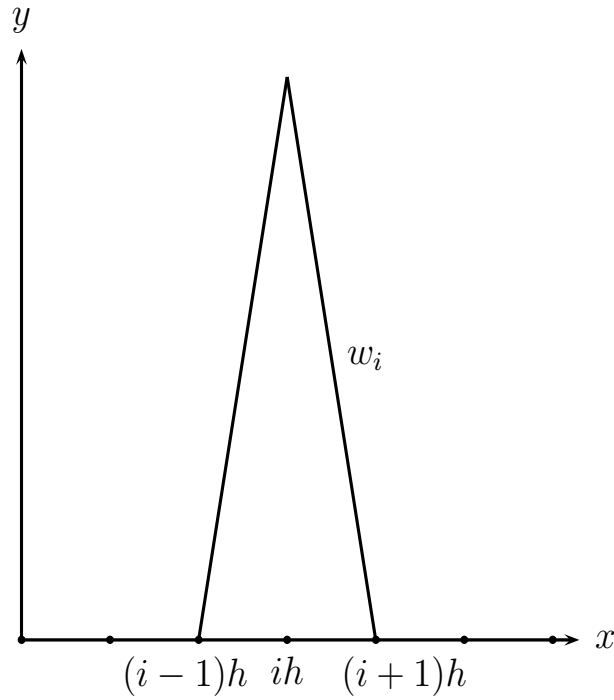


Figure 15.3: A basis “hat function”

In general, it is not hard to see that V_N^m has dimension $mN + 2(m - 1)$.

Going back to our problem (the bending of a beam), assuming that c and f are constant functions, it is not hard to show that the linear system $(*)$ becomes

$$\frac{1}{h} \begin{pmatrix} 2 + \frac{2c}{3}h^2 & -1 + \frac{c}{6}h^2 & & & \\ -1 + \frac{c}{6}h^2 & 2 + \frac{2c}{3}h^2 & -1 + \frac{c}{6}h^2 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 + \frac{c}{6}h^2 & 2 + \frac{2c}{3}h^2 & -1 + \frac{c}{6}h^2 \\ & & & -1 + \frac{c}{6}h^2 & 2 + \frac{2c}{3}h^2 \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_{N-1} \\ u_N \end{pmatrix} = h \begin{pmatrix} f \\ f \\ \vdots \\ f \\ f \end{pmatrix}.$$

We can also find a basis of $2N + 2$ cubic functions for V_N^1 consisting of functions with small support. This basis consists of the N functions w_i^0 and of the $N + 2$ functions w_i^1

uniquely determined by the following conditions:

$$\begin{aligned} w_i^0(x_j) &= \delta_{ij}, & 1 \leq j \leq N, 1 \leq i \leq N \\ (w_i^0)'(x_j) &= 0, & 0 \leq j \leq N+1, 1 \leq i \leq N \\ w_i^1(x_j) &= 0, & 1 \leq j \leq N, 0 \leq i \leq N+1 \\ (w_i^1)'(x_j) &= \delta_{ij}, & 0 \leq j \leq N+1, 0 \leq i \leq N+1 \end{aligned}$$

with $\delta_{ij} = 1$ iff $i = j$ and $\delta_{ij} = 0$ if $i \neq j$. Some of these functions are displayed in Figure 15.4. The function w_i^0 is given explicitly by

$$w_i^0(x) = \frac{1}{h^3}(x - (i-1)h)^2((2i+1)h - 2x), \quad (i-1)h \leq x \leq ih,$$

$$w_i^0(x) = \frac{1}{h^3}((i+1)h - x)^2(2x - (2i-1)h), \quad ih \leq x \leq (i+1)h,$$

for $i = 1, \dots, N$. The function w_j^1 is given explicitly by

$$w_j^1(x) = -\frac{1}{h^2}(ih - x)(x - (i-1)h)^2, \quad (i-1)h \leq x \leq ih,$$

and

$$w_j^1(x) = \frac{1}{h^2}((i+1)h - x)^2(x - ih), \quad ih \leq x \leq (i+1)h,$$

for $j = 0, \dots, N+1$. Furthermore, for every function $v \in V_N^1$, we have

$$v(x) = \sum_{i=1}^N v(ih)w_i^0(x) + \sum_{j=0}^{N+1} v'(jh)w_j^1(x), \quad x \in [0, 1].$$

If we order these basis functions as

$$w_0^1, w_1^0, w_1^1, w_2^0, w_2^1, \dots, w_N^0, w_N^1, w_{N+1}^1,$$

we find that if $c = 0$, the matrix A of the system (*) is tridiagonal by blocks, where the blocks are 2×2 , 2×1 , or 1×2 matrices, and with single entries in the top left and bottom right corner. A different order of the basis vectors would mess up the tridiagonal block structure of A . We leave the details as an exercise.

Let us now take a quick look at a two-dimensional problem, the bending of an elastic membrane.

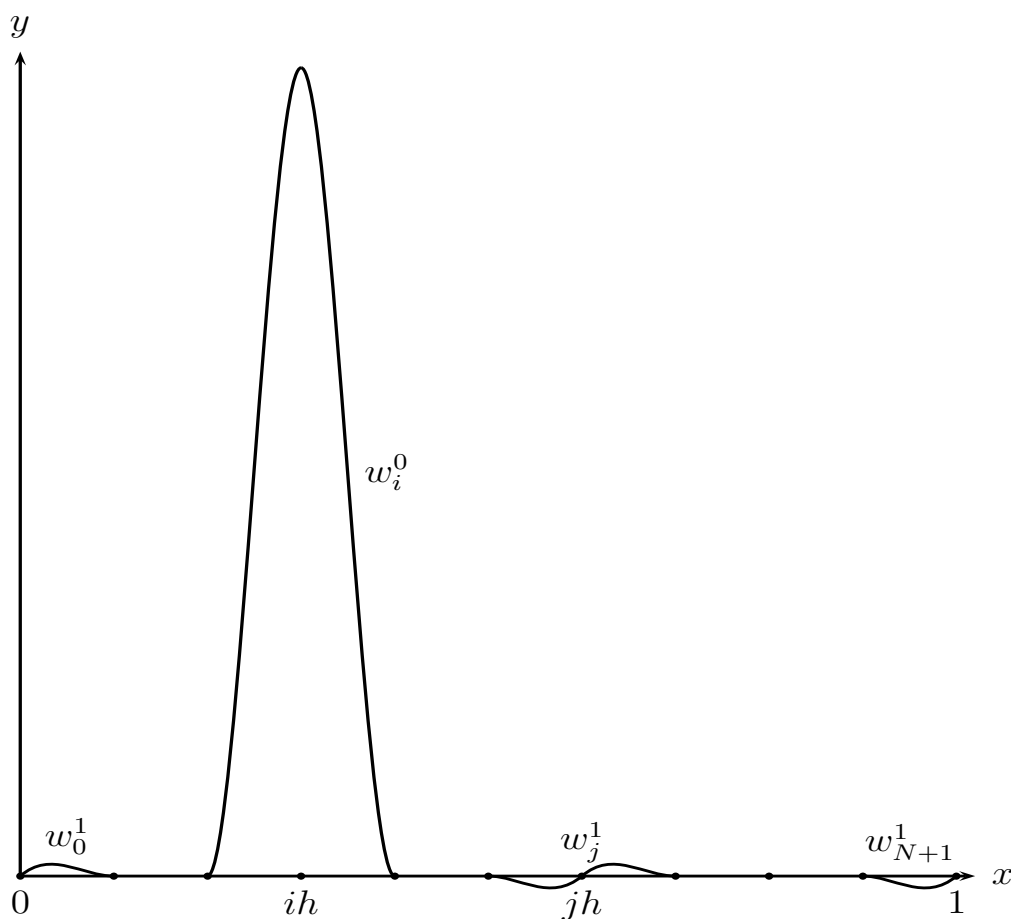


Figure 15.4: The basis functions w_i^0 and w_j^1

15.2 A Two-Dimensional Problem: An Elastic Membrane

Consider an elastic membrane attached to a round contour whose projection on the (x_1, x_2) -plane is the boundary Γ of an open, connected, bounded region Ω in the (x_1, x_2) -plane, as illustrated in Figure 15.5. In other words, we view the membrane as a surface consisting of the set of points (x, z) given by an equation of the form

$$z = u(x),$$

with $x = (x_1, x_2) \in \bar{\Omega}$, where $u: \bar{\Omega} \rightarrow \mathbb{R}$ is some sufficiently regular function, and we think of $u(x)$ as the vertical displacement of this membrane.

We assume that this membrane is under the action of a vertical force $\tau f(x)dx$ per surface element in the horizontal plane (where τ is the tension of the membrane). The problem is

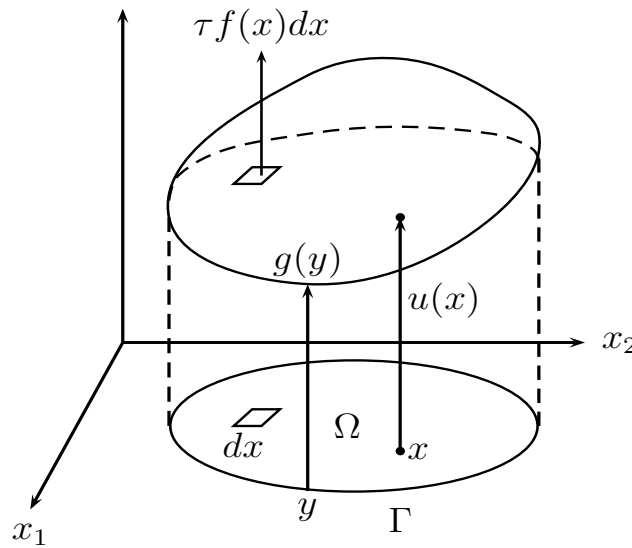


Figure 15.5: An elastic membrane

to find the vertical displacement u as a function of x , for $x \in \overline{\Omega}$. It can be shown (under some assumptions on Ω , Γ , and f), that $u(x)$ is given by a PDE with boundary condition, of the form

$$\begin{aligned} -\Delta u(x) &= f(x), & x \in \Omega \\ u(x) &= g(x), & x \in \Gamma, \end{aligned}$$

where $g: \Gamma \rightarrow \mathbb{R}$ represents the height of the contour of the membrane. We are looking for a function u in $C^2(\Omega) \cap C^1(\overline{\Omega})$. The operator Δ is the *Laplacian*, and it is given by

$$\Delta u(x) = \frac{\partial^2 u}{\partial x_1^2}(x) + \frac{\partial^2 u}{\partial x_2^2}(x).$$

This is an example of a *boundary problem*, since the solution u of the PDE must satisfy the condition $u(x) = g(x)$ on the boundary of the domain Ω . The above equation is known as *Poisson's equation*, and when $f = 0$ as *Laplace's equation*.

It can be proved that if the data f, g and Γ are sufficiently smooth, then the problem has a unique solution.

To get a weak formulation of the problem, first we have to make the boundary condition homogeneous, which means that $g(x) = 0$ on Γ . It turns out that g can be extended to the whole of $\overline{\Omega}$ as some sufficiently smooth function \hat{h} , so we can look for a solution of the form $u - \hat{h}$, but for simplicity, let us assume that the contour of Ω lies in a plane parallel to the

(x_1, x_2) - plane, so that $g = 0$. We let V be the subspace of $C^2(\Omega) \cap C^1(\overline{\Omega})$ consisting of functions v such that $v = 0$ on Γ .

As before, we multiply the PDE by a test function $v \in V$, getting

$$-\Delta u(x)v(x) = f(x)v(x),$$

and we “integrate by parts.” In this case, this means that we use a version of Stokes formula known as *Green’s first identity*, which says that

$$\int_{\Omega} -\Delta u v \, dx = \int_{\Omega} (\text{grad } u) \cdot (\text{grad } v) \, dx - \int_{\Gamma} (\text{grad } u) \cdot n v \, d\sigma$$

(where n denotes the outward pointing unit normal to the surface). Because $v = 0$ on Γ , the integral \int_{Γ} drops out, and we get an equation of the form

$$a(u, v) = \tilde{f}(v) \quad \text{for all } v \in V,$$

where a is the bilinear form given by

$$a(u, v) = \int_{\Omega} \left(\frac{\partial u}{\partial x_1} \frac{\partial v}{\partial x_1} + \frac{\partial u}{\partial x_2} \frac{\partial v}{\partial x_2} \right) dx$$

and \tilde{f} is the linear form given by

$$\tilde{f}(v) = \int_{\Omega} f v \, dx.$$

We get the same equation as in section 15.2, but over a set of functions defined on a two-dimensional domain. As before, we can choose a finite-dimensional subspace V_a of V and consider the discrete problem with respect to V_a . Again, if we pick a basis (w_1, \dots, w_n) of V_a , a vector $u = u_1 w_1 + \dots + u_n w_n$ is a solution of the Weak Formulation of our problem iff $\mathbf{u} = (u_1, \dots, u_n)$ is a solution of the linear system

$$A\mathbf{u} = b,$$

with $A = (a(w_i, w_j))$ and $b = (\tilde{f}(w_j))$. However, the integrals that give the entries in A and b are much more complicated.

An approach to deal with this problem is the *method of finite elements*. The idea is to also discretize the boundary curve Γ . If we assume that Γ is a *polygonal line*, then we can *triangulate* the domain Ω , and then we consider spaces of functions which are piecewise defined on the triangles of the triangulation of Ω . The simplest functions are piecewise affine and look like tents erected above groups of triangles. Again, we can define base functions with small support, so that the matrix A is tridiagonal by blocks.

The finite element method is a vast subject and it is presented in many books of various degrees of difficulty and obscurity. Let us simply state three important requirements of the finite element method:

1. “Good” triangulations must be found. This in itself is a vast research topic. Delaunay triangulations are good candidates.
2. “Good” spaces of functions must be found; typically piecewise polynomials and splines.
3. “Good” bases consisting of functions with small support must be found, so that integrals can be easily computed and sparse banded matrices arise.

We now consider boundary problems where the solution varies with time.

15.3 Time-Dependent Boundary Problems: The Wave Equation

Consider a homogeneous string (or rope) of constant cross-section, of length L , and stretched (in a vertical plane) between its two ends which are assumed to be fixed and located along the x -axis at $x = 0$ and at $x = L$.

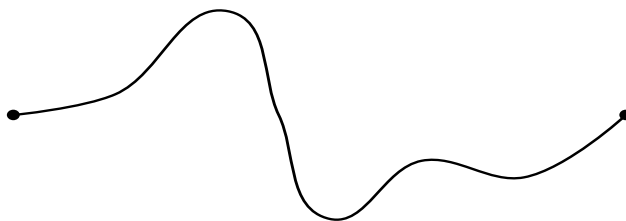


Figure 15.6: A vibrating string

The string is subjected to a transverse force $\tau f(x)dx$ per element of length dx (where τ is the tension of the string). We would like to investigate the small displacements of the string in the vertical plane, that is, how it vibrates.

Thus, we seek a function $u(x, t)$ defined for $t \geq 0$ and $x \in [0, L]$, such that $u(x, t)$ represents the vertical deformation of the string at the abscissa x and at time t .

It can be shown that u must satisfy the following PDE

$$\frac{1}{c^2} \frac{\partial^2 u}{\partial t^2}(x, t) - \frac{\partial^2 u}{\partial x^2}(x, t) = f(x, t), \quad 0 < x < L, \quad t > 0,$$

with $c = \sqrt{\tau/\rho}$, where ρ is the linear density of the string, known as the *one-dimensional wave equation*.

Furthermore, the initial shape of the string is known at $t = 0$, as well as the distribution of the initial velocities along the string; in other words, there are two functions $u_{i,0}$ and $u_{i,1}$ such that

$$\begin{aligned} u(x, 0) &= u_{i,0}(x), \quad 0 \leq x \leq L, \\ \frac{\partial u}{\partial t}(x, 0) &= u_{i,1}(x), \quad 0 \leq x \leq L. \end{aligned}$$

For example, if the string is simply released from its given starting position, we have $u_{i,1} = 0$. Lastly, because the ends of the string are fixed, we must have

$$u(0, t) = u(L, t) = 0, \quad t \geq 0.$$

Consequently, we look for a function $u: \mathbb{R}_+ \times [0, L] \rightarrow \mathbb{R}$ satisfying the following conditions:

$$\begin{aligned} \frac{1}{c^2} \frac{\partial^2 u}{\partial t^2}(x, t) - \frac{\partial^2 u}{\partial x^2}(x, t) &= f(x, t), \quad 0 < x < L, \quad t > 0, \\ u(0, t) &= u(L, t) = 0, \quad t \geq 0 \quad (\text{boundary condition}), \\ u(x, 0) &= u_{i,0}(x), \quad 0 \leq x \leq L \quad (\text{intitial condition}), \\ \frac{\partial u}{\partial t}(x, 0) &= u_{i,1}(x), \quad 0 \leq x \leq L \quad (\text{intitial condition}). \end{aligned}$$

This is an example of a *time-dependent boundary-value problem*, with two *initial conditions*.

To simplify the problem, assume that $f = 0$, which amounts to neglecting the effect of gravity. In this case, our PDE becomes

$$\frac{1}{c^2} \frac{\partial^2 u}{\partial t^2}(x, t) - \frac{\partial^2 u}{\partial x^2}(x, t) = 0, \quad 0 < x < L, \quad t > 0,$$

Let us try our trick of multiplying by a test function v depending only on x , C^1 on $[0, L]$, and such that $v(0) = v(L) = 0$, and integrate by parts. We get the equation

$$\int_0^L \frac{\partial^2 u}{\partial t^2}(x, t) v(x) dx - c^2 \int_0^L \frac{\partial^2 u}{\partial x^2}(x, t) v(x) dx = 0.$$

For the first term, we get

$$\begin{aligned} \int_0^L \frac{\partial^2 u}{\partial t^2}(x, t) v(x) dx &= \int_0^L \frac{\partial^2}{\partial t^2} [u(x, t) v(x)] dx \\ &= \frac{d^2}{dt^2} \int_0^L u(x, t) v(x) dx \\ &= \frac{d^2}{dt^2} \langle u, v \rangle, \end{aligned}$$

where $\langle u, v \rangle$ is the inner product in $L^2([0, L])$. The fact that it is legitimate to move $\partial^2/\partial t^2$ outside of the integral needs to be justified rigorously, but we won't do it here.

For the second term, we get

$$-\int_0^L \frac{\partial^2 u}{\partial x^2}(x, t)v(x)dx = -\left[\frac{\partial u}{\partial x}(x, t)v(x)\right]_{x=0}^{x=L} + \int_0^L \frac{\partial u}{\partial x}(x, t)\frac{dv}{dx}(x)dx,$$

and because $v \in V$, we have $v(0) = v(L) = 0$, so we obtain

$$-\int_0^L \frac{\partial^2 u}{\partial x^2}(x, t)v(x)dx = \int_0^L \frac{\partial u}{\partial x}(x, t)\frac{dv}{dx}(x)dx.$$

Our integrated equation becomes

$$\frac{d^2}{dt^2}\langle u, v \rangle + c^2 \int_0^L \frac{\partial u}{\partial x}(x, t)\frac{dv}{dx}(x)dx = 0, \quad \text{for all } v \in V \quad \text{and all } t \geq 0.$$

It is natural to introduce the bilinear form $a: V \times V \rightarrow \mathbb{R}$ given by

$$a(u, v) = \int_0^L \frac{\partial u}{\partial x}(x, t)\frac{\partial v}{\partial x}(x, t)dx,$$

where, for every $t \in \mathbb{R}_+$, the functions $u(x, t)$ and (v, t) belong to V . Actually, we have to replace V by the subspace of the Sobolev space $H_0^1(0, L)$ consisting of the functions such that $v(0) = v(L) = 0$. Then, the weak formulation (variational formulation) of our problem is this:

Find a function $u \in V$ such that

$$\begin{aligned} \frac{d^2}{dt^2}\langle u, v \rangle + a(u, v) &= 0, \quad \text{for all } v \in V \quad \text{and all } t \geq 0 \\ u(x, 0) &= u_{i,0}(x), \quad 0 \leq x \leq L \quad (\text{intitial condition}), \\ \frac{\partial u}{\partial t}(x, 0) &= u_{i,1}(x), \quad 0 \leq x \leq L \quad (\text{intitial condition}). \end{aligned}$$

It can be shown that there is a positive constant $\alpha > 0$ such that

$$a(u, u) \geq \alpha \|u\|_{H_0^1}^2 \quad \text{for all } u \in V$$

(Poincaré's inequality), which shows that a is positive definite on V . The above method is known as the method of *Rayleigh-Ritz*.

A study of the above equation requires some sophisticated tools of analysis which go far beyond the scope of these notes. Let us just say that there is a countable sequence of solutions with separated variables of the form

$$u_k^{(1)} = \sin\left(\frac{k\pi x}{L}\right) \cos\left(\frac{k\pi ct}{L}\right), \quad u_k^{(2)} = \sin\left(\frac{k\pi x}{L}\right) \sin\left(\frac{k\pi ct}{L}\right), \quad k \in \mathbb{N}_+,$$

called *modes* (or *normal modes*). Complete solutions of the problem are series obtained by combining the normal modes, and they are of the form

$$u(x, t) = \sum_{k=1}^{\infty} \sin\left(\frac{k\pi x}{L}\right) \left(A_k \cos\left(\frac{k\pi ct}{L}\right) + B_k \sin\left(\frac{k\pi ct}{L}\right) \right),$$

where the coefficients A_k, B_k are determined from the Fourier series of $u_{i,0}$ and $u_{i,1}$.

We now consider discrete approximations of our problem. As before, consider a finite dimensional subspace V_a of V and assume that we have approximations $u_{a,0}$ and $u_{a,1}$ of $u_{i,0}$ and $u_{i,1}$. If we pick a basis (w_1, \dots, w_n) of V_a , then we can write our unknown function $u(x, t)$ as

$$u(x, t) = u_1(t)w_1 + \dots + u_n(t)w_n,$$

where u_1, \dots, u_n are functions of t . Then, if we write $\mathbf{u} = (u_1, \dots, u_n)$, the discrete version of our problem is

$$\begin{aligned} A \frac{d^2 \mathbf{u}}{dt^2} + K \mathbf{u} &= 0, \\ u(x, 0) &= u_{a,0}(x), \quad 0 \leq x \leq L, \\ \frac{\partial u}{\partial t}(x, 0) &= u_{a,1}(x), \quad 0 \leq x \leq L, \end{aligned}$$

where $A = (\langle w_i, w_j \rangle)$ and $K = (a(w_i, w_j))$ are two symmetric matrices, called the *mass matrix* and the *stiffness matrix*, respectively. In fact, because a and the inner product $\langle -, - \rangle$ are positive definite, these matrices are also positive definite.

We have made some progress since we now have a system of ODE's, and we can solve it by analogy with the scalar case. So, we look for solutions of the form $\mathbf{U} \cos \omega t$ (or $\mathbf{U} \sin \omega t$), where \mathbf{U} is an n -dimensional vector. We find that we should have

$$(K - \omega^2 A) \mathbf{U} \cos \omega t = 0,$$

which implies that ω must be a solution of the equation

$$K \mathbf{U} = \omega^2 A \mathbf{U}.$$

Thus, we have to find some λ such that

$$K \mathbf{U} = \lambda A \mathbf{U},$$

a problem known as a *generalized eigenvalue problem*, since the ordinary eigenvalue problem for K is

$$K \mathbf{U} = \lambda \mathbf{U}.$$

Fortunately, because A is SPD, we can reduce this generalized eigenvalue problem to a standard eigenvalue problem. A good way to do so is to use a Cholesky decomposition of A as

$$A = LL^\top,$$

where L is a lower triangular matrix (see Theorem 6.10). Because A is SPD, it is invertible, so L is also invertible, and

$$K\mathbf{U} = \lambda A\mathbf{U} = \lambda LL^\top \mathbf{U}$$

yields

$$L^{-1}K\mathbf{U} = \lambda L^\top \mathbf{U},$$

which can also be written as

$$L^{-1}K(L^\top)^{-1}L^\top \mathbf{U} = \lambda L^\top \mathbf{U}.$$

Then, if we make the change of variable

$$\mathbf{Y} = L^\top \mathbf{U},$$

using the fact $(L^\top)^{-1} = (L^{-1})^\top$, the above equation is equivalent to

$$L^{-1}K(L^{-1})^\top \mathbf{Y} = \lambda \mathbf{Y},$$

a standard eigenvalue problem for the matrix $\hat{K} = L^{-1}K(L^{-1})^\top$. Furthermore, we know from Section 6.3 that since K is SPD and L^{-1} is invertible, the matrix $\hat{K} = L^{-1}K(L^{-1})^\top$ is also SPD.

Consequently, \hat{K} has positive real eigenvalues $(\omega_1^2, \dots, \omega_n^2)$ (not necessarily distinct) and it can be diagonalized with respect to an orthonormal basis of eigenvectors, say $\mathbf{Y}^1, \dots, \mathbf{Y}^n$. Then, since $\mathbf{Y} = L^\top \mathbf{U}$, the vectors

$$\mathbf{U}^i = (L^\top)^{-1} \mathbf{Y}^i, \quad i = 1, \dots, n,$$

are linearly independent and are solutions of the generalized eigenvalue problem; that is,

$$K\mathbf{U}^i = \omega_i^2 A\mathbf{U}^i, \quad i = 1, \dots, n.$$

More is true. Because the vectors $\mathbf{Y}^1, \dots, \mathbf{Y}^n$ are orthonormal, and because $\mathbf{Y}^i = L^\top \mathbf{U}^i$, from

$$(\mathbf{Y}^i)^\top \mathbf{Y}^j = \delta_{ij},$$

we get

$$(\mathbf{U}^i)^\top LL^\top \mathbf{U}^j = \delta_{ij}, \quad 1 \leq i, j \leq n,$$

and since $A = LL^\top$, this yields

$$(\mathbf{U}^i)^\top A\mathbf{U}^j = \delta_{ij}, \quad 1 \leq i, j \leq n.$$

This suggests defining the functions $U^i \in V_a$ by

$$U^i = \sum_{k=1}^n \mathbf{U}_k^i w_k.$$

Then, it is immediate to check that

$$a(U^i, U^j) = (\mathbf{U}^i)^\top A \mathbf{U}^j = \delta_{ij},$$

which means that the functions (U^1, \dots, U^n) form an orthonormal basis of V_a for the inner product a . The functions $U^i \in V_a$ are called *modes* (or *modal vectors*).

As a final step, let us look again for a solution of our discrete weak formulation of the problem, this time expressing the unknown solution $u(x, t)$ over the modal basis (U^1, \dots, U^n) , say

$$u = \sum_{j=1}^n \tilde{u}_j(t) U^j,$$

where each \tilde{u}_j is a function of t . Because

$$u = \sum_{j=1}^n \tilde{u}_j(t) U^j = \sum_{j=1}^n \tilde{u}_j(t) \left(\sum_{k=1}^n \mathbf{U}_k^j w_k \right) = \sum_{k=1}^n \left(\sum_{j=1}^n \tilde{u}_j(t) \mathbf{U}_k^j \right) w_k,$$

if we write $\mathbf{u} = (u_1, \dots, u_n)$ with $u_k = \sum_{j=1}^n \tilde{u}_j(t) \mathbf{U}_k^j$ for $k = 1, \dots, n$, we see that

$$\mathbf{u} = \sum_{j=1}^n \tilde{u}_j \mathbf{U}^j,$$

so using the fact that

$$K \mathbf{U}^j = \omega_j^2 A \mathbf{U}^j, \quad j = 1, \dots, n,$$

the equation

$$A \frac{d^2 \mathbf{u}}{dt^2} + K \mathbf{u} = 0$$

yields

$$\sum_{j=1}^n [(\tilde{u}_j)'' + \omega_j^2 \tilde{u}_j] A \mathbf{U}^j = 0.$$

Since A is invertible and since $(\mathbf{U}^1, \dots, \mathbf{U}^n)$ are linearly independent, the vectors $(A \mathbf{U}^1, \dots, A \mathbf{U}^n)$ are linearly independent, and consequently we get the system of n ODEs'

$$(\tilde{u}_j)'' + \omega_j^2 \tilde{u}_j = 0, \quad 1 \leq j \leq n.$$

Each of these equations has a well-known solution of the form

$$\tilde{u}_j = A_j \cos \omega_j t + B_j \sin \omega_j t.$$

Therefore, the solution of our approximation problem is given by

$$u = \sum_{j=1}^n (A_j \cos \omega_j t + B_j \sin \omega_j t) U^j,$$

and the constants A_j, B_j are obtained from the initial conditions

$$\begin{aligned} u(x, 0) &= u_{a,0}(x), \quad 0 \leq x \leq L, \\ \frac{\partial u}{\partial t}(x, 0) &= u_{a,1}(x), \quad 0 \leq x \leq L, \end{aligned}$$

by expressing $u_{a,0}$ and $u_{a,1}$ on the modal basis (U^1, \dots, U^n) . Furthermore, the modal functions (U^1, \dots, U^n) form an orthonormal basis of V_a for the inner product a .

If we use the vector space V_N^0 of piecewise affine functions, we find that the matrices A and K are familiar! Indeed,

$$A = \frac{1}{h} \begin{pmatrix} 2 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & 2 \end{pmatrix}$$

and

$$K = \frac{h}{6} \begin{pmatrix} 4 & 1 & 0 & 0 & 0 \\ 1 & 4 & 1 & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & 1 & 4 & 1 \\ 0 & 0 & 0 & 1 & 4 \end{pmatrix}.$$

To conclude this section, let us discuss briefly the wave equation for an elastic membrane, as described in Section 15.2. This time, we look for a function $u: \mathbb{R}_+ \times \Omega \rightarrow \mathbb{R}$ satisfying the following conditions:

$$\begin{aligned} \frac{1}{c^2} \frac{\partial^2 u}{\partial t^2}(x, t) - \Delta u(x, t) &= f(x, t), \quad x \in \Omega, \quad t > 0, \\ u(x, t) &= 0, \quad x \in \Gamma, \quad t \geq 0 \quad (\text{boundary condition}), \\ u(x, 0) &= u_{i,0}(x), \quad x \in \Omega \quad (\text{initial condition}), \\ \frac{\partial u}{\partial t}(x, 0) &= u_{i,1}(x), \quad x \in \Omega \quad (\text{initial condition}). \end{aligned}$$

Assuming that $f = 0$, we look for solutions in the subspace V of the Sobolev space $H_0^1(\bar{\Omega})$ consisting of functions v such that $v = 0$ on Γ . Multiplying by a test function $v \in V$ and using Green's first identity, we get the weak formulation of our problem:

Find a function $u \in V$ such that

$$\begin{aligned} \frac{d^2}{dt^2} \langle u, v \rangle + a(u, v) &= 0, \quad \text{for all } v \in V \text{ and all } t \geq 0 \\ u(x, 0) &= u_{i,0}(x), \quad x \in \Omega \quad (\text{intitial condition}), \\ \frac{\partial u}{\partial t}(x, 0) &= u_{i,1}(x), \quad x \in \Omega \quad (\text{intitial condition}), \end{aligned}$$

where $a: V \times V \rightarrow \mathbb{R}$ is the bilinear form given by

$$a(u, v) = \int_{\Omega} \left(\frac{\partial u}{\partial x_1} \frac{\partial v}{\partial x_1} + \frac{\partial u}{\partial x_2} \frac{\partial v}{\partial x_2} \right) dx,$$

and

$$\langle u, v \rangle = \int_{\Omega} uv dx.$$

As usual, we find approximations of our problem by using finite dimensional subspaces V_a of V . Picking some basis (w_1, \dots, w_n) of V_a , and triangulating Ω , as before, we obtain the equation

$$\begin{aligned} A \frac{d^2 \mathbf{u}}{dt^2} + K \mathbf{u} &= 0, \\ u(x, 0) &= u_{a,0}(x), \quad x \in \Gamma, \\ \frac{\partial u}{\partial t}(x, 0) &= u_{a,1}(x), \quad x \in \Gamma, \end{aligned}$$

where $A = (\langle w_i, w_j \rangle)$ and $K = (a(w_i, w_j))$ are two symmetric positive definite matrices.

In principle, the problem is solved, but, it may be difficult to find good spaces V_a , good triangulations of Ω , and good bases of V_a , to be able to compute the matrices A and K , and to ensure that they are sparse.

Chapter 16

Singular Value Decomposition and Polar Form

16.1 Singular Value Decomposition for Square Matrices

In this section, we assume that we are dealing with real Euclidean spaces. Let $f: E \rightarrow E$ be any linear map. In general, it may not be possible to diagonalize f . We show that every linear map can be diagonalized if we are willing to use *two* orthonormal bases. This is the celebrated *singular value decomposition (SVD)*. A close cousin of the SVD is the *polar form* of a linear map, which shows how a linear map can be decomposed into its purely rotational component (perhaps with a flip) and its purely stretching part.

The key observation is that $f^* \circ f$ is self-adjoint, since

$$\langle (f^* \circ f)(u), v \rangle = \langle f(u), f(v) \rangle = \langle u, (f^* \circ f)(v) \rangle.$$

Similarly, $f \circ f^*$ is self-adjoint.

The fact that $f^* \circ f$ and $f \circ f^*$ are self-adjoint is very important, because it implies that $f^* \circ f$ and $f \circ f^*$ can be diagonalized and that they have real eigenvalues. In fact, these eigenvalues are all nonnegative. Indeed, if u is an eigenvector of $f^* \circ f$ for the eigenvalue λ , then

$$\langle (f^* \circ f)(u), u \rangle = \langle f(u), f(u) \rangle$$

and

$$\langle (f^* \circ f)(u), u \rangle = \lambda \langle u, u \rangle,$$

and thus

$$\lambda \langle u, u \rangle = \langle f(u), f(u) \rangle,$$

which implies that $\lambda \geq 0$, since $\langle -, - \rangle$ is positive definite. A similar proof applies to $f \circ f^*$. Thus, the eigenvalues of $f^* \circ f$ are of the form $\sigma_1^2, \dots, \sigma_r^2$ or 0, where $\sigma_i > 0$, and similarly for $f \circ f^*$.

The above considerations also apply to any linear map $f: E \rightarrow F$ between two Euclidean spaces $(E, \langle -, - \rangle_1)$ and $(F, \langle -, - \rangle_2)$. Recall that the adjoint $f^*: F \rightarrow E$ of f is the unique linear map f^* such that

$$\langle f(u), v \rangle_2 = \langle u, f^*(v) \rangle_1, \quad \text{for all } u \in E \text{ and all } v \in F.$$

Then, $f^* \circ f$ and $f \circ f^*$ are self-adjoint (the proof is the same as in the previous case), and the eigenvalues of $f^* \circ f$ and $f \circ f^*$ are nonnegative. If λ is an eigenvalue of $f^* \circ f$ and $u (\neq 0)$ is a corresponding eigenvector, we have

$$\langle (f^* \circ f)(u), u \rangle_1 = \langle f(u), f(u) \rangle_2,$$

and also

$$\langle (f^* \circ f)(u), u \rangle_1 = \lambda \langle u, u \rangle_1,$$

so

$$\lambda \langle u, u \rangle_1 = \langle f(u), f(u) \rangle_2,$$

which implies that $\lambda \geq 0$. A similar proof applies to $f \circ f^*$. The situation is even better, since we will show shortly that $f^* \circ f$ and $f \circ f^*$ have the same nonzero eigenvalues.

Remark: Given any two linear maps $f: E \rightarrow F$ and $g: F \rightarrow E$, where $\dim(E) = n$ and $\dim(F) = m$, it can be shown that

$$\lambda^m \det(\lambda I_n - g \circ f) = \lambda^n \det(\lambda I_m - f \circ g),$$

and thus $g \circ f$ and $f \circ g$ always have the same nonzero eigenvalues!

Definition 16.1. Given any linear map $f: E \rightarrow F$, the square roots $\sigma_i > 0$ of the positive eigenvalues of $f^* \circ f$ (and $f \circ f^*$) are called the *singular values* of f .

Definition 16.2. A self-adjoint linear map $f: E \rightarrow E$ whose eigenvalues are nonnegative is called *positive semidefinite* (or *positive*), and if f is also invertible, f is said to be *positive definite*. In the latter case, every eigenvalue of f is strictly positive.

If $f: E \rightarrow F$ is any linear map, we just showed that $f^* \circ f$ and $f \circ f^*$ are positive semidefinite self-adjoint linear maps. This fact has the remarkable consequence that every linear map has two important decompositions:

1. The polar form.
2. The singular value decomposition (SVD).

The wonderful thing about the singular value decomposition is that there exist two orthonormal bases (u_1, \dots, u_n) and (v_1, \dots, v_m) such that, with respect to these bases, f is a diagonal matrix consisting of the singular values of f , or 0. Thus, in some sense, f can always be diagonalized with respect to *two* orthonormal bases. The SVD is also a useful tool for solving overdetermined linear systems in the least squares sense and for data analysis, as we show later on.

First, we show some useful relationships between the kernels and the images of f , f^* , $f^* \circ f$, and $f \circ f^*$. Recall that if $f: E \rightarrow F$ is a linear map, the *image* $\text{Im } f$ of f is the subspace $f(E)$ of F , and the *rank* of f is the dimension $\dim(\text{Im } f)$ of its image. Also recall that (Theorem 4.11))

$$\dim(\text{Ker } f) + \dim(\text{Im } f) = \dim(E),$$

and that (Propositions 10.9 and 12.10) for every subspace W of E ,

$$\dim(W) + \dim(W^\perp) = \dim(E).$$

Proposition 16.1. *Given any two Euclidean spaces E and F , where E has dimension n and F has dimension m , for any linear map $f: E \rightarrow F$, we have*

$$\begin{aligned} \text{Ker } f &= \text{Ker } (f^* \circ f), \\ \text{Ker } f^* &= \text{Ker } (f \circ f^*), \\ \text{Ker } f &= (\text{Im } f^*)^\perp, \\ \text{Ker } f^* &= (\text{Im } f)^\perp, \\ \dim(\text{Im } f) &= \dim(\text{Im } f^*), \end{aligned}$$

and f , f^* , $f^* \circ f$, and $f \circ f^*$ have the same rank.

Proof. To simplify the notation, we will denote the inner products on E and F by the same symbol $\langle -, - \rangle$ (to avoid subscripts). If $f(u) = 0$, then $(f^* \circ f)(u) = f^*(f(u)) = f^*(0) = 0$, and so $\text{Ker } f \subseteq \text{Ker } (f^* \circ f)$. By definition of f^* , we have

$$\langle f(u), f(u) \rangle = \langle (f^* \circ f)(u), u \rangle$$

for all $u \in E$. If $(f^* \circ f)(u) = 0$, since $\langle -, - \rangle$ is positive definite, we must have $f(u) = 0$, and so $\text{Ker } (f^* \circ f) \subseteq \text{Ker } f$. Therefore,

$$\text{Ker } f = \text{Ker } (f^* \circ f).$$

The proof that $\text{Ker } f^* = \text{Ker } (f \circ f^*)$ is similar.

By definition of f^* , we have

$$\langle f(u), v \rangle = \langle u, f^*(v) \rangle \quad \text{for all } u \in E \text{ and all } v \in F. \quad (*)$$

This immediately implies that

$$\text{Ker } f = (\text{Im } f^*)^\perp \quad \text{and} \quad \text{Ker } f^* = (\text{Im } f)^\perp.$$

Let us explain why $\text{Ker } f = (\text{Im } f^*)^\perp$, the proof of the other equation being similar.

Because the inner product is positive definite, for every $u \in E$, we have
 $u \in \text{Ker } f$
iff $f(u) = 0$
iff $\langle f(u), v \rangle = 0$ for all v ,
by (*) iff $\langle u, f^*(v) \rangle = 0$ for all v ,
iff $u \in (\text{Im } f^*)^\perp$.

Since

$$\dim(\text{Im } f) = n - \dim(\text{Ker } f)$$

and

$$\dim(\text{Im } f^*) = n - \dim((\text{Im } f^*)^\perp),$$

from

$$\text{Ker } f = (\text{Im } f^*)^\perp$$

we also have

$$\dim(\text{Ker } f) = \dim((\text{Im } f^*)^\perp),$$

from which we obtain

$$\dim(\text{Im } f) = \dim(\text{Im } f^*).$$

Since

$$\dim(\text{Ker } (f^* \circ f)) + \dim(\text{Im } (f^* \circ f)) = \dim(E),$$

$\text{Ker } (f^* \circ f) = \text{Ker } f$ and $\text{Ker } f = (\text{Im } f^*)^\perp$, we get

$$\dim((\text{Im } f^*)^\perp) + \dim(\text{Im } (f^* \circ f)) = \dim(E).$$

Since

$$\dim((\text{Im } f^*)^\perp) + \dim(\text{Im } f^*) = \dim(E),$$

we deduce that

$$\dim(\text{Im } f^*) = \dim(\text{Im } (f^* \circ f)).$$

A similar proof shows that

$$\dim(\text{Im } f) = \dim(\text{Im } (f \circ f^*)).$$

Consequently, f , f^* , $f^* \circ f$, and $f \circ f^*$ have the same rank. □

We will now prove that every square matrix has an SVD. Stronger results can be obtained if we first consider the polar form and then derive the SVD from it (there are uniqueness properties of the polar decomposition). For our purposes, uniqueness results are not as important so we content ourselves with existence results, whose proofs are simpler. Readers interested in a more general treatment are referred to [44].

The early history of the singular value decomposition is described in a fascinating paper by Stewart [101]. The SVD is due to Beltrami and Camille Jordan independently (1873, 1874). Gauss is the grandfather of all this, for his work on least squares (1809, 1823) (but Legendre also published a paper on least squares!). Then come Sylvester, Schmidt, and Hermann Weyl. Sylvester's work was apparently "opaque." He gave a computational method to find an SVD. Schmidt's work really has to do with integral equations and symmetric and asymmetric kernels (1907). Weyl's work has to do with perturbation theory (1912). Autonne came up with the polar decomposition (1902, 1915). Eckart and Young extended SVD to rectangular matrices (1936, 1939).

Theorem 16.2. (*Singular value decomposition*) *For every real $n \times n$ matrix A there are two orthogonal matrices U and V and a diagonal matrix D such that $A = VDU^\top$, where D is of the form*

$$D = \begin{pmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ \vdots & \vdots & & \vdots \\ & & & \dots & \sigma_n \end{pmatrix},$$

where $\sigma_1, \dots, \sigma_r$ are the singular values of f , i.e., the (positive) square roots of the nonzero eigenvalues of $A^\top A$ and AA^\top , and $\sigma_{r+1} = \dots = \sigma_n = 0$. The columns of U are eigenvectors of $A^\top A$, and the columns of V are eigenvectors of AA^\top .

Proof. Since $A^\top A$ is a symmetric matrix, in fact, a positive semidefinite matrix, there exists an orthogonal matrix U such that

$$A^\top A = UD^2U^\top,$$

with $D = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0)$, where $\sigma_1^2, \dots, \sigma_r^2$ are the nonzero eigenvalues of $A^\top A$, and where r is the rank of A ; that is, $\sigma_1, \dots, \sigma_r$ are the singular values of A . It follows that

$$U^\top A^\top A U = (AU)^\top A U = D^2,$$

and if we let f_j be the j th column of AU for $j = 1, \dots, n$, then we have

$$\langle f_i, f_j \rangle = \sigma_i^2 \delta_{ij}, \quad 1 \leq i, j \leq r$$

and

$$f_j = 0, \quad r+1 \leq j \leq n.$$

If we define (v_1, \dots, v_r) by

$$v_j = \sigma_j^{-1} f_j, \quad 1 \leq j \leq r,$$

then we have

$$\langle v_i, v_j \rangle = \delta_{ij}, \quad 1 \leq i, j \leq r,$$

so complete (v_1, \dots, v_r) into an orthonormal basis $(v_1, \dots, v_r, v_{r+1}, \dots, v_n)$ (for example, using Gram–Schmidt). Now, since $f_j = \sigma_j v_j$ for $j = 1, \dots, r$, we have

$$\langle v_i, f_j \rangle = \sigma_j \langle v_i, v_j \rangle = \sigma_j \delta_{ij}, \quad 1 \leq i \leq n, \quad 1 \leq j \leq r$$

and since $f_j = 0$ for $j = r+1, \dots, n$,

$$\langle v_i, f_j \rangle = 0 \quad 1 \leq i \leq n, \quad r+1 \leq j \leq n.$$

If V is the matrix whose columns are v_1, \dots, v_n , then V is orthogonal and the above equations prove that

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

which yields $A = V D U^\top$, as required.

The equation $A = V D U^\top$ implies that

$$A^\top A = U D^2 U^\top, \quad A A^\top = V D^2 V^\top,$$

which shows that $A^\top A$ and $A A^\top$ have the same eigenvalues, that the columns of U are eigenvectors of $A^\top A$, and that the columns of V are eigenvectors of $A A^\top$. \square

Theorem 16.2 suggests the following definition.

Definition 16.3. A triple (U, D, V) such that $A = V D U^\top$, where U and V are orthogonal and D is a diagonal matrix whose entries are nonnegative (it is positive semidefinite) is called a *singular value decomposition (SVD)* of A .

The proof of Theorem 16.2 shows that there are two orthonormal bases (u_1, \dots, u_n) and (v_1, \dots, v_n) , where (u_1, \dots, u_n) are eigenvectors of $A^\top A$ and (v_1, \dots, v_n) are eigenvectors of $A A^\top$. Furthermore, (u_1, \dots, u_r) is an orthonormal basis of $\text{Im } A^\top$, (u_{r+1}, \dots, u_n) is an orthonormal basis of $\text{Ker } A$, (v_1, \dots, v_r) is an orthonormal basis of $\text{Im } A$, and (v_{r+1}, \dots, v_n) is an orthonormal basis of $\text{Ker } A^\top$.

Using a remark made in Chapter 3, if we denote the columns of U by u_1, \dots, u_n and the columns of V by v_1, \dots, v_n , then we can write

$$A = V D U^\top = \sigma_1 v_1 u_1^\top + \dots + \sigma_r v_r u_r^\top.$$

As a consequence, if r is a lot smaller than n (we write $r \ll n$), we see that A can be reconstructed from U and V using a much smaller number of elements. This idea will be used to provide “low-rank” approximations of a matrix. The idea is to keep only the k top singular values for some suitable $k \ll r$ for which $\sigma_{k+1}, \dots, \sigma_r$ are very small.

Remarks:

- (1) In Strang [105] the matrices U, V, D are denoted by $U = Q_2$, $V = Q_1$, and $D = \Sigma$, and an SVD is written as $A = Q_1 \Sigma Q_2^\top$. This has the advantage that Q_1 comes before Q_2 in $A = Q_1 \Sigma Q_2^\top$. This has the disadvantage that A maps the columns of Q_2 (eigenvectors of $A^\top A$) to multiples of the columns of Q_1 (eigenvectors of $A A^\top$).
- (2) Algorithms for actually computing the SVD of a matrix are presented in Golub and Van Loan [49], Demmel [27], and Trefethen and Bau [110], where the SVD and its applications are also discussed quite extensively.
- (3) The SVD also applies to complex matrices. In this case, for every complex $n \times n$ matrix A , there are two unitary matrices U and V and a diagonal matrix D such that

$$A = V D U^*,$$

where D is a diagonal matrix consisting of real entries $\sigma_1, \dots, \sigma_n$, where $\sigma_1, \dots, \sigma_r$ are the singular values of A , i.e., the positive square roots of the nonzero eigenvalues of $A^* A$ and $A A^*$, and $\sigma_{r+1} = \dots = \sigma_n = 0$.

A notion closely related to the SVD is the polar form of a matrix.

Definition 16.4. A pair (R, S) such that $A = RS$ with R orthogonal and S symmetric positive semidefinite is called a *polar decomposition* of A .

Theorem 16.2 implies that for every real $n \times n$ matrix A , there is some orthogonal matrix R and some positive semidefinite symmetric matrix S such that

$$A = RS.$$

This is easy to show and we will prove it below. Furthermore, R, S are unique if A is invertible, but this is harder to prove.

For example, the matrix

$$A = \frac{1}{2} \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & 1 & -1 \\ 1 & -1 & -1 & 1 \end{pmatrix}$$

is both orthogonal and symmetric, and $A = RS$ with $R = A$ and $S = I$, which implies that some of the eigenvalues of A are negative.

Remark: In the complex case, the polar decomposition states that for every complex $n \times n$ matrix A , there is some unitary matrix U and some positive semidefinite Hermitian matrix H such that

$$A = UH.$$

It is easy to go from the polar form to the SVD, and conversely.

Given an SVD decomposition $A = VDU^\top$, let $R = VU^\top$ and $S = UDU^\top$. It is clear that R is orthogonal and that S is positive semidefinite symmetric, and

$$RS = VU^\top UDU^\top = VDU^\top = A.$$

Going the other way, given a polar decomposition $A = R_1S$, where R_1 is orthogonal and S is positive semidefinite symmetric, there is an orthogonal matrix R_2 and a positive semidefinite diagonal matrix D such that $S = R_2DR_2^\top$, and thus

$$A = R_1R_2DR_2^\top = VDU^\top,$$

where $V = R_1R_2$ and $U = R_2$ are orthogonal.

The eigenvalues and the singular values of a matrix are typically not related in any obvious way. For example, the $n \times n$ matrix

$$A = \begin{pmatrix} 1 & 2 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & 2 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & 2 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 1 & 2 & 0 \\ 0 & 0 & \dots & 0 & 0 & 1 & 2 \\ 0 & 0 & \dots & 0 & 0 & 0 & 1 \end{pmatrix}$$

has the eigenvalue 1 with multiplicity n , but its singular values, $\sigma_1 \geq \dots \geq \sigma_n$, which are the positive square roots of the eigenvalues of the matrix $B = A^\top A$ with

$$B = \begin{pmatrix} 1 & 2 & 0 & 0 & \dots & 0 & 0 \\ 2 & 5 & 2 & 0 & \dots & 0 & 0 \\ 0 & 2 & 5 & 2 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 2 & 5 & 2 & 0 \\ 0 & 0 & \dots & 0 & 2 & 5 & 2 \\ 0 & 0 & \dots & 0 & 0 & 2 & 5 \end{pmatrix}$$

have a wide spread, since

$$\frac{\sigma_1}{\sigma_n} = \text{cond}_2(A) \geq 2^{n-1}.$$

If A is a complex $n \times n$ matrix, the eigenvalues $\lambda_1, \dots, \lambda_n$ and the singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$ of A are not unrelated, since

$$\sigma_1^2 \dots \sigma_n^2 = \det(A^*A) = |\det(A)|^2$$

and

$$|\lambda_1| \cdots |\lambda_n| = |\det(A)|,$$

so we have

$$|\lambda_1| \cdots |\lambda_n| = \sigma_1 \cdots \sigma_n.$$

More generally, Hermann Weyl proved the following remarkable theorem:

Theorem 16.3. (*Weyl's inequalities, 1949*) For any complex $n \times n$ matrix, A , if $\lambda_1, \dots, \lambda_n \in \mathbb{C}$ are the eigenvalues of A and $\sigma_1, \dots, \sigma_n \in \mathbb{R}_+$ are the singular values of A , listed so that $|\lambda_1| \geq \cdots \geq |\lambda_n|$ and $\sigma_1 \geq \cdots \geq \sigma_n \geq 0$, then

$$\begin{aligned} |\lambda_1| \cdots |\lambda_n| &= \sigma_1 \cdots \sigma_n \quad \text{and} \\ |\lambda_1| \cdots |\lambda_k| &\leq \sigma_1 \cdots \sigma_k, \quad \text{for } k = 1, \dots, n-1. \end{aligned}$$

A proof of Theorem 16.3 can be found in Horn and Johnson [58], Chapter 3, Section 3.3, where more inequalities relating the eigenvalues and the singular values of a matrix are given.

Theorem 16.2 can be easily extended to rectangular $m \times n$ matrices, as we show in the next section (for various versions of the SVD for rectangular matrices, see Strang [105] Golub and Van Loan [49], Demmel [27], and Trefethen and Bau [110]).

16.2 Singular Value Decomposition for Rectangular Matrices

Here is the generalization of Theorem 16.2 to rectangular matrices.

Theorem 16.4. (*Singular value decomposition*) For every real $m \times n$ matrix A , there are two orthogonal matrices U ($n \times n$) and V ($m \times m$) and a diagonal $m \times n$ matrix D such that $A = VDU^\top$, where D is of the form

$$D = \begin{pmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_n \\ 0 & & & & 0 \\ & & & & & \ddots \\ & & & & & & 0 \end{pmatrix} \quad \text{or} \quad D = \begin{pmatrix} \sigma_1 & & & 0 & \cdots & 0 \\ & \sigma_2 & & 0 & \cdots & 0 \\ & & \ddots & & & \\ & & & \sigma_m & & 0 \\ & & & & & 0 \end{pmatrix},$$

where $\sigma_1, \dots, \sigma_r$ are the singular values of f , i.e. the (positive) square roots of the nonzero eigenvalues of $A^\top A$ and AA^\top , and $\sigma_{r+1} = \dots = \sigma_p = 0$, where $p = \min(m, n)$. The columns of U are eigenvectors of $A^\top A$, and the columns of V are eigenvectors of AA^\top .

Proof. As in the proof of Theorem 16.2, since $A^\top A$ is symmetric positive semidefinite, there exists an $n \times n$ orthogonal matrix U such that

$$A^\top A = U \Sigma^2 U^\top,$$

with $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0)$, where $\sigma_1^2, \dots, \sigma_r^2$ are the nonzero eigenvalues of $A^\top A$, and where r is the rank of A . Observe that $r \leq \min\{m, n\}$, and AU is an $m \times n$ matrix. It follows that

$$U^\top A^\top A U = (AU)^\top A U = \Sigma^2,$$

and if we let $f_j \in \mathbb{R}^m$ be the j th column of AU for $j = 1, \dots, n$, then we have

$$\langle f_i, f_j \rangle = \sigma_i^2 \delta_{ij}, \quad 1 \leq i, j \leq r$$

and

$$f_j = 0, \quad r+1 \leq j \leq n.$$

If we define (v_1, \dots, v_r) by

$$v_j = \sigma_j^{-1} f_j, \quad 1 \leq j \leq r,$$

then we have

$$\langle v_i, v_j \rangle = \delta_{ij}, \quad 1 \leq i, j \leq r,$$

so complete (v_1, \dots, v_r) into an orthonormal basis $(v_1, \dots, v_r, v_{r+1}, \dots, v_m)$ (for example, using Gram-Schmidt).

Now, since $f_j = \sigma_j v_j$ for $j = 1, \dots, r$, we have

$$\langle v_i, f_j \rangle = \sigma_j \langle v_i, v_j \rangle = \sigma_j \delta_{ij}, \quad 1 \leq i \leq m, 1 \leq j \leq r$$

and since $f_j = 0$ for $j = r+1, \dots, n$, we have

$$\langle v_i, f_j \rangle = 0 \quad 1 \leq i \leq m, r+1 \leq j \leq n.$$

If V is the matrix whose columns are v_1, \dots, v_m , then V is an $m \times m$ orthogonal matrix and if $m \geq n$, we let

$$D = \begin{pmatrix} \Sigma \\ 0_{m-n} \end{pmatrix} = \begin{pmatrix} \sigma_1 & \dots & & \\ & \sigma_2 & \dots & \\ \vdots & \vdots & \ddots & \vdots \\ & & \dots & \sigma_n \\ 0 & \vdots & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \vdots & \dots & 0 \end{pmatrix},$$

else if $n \geq m$, then we let

$$D = \begin{pmatrix} \sigma_1 & \dots & 0 & \dots & 0 \\ & \sigma_2 & \dots & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & 0 & \vdots & 0 \\ & & \dots & \sigma_m & 0 & \dots & 0 \end{pmatrix}.$$

In either case, the above equations prove that

$$V^\top AU = D,$$

which yields $A = VDU^\top$, as required.

The equation $A = VDU^\top$ implies that

$$A^\top A = UD^\top DU^\top = U \operatorname{diag}(\sigma_1^2, \dots, \sigma_r^2, \underbrace{0, \dots, 0}_{n-r}) U^\top$$

and

$$AA^\top = VDD^\top V^\top = V \operatorname{diag}(\sigma_1^2, \dots, \sigma_r^2, \underbrace{0, \dots, 0}_{m-r}) V^\top,$$

which shows that $A^\top A$ and AA^\top have the same nonzero eigenvalues, that the columns of U are eigenvectors of $A^\top A$, and that the columns of V are eigenvectors of AA^\top . \square

A triple (U, D, V) such that $A = VDU^\top$ is called a *singular value decomposition (SVD)* of A .

Even though the matrix D is an $m \times n$ rectangular matrix, since its only nonzero entries are on the descending diagonal, we still say that D is a diagonal matrix.

If we view A as the representation of a linear map $f: E \rightarrow F$, where $\dim(E) = n$ and $\dim(F) = m$, the proof of Theorem 16.4 shows that there are two orthonormal bases (u_1, \dots, u_n) and (v_1, \dots, v_m) for E and F , respectively, where (u_1, \dots, u_n) are eigenvectors of $f^* \circ f$ and (v_1, \dots, v_m) are eigenvectors of $f \circ f^*$. Furthermore, (u_1, \dots, u_r) is an orthonormal basis of $\operatorname{Im} f^*$, (u_{r+1}, \dots, u_n) is an orthonormal basis of $\operatorname{Ker} f$, (v_1, \dots, v_r) is an orthonormal basis of $\operatorname{Im} f$, and (v_{r+1}, \dots, v_m) is an orthonormal basis of $\operatorname{Ker} f^*$.

The SVD of matrices can be used to define the pseudo-inverse of a rectangular matrix; we will do so in Chapter 17. The reader may also consult Strang [105], Demmel [27], Trefethen and Bau [110], and Golub and Van Loan [49].

One of the spectral theorems states that a symmetric matrix can be diagonalized by an orthogonal matrix. There are several numerical methods to compute the eigenvalues of a symmetric matrix A . One method consists in *tridiagonalizing* A , which means that there exists some orthogonal matrix P and some symmetric tridiagonal matrix T such that $A = PTP^\top$. In fact, this can be done using Householder transformations. It is then possible to compute the eigenvalues of T using a bisection method based on Sturm sequences. One can also use Jacobi's method. For details, see Golub and Van Loan [49], Chapter 8, Demmel [27], Trefethen and Bau [110], Lecture 26, or Ciarlet [24]. Computing the SVD of a matrix A is more involved. Most methods begin by finding orthogonal matrices U and V and a *bidagonal* matrix B such that $A = VBU^\top$. This can also be done using Householder transformations. Observe that $B^\top B$ is symmetric tridiagonal. Thus, in principle, the previous method to diagonalize a symmetric tridiagonal matrix can be applied. However, it is unwise to compute

$B^\top B$ explicitly, and more subtle methods are used for this last step. Again, see Golub and Van Loan [49], Chapter 8, Demmel [27], and Trefethen and Bau [110], Lecture 31.

The polar form has applications in continuum mechanics. Indeed, in any deformation it is important to separate stretching from rotation. This is exactly what QS achieves. The orthogonal part Q corresponds to rotation (perhaps with an additional reflection), and the symmetric matrix S to stretching (or compression). The real eigenvalues $\sigma_1, \dots, \sigma_r$ of S are the stretch factors (or compression factors) (see Marsden and Hughes [76]). The fact that S can be diagonalized by an orthogonal matrix corresponds to a natural choice of axes, the principal axes.

The SVD has applications to data compression, for instance in image processing. The idea is to retain only singular values whose magnitudes are significant enough. The SVD can also be used to determine the rank of a matrix when other methods such as Gaussian elimination produce very small pivots. One of the main applications of the SVD is the computation of the pseudo-inverse. Pseudo-inverses are the key to the solution of various optimization problems, in particular the method of least squares. This topic is discussed in the next chapter (Chapter 17). Applications of the material of this chapter can be found in Strang [105, 104]; Ciarlet [24]; Golub and Van Loan [49], which contains many other references; Demmel [27]; and Trefethen and Bau [110].

16.3 Ky Fan Norms and Schatten Norms

The singular values of a matrix can be used to define various norms on matrices which have found recent applications in quantum information theory and in spectral graph theory. Following Horn and Johnson [58] (Section 3.4) we can make the following definitions:

Definition 16.5. For any matrix $A \in M_{m,n}(\mathbb{C})$, let $q = \min\{m, n\}$, and if $\sigma_1 \geq \dots \geq \sigma_q$ are the singular values of A , for any k with $1 \leq k \leq q$, let

$$N_k(A) = \sigma_1 + \dots + \sigma_k,$$

called the *Ky Fan k -norm* of A .

More generally, for any $p \geq 1$ and any k with $1 \leq k \leq q$, let

$$N_{k;p}(A) = (\sigma_1^p + \dots + \sigma_k^p)^{1/p},$$

called the *Ky Fan p - k -norm* of A . When $k = q$, $N_{q;p}$ is also called the *Schatten p -norm*.

Observe that when $k = 1$, $N_1(A) = \sigma_1$, and the Ky Fan norm N_1 is simply the *spectral norm* from Chapter 7, which is the subordinate matrix norm associated with the Euclidean norm. When $k = q$, the Ky Fan norm N_q is given by

$$N_q(A) = \sigma_1 + \dots + \sigma_q = \text{tr}((A^*A)^{1/2})$$

and is called the *trace norm* or *nuclear norm*. When $p = 2$ and $k = q$, the Ky Fan $N_{q;2}$ norm is given by

$$N_{k;2}(A) = (\sigma_1^2 + \cdots + \sigma_q^2)^{1/2} = \sqrt{\operatorname{tr}(A^*A)} = \|A\|_F,$$

which is the *Frobenius norm* of A .

It can be shown that N_k and $N_{k;p}$ are unitarily invariant norms, and that when $m = n$, they are matrix norms; see Horn and Johnson [58] (Section 3.4, Corollary 3.4.4 and Problem 3).

16.4 Summary

The main concepts and results of this chapter are listed below:

- For any linear map $f: E \rightarrow E$ on a Euclidean space E , the maps $f^* \circ f$ and $f \circ f^*$ are self-adjoint and positive semidefinite.
- The *singular values* of a linear map.
- *Positive semidefinite* and *positive definite* self-adjoint maps.
- Relationships between $\operatorname{Im} f$, $\operatorname{Ker} f$, $\operatorname{Im} f^*$, and $\operatorname{Ker} f^*$.
- The *singular value decomposition theorem* for square matrices (Theorem 16.2).
- The *SVD* of matrix.
- The *polar decomposition* of a matrix.
- The *Weyl inequalities*.
- The *singular value decomposition theorem* for $m \times n$ matrices (Theorem 16.4).
- Ky Fan k -norms, Ky Fan p - k -norms, Schatten p -norms.

Chapter 17

Applications of SVD and Pseudo-Inverses

De tous les principes qu'on peut proposer pour cet objet, je pense qu'il n'en est pas de plus général, de plus exact, ni d'une application plus facile, que celui dont nous avons fait usage dans les recherches précédentes, et qui consiste à rendre *minimum* la somme des carrés des erreurs. Par ce moyen il s'établit entre les erreurs une sorte d'équilibre qui, empêchant les extrêmes de prévaloir, est très propre à faire connaître l'état du système le plus proche de la vérité.

—**Legendre, 1805**, *Nouvelles Méthodes pour la détermination des Orbites des Comètes*

17.1 Least Squares Problems and the Pseudo-Inverse

This chapter presents several applications of SVD. The first one is the pseudo-inverse, which plays a crucial role in solving linear systems by the method of least squares. The second application is data compression. The third application is principal component analysis (PCA), whose purpose is to identify patterns in data and understand the variance–covariance structure of the data. The fourth application is the best affine approximation of a set of data, a problem closely related to PCA.

The method of least squares is a way of “solving” an overdetermined system of linear equations

$$Ax = b,$$

i.e., a system in which A is a rectangular $m \times n$ matrix with more equations than unknowns (when $m > n$). Historically, the method of least squares was used by Gauss and Legendre to solve problems in astronomy and geodesy. The method was first published by Legendre in 1805 in a paper on methods for determining the orbits of comets. However, Gauss had already used the method of least squares as early as 1801 to determine the orbit of the asteroid

Ceres, and he published a paper about it in 1810 after the discovery of the asteroid Pallas. Incidentally, it is in that same paper that Gaussian elimination using pivots is introduced.

The reason why more equations than unknowns arise in such problems is that repeated measurements are taken to minimize errors. This produces an overdetermined and often inconsistent system of linear equations. For example, Gauss solved a system of eleven equations in six unknowns to determine the orbit of the asteroid Pallas. As a concrete illustration, suppose that we observe the motion of a small object, assimilated to a point, in the plane. From our observations, we suspect that this point moves along a straight line, say of equation $y = dx + c$. Suppose that we observed the moving point at three different locations (x_1, y_1) , (x_2, y_2) , and (x_3, y_3) . Then we should have

$$\begin{aligned}c + dx_1 &= y_1, \\c + dx_2 &= y_2, \\c + dx_3 &= y_3.\end{aligned}$$

If there were no errors in our measurements, these equations would be compatible, and c and d would be determined by only two of the equations. However, in the presence of errors, the system may be inconsistent. Yet we would like to find c and d !

The idea of the method of least squares is to determine (c, d) such that it minimizes the sum of the squares of the errors, namely,

$$(c + dx_1 - y_1)^2 + (c + dx_2 - y_2)^2 + (c + dx_3 - y_3)^2.$$

In general, for an overdetermined $m \times n$ system $Ax = b$, what Gauss and Legendre discovered is that there are solutions x minimizing

$$\|Ax - b\|_2^2$$

(where $\|u\|_2^2 = u_1^2 + \cdots + u_n^2$, the square of the Euclidean norm of the vector $u = (u_1, \dots, u_n)$), and that these solutions are given by the square $n \times n$ system

$$A^\top Ax = A^\top b,$$

called the *normal equations*. Furthermore, when the columns of A are linearly independent, it turns out that $A^\top A$ is invertible, and so x is unique and given by

$$x = (A^\top A)^{-1} A^\top b.$$

Note that $A^\top A$ is a symmetric matrix, one of the nice features of the normal equations of a least squares problem. For instance, the normal equations for the above problem are

$$\begin{pmatrix} 3 & x_1 + x_2 + x_3 \\ x_1 + x_2 + x_3 & x_1^2 + x_2^2 + x_3^2 \end{pmatrix} \begin{pmatrix} c \\ d \end{pmatrix} = \begin{pmatrix} y_1 + y_2 + y_3 \\ x_1 y_1 + x_2 y_2 + x_3 y_3 \end{pmatrix}.$$

In fact, given any real $m \times n$ matrix A , there is always a unique x^+ of minimum norm that minimizes $\|Ax - b\|_2^2$, even when the columns of A are linearly dependent. How do we prove this, and how do we find x^+ ?

Theorem 17.1. *Every linear system $Ax = b$, where A is an $m \times n$ matrix, has a unique least squares solution x^+ of smallest norm.*

Proof. Geometry offers a nice proof of the existence and uniqueness of x^+ . Indeed, we can interpret b as a point in the Euclidean (affine) space \mathbb{R}^m , and the image subspace of A (also called the column space of A) as a subspace U of \mathbb{R}^m (passing through the origin). Then, it is clear that

$$\inf_{x \in \mathbb{R}^n} \|Ax - b\|_2^2 = \inf_{y \in U} \|y - b\|_2^2,$$

with $U = \text{Im } A$, and we claim that x minimizes $\|Ax - b\|_2^2$ iff $Ax = p$, where p the orthogonal projection of b onto the subspace U .

Recall from Section 11.1 that the orthogonal projection $p_U: U \oplus U^\perp \rightarrow U$ is the linear map given by

$$p_U(u + v) = u,$$

with $u \in U$ and $v \in U^\perp$. If we let $p = p_U(b) \in U$, then for any point $y \in U$, the vectors $\overrightarrow{py} = y - p \in U$ and $\overrightarrow{bp} = p - b \in U^\perp$ are orthogonal, which implies that

$$\|\overrightarrow{by}\|_2^2 = \|\overrightarrow{bp}\|_2^2 + \|\overrightarrow{py}\|_2^2,$$

where $\overrightarrow{by} = y - b$. Thus, p is indeed the unique point in U that minimizes the distance from b to any point in U .

Thus, the problem has been reduced to proving that there is a unique x^+ of minimum norm such that $Ax^+ = p$, with $p = p_U(b) \in U$, the orthogonal projection of b onto U . We use the fact that

$$\mathbb{R}^n = \text{Ker } A \oplus (\text{Ker } A)^\perp.$$

Consequently, every $x \in \mathbb{R}^n$ can be written uniquely as $x = u + v$, where $u \in \text{Ker } A$ and $v \in (\text{Ker } A)^\perp$, and since u and v are orthogonal,

$$\|x\|_2^2 = \|u\|_2^2 + \|v\|_2^2.$$

Furthermore, since $u \in \text{Ker } A$, we have $Au = 0$, and thus $Ax = p$ iff $Av = p$, which shows that the solutions of $Ax = p$ for which x has minimum norm must belong to $(\text{Ker } A)^\perp$. However, the restriction of A to $(\text{Ker } A)^\perp$ is injective. This is because if $Av_1 = Av_2$, where $v_1, v_2 \in (\text{Ker } A)^\perp$, then $A(v_2 - v_1) = 0$, which implies $v_2 - v_1 \in \text{Ker } A$, and since $v_1, v_2 \in (\text{Ker } A)^\perp$, we also have $v_2 - v_1 \in (\text{Ker } A)^\perp$, and consequently, $v_2 - v_1 = 0$. This shows that there is a unique x^+ of minimum norm such that $Ax^+ = p$, and that x^+ must belong to $(\text{Ker } A)^\perp$. By our previous reasoning, x^+ is the unique vector of minimum norm minimizing $\|Ax - b\|_2^2$. \square

The proof also shows that x minimizes $\|Ax - b\|_2^2$ iff $\overrightarrow{pb} = b - Ax$ is orthogonal to U , which can be expressed by saying that $b - Ax$ is orthogonal to every column of A . However, this is equivalent to

$$A^\top(b - Ax) = 0, \quad \text{i.e.,} \quad A^\top Ax = A^\top b.$$

Finally, it turns out that the minimum norm least squares solution x^+ can be found in terms of the pseudo-inverse A^+ of A , which is itself obtained from any SVD of A .

Definition 17.1. Given any nonzero $m \times n$ matrix A of rank r , if $A = VDU^\top$ is an SVD of A such that

$$D = \begin{pmatrix} \Lambda & 0_{r,n-r} \\ 0_{m-r,r} & 0_{m-r,n-r} \end{pmatrix},$$

with

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_r)$$

an $r \times r$ diagonal matrix consisting of the nonzero singular values of A , then if we let D^+ be the $n \times m$ matrix

$$D^+ = \begin{pmatrix} \Lambda^{-1} & 0_{r,m-r} \\ 0_{n-r,r} & 0_{n-r,m-r} \end{pmatrix},$$

with

$$\Lambda^{-1} = \text{diag}(1/\lambda_1, \dots, 1/\lambda_r),$$

the *pseudo-inverse* of A is defined by

$$A^+ = UD^+V^\top.$$

If $A = 0_{m,n}$ is the zero matrix, we set $A^+ = 0_{n,m}$. Observe that D^+ is obtained from D by inverting the nonzero diagonal entries of D , leaving all zeros in place, and then transposing the matrix. The pseudo-inverse of a matrix is also known as the *Moore–Penrose pseudo-inverse*.

Actually, it seems that A^+ depends on the specific choice of U and V in an SVD (U, D, V) for A , but the next theorem shows that this is not so.

Theorem 17.2. *The least squares solution of smallest norm of the linear system $Ax = b$, where A is an $m \times n$ matrix, is given by*

$$x^+ = A^+b = UD^+V^\top b.$$

Proof. First, assume that A is a (rectangular) diagonal matrix D , as above. Then, since x minimizes $\|Dx - b\|_2^2$ iff Dx is the projection of b onto the image subspace F of D , it is fairly obvious that $x^+ = D^+b$. Otherwise, we can write

$$A = VDU^\top,$$

where U and V are orthogonal. However, since V is an isometry,

$$\|Ax - b\|_2 = \|VDU^\top x - b\|_2 = \|DU^\top x - V^\top b\|_2.$$

Letting $y = U^\top x$, we have $\|x\|_2 = \|y\|_2$, since U is an isometry, and since U is surjective, $\|Ax - b\|_2$ is minimized iff $\|Dy - V^\top b\|_2$ is minimized, and we have shown that the least solution is

$$y^+ = D^+ V^\top b.$$

Since $y = U^\top x$, with $\|x\|_2 = \|y\|_2$, we get

$$x^+ = U D^+ V^\top b = A^+ b.$$

Thus, the pseudo-inverse provides the optimal solution to the least squares problem. \square

By Proposition 17.2 and Theorem 17.1, $A^+ b$ is uniquely defined by every b , and thus A^+ depends only on A .

When A has full rank, the pseudo-inverse A^+ can be expressed as $A^+ = (A^\top A)^{-1} A^\top$ when $m \geq n$, and as $A^+ = A^\top (A A^\top)^{-1}$ when $n \geq m$. In the first case ($m \geq n$), observe that $A^+ A = I$, so A^+ is a left inverse of A ; in the second case ($n \geq m$), we have $A A^+ = I$, so A^+ is a right inverse of A .

Proof. If $m \geq n$ and A has full rank $\text{rank } n$, we have

$$A = V \begin{pmatrix} \Lambda \\ 0_{m-n,n} \end{pmatrix} U^\top$$

with Λ an $n \times n$ diagonal invertible matrix (with positive entries), so

$$A^+ = U \begin{pmatrix} \Lambda^{-1} & 0_{n,m-n} \end{pmatrix} V^\top.$$

We find that

$$A^\top A = U \begin{pmatrix} \Lambda & 0_{n,m-n} \end{pmatrix} V^\top V \begin{pmatrix} \Lambda \\ 0_{m-n,n} \end{pmatrix} U^\top = U \Lambda^2 U^\top,$$

which yields

$$(A^\top A)^{-1} A^\top = U \Lambda^{-2} U^\top U \begin{pmatrix} \Lambda & 0_{n,m-n} \end{pmatrix} V^\top V = U \begin{pmatrix} \Lambda^{-1} & 0_{n,m-n} \end{pmatrix} V^\top = A^+.$$

Therefore, if $m \geq n$ and A has full rank $\text{rank } n$, then

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

If $n \geq m$ and A has full rank $\text{rank } m$, then

$$A = V \begin{pmatrix} \Lambda & 0_{m,n-m} \end{pmatrix} U^\top$$

with Λ an $m \times m$ diagonal invertible matrix (with positive entries), so

$$A^+ = U \begin{pmatrix} \Lambda^{-1} \\ 0_{n-m,m} \end{pmatrix} V^\top.$$

We find that

$$AA^\top = V \begin{pmatrix} \Lambda & 0_{m,n-m} \end{pmatrix} U^\top U \begin{pmatrix} \Lambda \\ 0_{n-m,m} \end{pmatrix} V^\top = V \Lambda^2 V^\top,$$

which yields

$$A^\top (AA^\top)^{-1} = U \begin{pmatrix} \Lambda \\ 0_{n-m,m} \end{pmatrix} V^\top V \Lambda^{-2} V^\top = U \begin{pmatrix} \Lambda^{-1} \\ 0_{n-m,m} \end{pmatrix} V^\top = A^+.$$

Therefore, if $n \geq m$ and A has full rank $\text{rank } m$, then $A^+ = A^\top (AA^\top)^{-1}$. \square

Let $A = U\Sigma V^\top$ be an SVD for A . It is easy to check that

$$\begin{aligned} AA^+A &= A, \\ A^+AA^+ &= A^+, \end{aligned}$$

and both AA^+ and A^+A are symmetric matrices. In fact,

$$AA^+ = U\Sigma V^\top V\Sigma^+ U^\top = U\Sigma\Sigma^+ U^\top = U \begin{pmatrix} I_r & 0 \\ 0 & 0_{m-r} \end{pmatrix} U^\top$$

and

$$A^+A = V\Sigma^+ U^\top U\Sigma V^\top = V\Sigma^+\Sigma V^\top = V \begin{pmatrix} I_r & 0 \\ 0 & 0_{n-r} \end{pmatrix} V^\top.$$

We immediately get

$$\begin{aligned} (AA^+)^2 &= AA^+, \\ (A^+A)^2 &= A^+A, \end{aligned}$$

so both AA^+ and A^+A are orthogonal projections (since they are both symmetric). *We claim that AA^+ is the orthogonal projection onto the range of A and A^+A is the orthogonal projection onto $\text{Ker}(A)^\perp = \text{Im}(A^\top)$, the range of A^\top .*

Obviously, we have $\text{range}(AA^+) \subseteq \text{range}(A)$, and for any $y = Ax \in \text{range}(A)$, since $AA^+A = A$, we have

$$AA^+y = AA^+Ax = Ax = y,$$

so the image of AA^+ is indeed the range of A . It is also clear that $\text{Ker}(A) \subseteq \text{Ker}(A^+A)$, and since $AA^+A = A$, we also have $\text{Ker}(A^+A) \subseteq \text{Ker}(A)$, and so

$$\text{Ker}(A^+A) = \text{Ker}(A).$$

Since A^+A is Hermitian, $\text{range}(A^+A) = \text{range}((A^+A)^\top) = \text{Ker}(A^+A)^\perp = \text{Ker}(A)^\perp$, as claimed.

It will also be useful to see that $\text{range}(A) = \text{range}(AA^+)$ consists of all vectors $y \in \mathbb{R}^m$ such that

$$U^\top y = \begin{pmatrix} z \\ 0 \end{pmatrix},$$

with $z \in \mathbb{R}^r$.

Indeed, if $y = Ax$, then

$$U^\top y = U^\top Ax = U^\top U \Sigma V^\top x = \Sigma V^\top x = \begin{pmatrix} \Sigma_r & 0 \\ 0 & 0_{m-r} \end{pmatrix} V^\top x = \begin{pmatrix} z \\ 0 \end{pmatrix},$$

where Σ_r is the $r \times r$ diagonal matrix $\text{diag}(\sigma_1, \dots, \sigma_r)$. Conversely, if $U^\top y = \begin{pmatrix} z \\ 0 \end{pmatrix}$, then $y = U \begin{pmatrix} z \\ 0 \end{pmatrix}$, and

$$\begin{aligned} AA^+y &= U \begin{pmatrix} I_r & 0 \\ 0 & 0_{m-r} \end{pmatrix} U^\top y \\ &= U \begin{pmatrix} I_r & 0 \\ 0 & 0_{m-r} \end{pmatrix} U^\top U \begin{pmatrix} z \\ 0 \end{pmatrix} \\ &= U \begin{pmatrix} I_r & 0 \\ 0 & 0_{m-r} \end{pmatrix} \begin{pmatrix} z \\ 0 \end{pmatrix} \\ &= U \begin{pmatrix} z \\ 0 \end{pmatrix} = y, \end{aligned}$$

which shows that y belongs to the range of A .

Similarly, we claim that $\text{range}(A^+A) = \text{Ker}(A)^\perp$ consists of all vectors $y \in \mathbb{R}^n$ such that

$$V^\top y = \begin{pmatrix} z \\ 0 \end{pmatrix},$$

with $z \in \mathbb{R}^r$.

If $y = A^+Au$, then

$$y = A^+Au = V \begin{pmatrix} I_r & 0 \\ 0 & 0_{n-r} \end{pmatrix} V^\top u = V \begin{pmatrix} z \\ 0 \end{pmatrix},$$

for some $z \in \mathbb{R}^r$. Conversely, if $V^\top y = \begin{pmatrix} z \\ 0 \end{pmatrix}$, then $y = V \begin{pmatrix} z \\ 0 \end{pmatrix}$, and so

$$\begin{aligned} A^+AV \begin{pmatrix} z \\ 0 \end{pmatrix} &= V \begin{pmatrix} I_r & 0 \\ 0 & 0_{n-r} \end{pmatrix} V^\top V \begin{pmatrix} z \\ 0 \end{pmatrix} \\ &= V \begin{pmatrix} I_r & 0 \\ 0 & 0_{n-r} \end{pmatrix} \begin{pmatrix} z \\ 0 \end{pmatrix} \\ &= V \begin{pmatrix} z \\ 0 \end{pmatrix} = y, \end{aligned}$$

which shows that $y \in \text{range}(A^+A)$.

If A is a symmetric matrix, then in general, there is no SVD $U\Sigma V^\top$ of A with $U = V$. However, if A is positive semidefinite, then the eigenvalues of A are nonnegative, and so the nonzero eigenvalues of A are equal to the singular values of A and SVDs of A are of the form

$$A = U\Sigma U^\top.$$

Analogous results hold for complex matrices, but in this case, U and V are unitary matrices and AA^+ and A^+A are Hermitian orthogonal projections.

If A is a normal matrix, which means that $AA^\top = A^\top A$, then there is an intimate relationship between SVD's of A and block diagonalizations of A . As a consequence, the pseudo-inverse of a normal matrix A can be obtained directly from a block diagonalization of A .

If A is a (real) normal matrix, then we know from Theorem 13.16 that A can be block diagonalized with respect to an orthogonal matrix U as

$$A = U\Lambda U^\top,$$

where Λ is the (real) block diagonal matrix

$$\Lambda = \text{diag}(B_1, \dots, B_n),$$

consisting either of 2×2 blocks of the form

$$B_j = \begin{pmatrix} \lambda_j & -\mu_j \\ \mu_j & \lambda_j \end{pmatrix}$$

with $\mu_j \neq 0$, or of one-dimensional blocks $B_k = (\lambda_k)$. Then we have the following proposition:

Proposition 17.3. *For any (real) normal matrix A and any block diagonalization $A = U\Lambda U^\top$ of A as above, the pseudo-inverse of A is given by*

$$A^+ = U\Lambda^+ U^\top,$$

where Λ^+ is the pseudo-inverse of Λ . Furthermore, if

$$\Lambda = \begin{pmatrix} \Lambda_r & 0 \\ 0 & 0 \end{pmatrix},$$

where Λ_r has rank r , then

$$\Lambda^+ = \begin{pmatrix} \Lambda_r^{-1} & 0 \\ 0 & 0 \end{pmatrix}.$$

Proof. Assume that B_1, \dots, B_p are 2×2 blocks and that $\lambda_{2p+1}, \dots, \lambda_n$ are the scalar entries. We know that the numbers $\lambda_j \pm i\mu_j$, and the λ_{2p+k} are the eigenvalues of A . Let $\rho_{2j-1} =$

$\rho_{2j} = \sqrt{\lambda_j^2 + \mu_j^2}$ for $j = 1, \dots, p$, $\rho_{2p+j} = \lambda_j$ for $j = 1, \dots, n-2p$, and assume that the blocks are ordered so that $\rho_1 \geq \rho_2 \geq \dots \geq \rho_n$. Then it is easy to see that

$$UU^\top = U^\top U = U\Lambda U^\top U\Lambda^\top U^\top = U\Lambda\Lambda^\top U^\top,$$

with

$$\Lambda\Lambda^\top = \text{diag}(\rho_1^2, \dots, \rho_n^2),$$

so the singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$ of A , which are the nonnegative square roots of the eigenvalues of AA^\top , are such that

$$\sigma_j = \rho_j, \quad 1 \leq j \leq n.$$

We can define the diagonal matrices

$$\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0),$$

where $r = \text{rank}(A)$, $\sigma_1 \geq \dots \geq \sigma_r > 0$ and

$$\Theta = \text{diag}(\sigma_1^{-1}B_1, \dots, \sigma_{2p}^{-1}B_p, 1, \dots, 1),$$

so that Θ is an orthogonal matrix and

$$\Lambda = \Theta\Sigma = (B_1, \dots, B_p, \lambda_{2p+1}, \dots, \lambda_r, 0, \dots, 0).$$

But then we can write

$$A = U\Lambda U^\top = U\Theta\Sigma U^\top,$$

and we if let $V = U\Theta$, since U is orthogonal and Θ is also orthogonal, V is also orthogonal and $A = V\Sigma U^\top$ is an SVD for A . Now we get

$$A^+ = U\Sigma^+V^\top = U\Sigma^+\Theta^\top U^\top.$$

However, since Θ is an orthogonal matrix, $\Theta^\top = \Theta^{-1}$, and a simple calculation shows that

$$\Sigma^+\Theta^\top = \Sigma^+\Theta^{-1} = \Lambda^+,$$

which yields the formula

$$A^+ = U\Lambda^+U^\top.$$

Also observe that if we write

$$\Lambda_r = (B_1, \dots, B_p, \lambda_{2p+1}, \dots, \lambda_r),$$

then Λ_r is invertible and

$$\Lambda^+ = \begin{pmatrix} \Lambda_r^{-1} & 0 \\ 0 & 0 \end{pmatrix}.$$

Therefore, the pseudo-inverse of a normal matrix can be computed directly from any block diagonalization of A , as claimed. \square

The following properties, due to Penrose, characterize the pseudo-inverse of a matrix. We have already proved that the pseudo-inverse satisfies these equations. For a proof of the converse, see Kincaid and Cheney [63].

Proposition 17.4. *Given any $m \times n$ matrix A (real or complex), the pseudo-inverse A^+ of A is the unique $n \times m$ matrix satisfying the following properties:*

$$\begin{aligned} AA^+A &= A, \\ A^+AA^+ &= A^+, \\ (AA^+)^T &= AA^+, \\ (A^+A)^T &= A^+A. \end{aligned}$$

If A is an $m \times n$ matrix of rank n (and so $m \geq n$), it is immediately shown that the QR -decomposition in terms of Householder transformations applies as follows:

There are n $m \times m$ matrices H_1, \dots, H_n , Householder matrices or the identity, and an upper triangular $m \times n$ matrix R of rank n such that

$$A = H_1 \cdots H_n R.$$

Then, because each H_i is an isometry,

$$\|Ax - b\|_2 = \|Rx - H_n \cdots H_1 b\|_2,$$

and the least squares problem $Ax = b$ is equivalent to the system

$$Rx = H_n \cdots H_1 b.$$

Now, the system

$$Rx = H_n \cdots H_1 b$$

is of the form

$$\begin{pmatrix} R_1 \\ 0_{m-n} \end{pmatrix} x = \begin{pmatrix} c \\ d \end{pmatrix},$$

where R_1 is an invertible $n \times n$ matrix (since A has rank n), $c \in \mathbb{R}^n$, and $d \in \mathbb{R}^{m-n}$, and the least squares solution of smallest norm is

$$x^+ = R_1^{-1}c.$$

Since R_1 is a triangular matrix, it is very easy to invert R_1 .

The method of least squares is one of the most effective tools of the mathematical sciences. There are entire books devoted to it. Readers are advised to consult Strang [105], Golub and Van Loan [49], Demmel [27], and Trefethen and Bau [110], where extensions and applications of least squares (such as weighted least squares and recursive least squares) are described. Golub and Van Loan [49] also contains a very extensive bibliography, including a list of books on least squares.

17.2 Data Compression and SVD

Among the many applications of SVD, a very useful one is *data compression*, notably for images. In order to make precise the notion of closeness of matrices, we use the notion of *matrix norm*. This concept is defined in Chapter 7 and the reader may want to review it before reading any further.

Given an $m \times n$ matrix of rank r , we would like to find a best approximation of A by a matrix B of rank $k \leq r$ (actually, $k < r$) so that $\|A - B\|_2$ (or $\|A - B\|_F$) is minimized.

Proposition 17.5. *Let A be an $m \times n$ matrix of rank r and let $VDU^\top = A$ be an SVD for A . Write u_i for the columns of U , v_i for the columns of V , and $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p$ for the singular values of A ($p = \min(m, n)$). Then a matrix of rank $k < r$ closest to A (in the $\|\cdot\|_2$ norm) is given by*

$$A_k = \sum_{i=1}^k \sigma_i v_i u_i^\top = V \operatorname{diag}(\sigma_1, \dots, \sigma_k) U^\top$$

and $\|A - A_k\|_2 = \sigma_{k+1}$.

Proof. By construction, A_k has rank k , and we have

$$\|A - A_k\|_2 = \left\| \sum_{i=k+1}^p \sigma_i v_i u_i^\top \right\|_2 = \|V \operatorname{diag}(0, \dots, 0, \sigma_{k+1}, \dots, \sigma_p) U^\top\|_2 = \sigma_{k+1}.$$

It remains to show that $\|A - B\|_2 \geq \sigma_{k+1}$ for all rank- k matrices B . Let B be any rank- k matrix, so its kernel has dimension $n - k$. The subspace U_{k+1} spanned by (u_1, \dots, u_{k+1}) has dimension $k + 1$, and because the sum of the dimensions of the kernel of B and of U_{k+1} is $(n - k) + k + 1 = n + 1$, these two subspaces must intersect in a subspace of dimension at least 1. Pick any unit vector h in $\operatorname{Ker}(B) \cap U_{k+1}$. Then since $Bh = 0$, we have

$$\|A - B\|_2^2 \geq \|(A - B)h\|_2^2 = \|Ah\|_2^2 = \|VDU^\top h\|_2^2 = \|DU^\top h\|_2^2 \geq \sigma_{k+1}^2 \|U^\top h\|_2^2 = \sigma_{k+1}^2,$$

which proves our claim. \square

Note that A_k can be stored using $(m + n)k$ entries, as opposed to mn entries. When $k \ll m$, this is a substantial gain.

A nice example of the use of Proposition 17.5 in image compression is given in Demmel [27], Chapter 3, Section 3.2.3, pages 113–115; see the Matlab demo.

An interesting topic that we have not addressed is the actual computation of an SVD. This is a very interesting but tricky subject. Most methods reduce the computation of an SVD to the diagonalization of a well-chosen symmetric matrix (which is not $A^\top A$). Interested readers should read Section 5.4 of Demmel's excellent book [27], which contains an overview of most known methods and an extensive list of references.

17.3 Principal Components Analysis (PCA)

Suppose we have a set of data consisting of n points X_1, \dots, X_n , with each $X_i \in \mathbb{R}^d$ viewed as a row vector.

Think of the X_i 's as persons, and if $X_i = (x_{i1}, \dots, x_{id})$, each x_{ij} is the value of some *feature* (or *attribute*) of that person. For example, the X_i 's could be mathematicians, $d = 2$, and the first component, x_{i1} , of X_i could be the year that X_i was born, and the second component, x_{i2} , the length of the beard of X_i in centimeters. Here is a small data set:

Name	year	length
Carl Friedrich Gauss	1777	0
Camille Jordan	1838	12
Adrien-Marie Legendre	1752	0
Bernhard Riemann	1826	15
David Hilbert	1862	2
Henri Poincaré	1854	5
Emmy Noether	1882	0
Karl Weierstrass	1815	0
Eugenio Beltrami	1835	2
Hermann Schwarz	1843	20

We usually form the $n \times d$ matrix X whose i th row is X_i , with $1 \leq i \leq n$. Then the j th column is denoted by C_j ($1 \leq j \leq d$). It is sometimes called a *feature vector*, but this terminology is far from being universally accepted. In fact, many people in computer vision call the data points X_i feature vectors!

The purpose of *principal components analysis*, for short *PCA*, is to identify patterns in data and understand the *variance-covariance* structure of the data. This is useful for the following tasks:

1. Data reduction: Often much of the variability of the data can be accounted for by a smaller number of *principal components*.
2. Interpretation: PCA can show relationships that were not previously suspected.

Given a vector (a *sample* of measurements) $x = (x_1, \dots, x_n) \in \mathbb{R}^n$, recall that the *mean* (or *average*) \bar{x} of x is given by

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}.$$

We let $x - \bar{x}$ denote the *centered data point*

$$x - \bar{x} = (x_1 - \bar{x}, \dots, x_n - \bar{x}).$$

In order to *measure the spread* of the x_i 's around the mean, we define the *sample variance* (for short, *variance*) $\text{var}(x)$ (or s^2) of the sample x by

$$\text{var}(x) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}.$$

There is a reason for using $n-1$ instead of n . The above definition makes $\text{var}(x)$ an unbiased estimator of the variance of the random variable being sampled. However, we don't need to worry about this. Curious readers will find an explanation of these peculiar definitions in Epstein [35] (Chapter 14, Section 14.5), or in any decent statistics book.

Given two vectors $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_n)$, the *sample covariance* (for short, *covariance*) of x and y is given by

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}.$$

The covariance of x and y measures how x and y vary from the mean with respect to each other. Obviously, $\text{cov}(x, y) = \text{cov}(y, x)$ and $\text{cov}(x, x) = \text{var}(x)$.

Note that

$$\text{cov}(x, y) = \frac{(x - \bar{x})^\top (y - \bar{y})}{n-1}.$$

We say that x and y are *uncorrelated* iff $\text{cov}(x, y) = 0$.

Finally, given an $n \times d$ matrix X of n points X_i , for PCA to be meaningful, it will be necessary to translate the origin to the *centroid* (or *center of gravity*) μ of the X_i 's, defined by

$$\mu = \frac{1}{n}(X_1 + \dots + X_n).$$

Observe that if $\mu = (\mu_1, \dots, \mu_d)$, then μ_j is the mean of the vector C_j (the j th column of X).

We let $X - \mu$ denote the *matrix* whose i th row is the centered data point $X_i - \mu$ ($1 \leq i \leq n$). Then, the *sample covariance matrix* (for short, *covariance matrix*) of X is the $d \times d$ symmetric matrix

$$\Sigma = \frac{1}{n-1}(X - \mu)^\top (X - \mu) = (\text{cov}(C_i, C_j)).$$

Remark: The factor $\frac{1}{n-1}$ is irrelevant for our purposes and can be ignored.

Here is the matrix $X - \mu$ in the case of our bearded mathematicians: Since

$$\mu_1 = 1828.4, \quad \mu_2 = 5.6,$$

we get

Name	year	length
Carl Friedrich Gauss	−51.4	−5.6
Camille Jordan	9.6	6.4
Adrien-Marie Legendre	−76.4	−5.6
Bernhard Riemann	−2.4	9.4
David Hilbert	33.6	−3.6
Henri Poincaré	25.6	−0.6
Emmy Noether	53.6	−5.6
Karl Weierstrass	13.4	−5.6
Eugenio Beltrami	6.6	−3.6
Hermann Schwarz	14.6	14.4

We can think of the vector C_j as representing the features of X in the direction e_j (the j th canonical basis vector in \mathbb{R}^d , namely $e_j = (0, \dots, 1, \dots, 0)$, with a 1 in the j th position).

If $v \in \mathbb{R}^d$ is a unit vector, we wish to consider the projection of the data points X_1, \dots, X_n onto the line spanned by v . Recall from Euclidean geometry that if $x \in \mathbb{R}^d$ is any vector and $v \in \mathbb{R}^d$ is a unit vector, the projection of x onto the line spanned by v is

$$\langle x, v \rangle v.$$

Thus, with respect to the basis v , the projection of x has coordinate $\langle x, v \rangle$. If x is represented by a row vector and v by a column vector, then

$$\langle x, v \rangle = xv.$$

Therefore, the vector $Y \in \mathbb{R}^n$ consisting of the coordinates of the projections of X_1, \dots, X_n onto the line spanned by v is given by $Y = Xv$, and this is the linear combination

$$Xv = v_1 C_1 + \dots + v_d C_d$$

of the columns of X (with $v = (v_1, \dots, v_d)$).

Observe that because μ_j is the mean of the vector C_j (the j th column of X), we get

$$\bar{Y} = \overline{Xv} = v_1 \mu_1 + \dots + v_d \mu_d,$$

and so the centered point $Y - \bar{Y}$ is given by

$$Y - \bar{Y} = v_1 (C_1 - \mu_1) + \dots + v_d (C_d - \mu_d) = (X - \mu)v.$$

Furthermore, if $Y = Xv$ and $Z = Xw$, then

$$\begin{aligned} \text{cov}(Y, Z) &= \frac{((X - \mu)v)^\top (X - \mu)w}{n - 1} \\ &= v^\top \frac{1}{n - 1} (X - \mu)^\top (X - \mu)w \\ &= v^\top \Sigma w, \end{aligned}$$

where Σ is the covariance matrix of X . Since $Y - \bar{Y}$ has zero mean, we have

$$\text{var}(Y) = \text{var}(Y - \bar{Y}) = v^\top \frac{1}{n-1} (X - \mu)^\top (X - \mu) v.$$

The above suggests that we should move the origin to the centroid μ of the X_i 's and consider the matrix $X - \mu$ of the centered data points $X_i - \mu$.

From now on, beware that we denote the columns of $X - \mu$ by C_1, \dots, C_d and that Y denotes the *centered* point $Y = (X - \mu)v = \sum_{j=1}^d v_j C_j$, where v is a unit vector.

Basic idea of PCA: The principal components of X are *uncorrelated* projections Y of the data points X_1, \dots, X_n onto some directions v (where the v 's are unit vectors) such that $\text{var}(Y)$ is maximal.

This suggests the following definition:

Definition 17.2. Given an $n \times d$ matrix X of data points X_1, \dots, X_n , if μ is the centroid of the X_i 's, then a *first principal component of X (first PC)* is a centered point $Y_1 = (X - \mu)v_1$, the projection of X_1, \dots, X_n onto a direction v_1 such that $\text{var}(Y_1)$ is maximized, where v_1 is a unit vector (recall that $Y_1 = (X - \mu)v_1$ is a linear combination of the C_j 's, the columns of $X - \mu$).

More generally, if Y_1, \dots, Y_k are k principal components of X along some unit vectors v_1, \dots, v_k , where $1 \leq k < d$, a *$(k+1)$ th principal component of X ($(k+1)$ th PC)* is a centered point $Y_{k+1} = (X - \mu)v_{k+1}$, the projection of X_1, \dots, X_n onto some direction v_{k+1} such that $\text{var}(Y_{k+1})$ is maximized, subject to $\text{cov}(Y_h, Y_{k+1}) = 0$ for all h with $1 \leq h \leq k$, and where v_{k+1} is a unit vector (recall that $Y_h = (X - \mu)v_h$ is a linear combination of the C_j 's). The v_h are called *principal directions*.

The following proposition is the key to the main result about PCA:

Proposition 17.6. If A is a symmetric $d \times d$ matrix with eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ and if (u_1, \dots, u_d) is any orthonormal basis of eigenvectors of A , where u_i is a unit eigenvector associated with λ_i , then

$$\max_{x \neq 0} \frac{x^\top A x}{x^\top x} = \lambda_1$$

(with the maximum attained for $x = u_1$) and

$$\max_{x \neq 0, x \in \{u_1, \dots, u_k\}^\perp} \frac{x^\top A x}{x^\top x} = \lambda_{k+1}$$

(with the maximum attained for $x = u_{k+1}$), where $1 \leq k \leq d-1$.

Proof. First, observe that

$$\max_{x \neq 0} \frac{x^\top Ax}{x^\top x} = \max_x \{x^\top Ax \mid x^\top x = 1\},$$

and similarly,

$$\max_{x \neq 0, x \in \{u_1, \dots, u_k\}^\perp} \frac{x^\top Ax}{x^\top x} = \max_x \{x^\top Ax \mid (x \in \{u_1, \dots, u_k\}^\perp) \wedge (x^\top x = 1)\}.$$

Since A is a symmetric matrix, its eigenvalues are real and it can be diagonalized with respect to an orthonormal basis of eigenvectors, so let (u_1, \dots, u_d) be such a basis. If we write

$$x = \sum_{i=1}^d x_i u_i,$$

a simple computation shows that

$$x^\top Ax = \sum_{i=1}^d \lambda_i x_i^2.$$

If $x^\top x = 1$, then $\sum_{i=1}^d x_i^2 = 1$, and since we assumed that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$, we get

$$x^\top Ax = \sum_{i=1}^d \lambda_i x_i^2 \leq \lambda_1 \left(\sum_{i=1}^d x_i^2 \right) = \lambda_1.$$

Thus,

$$\max_x \{x^\top Ax \mid x^\top x = 1\} \leq \lambda_1,$$

and since this maximum is achieved for $e_1 = (1, 0, \dots, 0)$, we conclude that

$$\max_x \{x^\top Ax \mid x^\top x = 1\} = \lambda_1.$$

Next, observe that $x \in \{u_1, \dots, u_k\}^\perp$ and $x^\top x = 1$ iff $x_1 = \dots = x_k = 0$ and $\sum_{i=1}^d x_i^2 = 1$. Consequently, for such an x , we have

$$x^\top Ax = \sum_{i=k+1}^d \lambda_i x_i^2 \leq \lambda_{k+1} \left(\sum_{i=k+1}^d x_i^2 \right) = \lambda_{k+1}.$$

Thus,

$$\max_x \{x^\top Ax \mid (x \in \{u_1, \dots, u_k\}^\perp) \wedge (x^\top x = 1)\} \leq \lambda_{k+1},$$

and since this maximum is achieved for $e_{k+1} = (0, \dots, 0, 1, 0, \dots, 0)$ with a 1 in position $k+1$, we conclude that

$$\max_x \{x^\top Ax \mid (x \in \{u_1, \dots, u_k\}^\perp) \wedge (x^\top x = 1)\} = \lambda_{k+1},$$

as claimed. □

The quantity

$$\frac{x^\top Ax}{x^\top x}$$

is known as the *Rayleigh–Ritz ratio* and Proposition 17.6 is often known as part of the *Rayleigh–Ritz theorem*.

Proposition 17.6 also holds if A is a Hermitian matrix and if we replace $x^\top Ax$ by x^*Ax and $x^\top x$ by x^*x . The proof is unchanged, since a Hermitian matrix has real eigenvalues and is diagonalized with respect to an orthonormal basis of eigenvectors (with respect to the Hermitian inner product).

We then have the following fundamental result showing how *the SVD of X yields the PCs*:

Theorem 17.7. (*SVD yields PCA*) *Let X be an $n \times d$ matrix of data points X_1, \dots, X_n , and let μ be the centroid of the X_i 's. If $X - \mu = VDU^\top$ is an SVD decomposition of $X - \mu$ and if the main diagonal of D consists of the singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d$, then the centered points Y_1, \dots, Y_d , where*

$$Y_k = (X - \mu)u_k = \text{\textit{kth column of } } VD$$

and u_k is the k th column of U , are d principal components of X . Furthermore,

$$\text{var}(Y_k) = \frac{\sigma_k^2}{n-1}$$

and $\text{cov}(Y_h, Y_k) = 0$, whenever $h \neq k$ and $1 \leq k, h \leq d$.

Proof. Recall that for any unit vector v , the centered projection of the points X_1, \dots, X_n onto the line of direction v is $Y = (X - \mu)v$ and that the variance of Y is given by

$$\text{var}(Y) = v^\top \frac{1}{n-1} (X - \mu)^\top (X - \mu) v.$$

Since $X - \mu = VDU^\top$, we get

$$\begin{aligned} \text{var}(Y) &= v^\top \frac{1}{(n-1)} (X - \mu)^\top (X - \mu) v \\ &= v^\top \frac{1}{(n-1)} U D V^\top V D U^\top v \\ &= v^\top U \frac{1}{(n-1)} D^2 U^\top v. \end{aligned}$$

Similarly, if $Y = (X - \mu)v$ and $Z = (X - \mu)w$, then the covariance of Y and Z is given by

$$\text{cov}(Y, Z) = v^\top U \frac{1}{(n-1)} D^2 U^\top w.$$

Obviously, $U \frac{1}{(n-1)} D^2 U^\top$ is a symmetric matrix whose eigenvalues are $\frac{\sigma_1^2}{n-1} \geq \dots \geq \frac{\sigma_d^2}{n-1}$, and the columns of U form an orthonormal basis of unit eigenvectors.

We proceed by induction on k . For the base case, $k = 1$, maximizing $\text{var}(Y)$ is equivalent to maximizing

$$v^\top U \frac{1}{(n-1)} D^2 U^\top v,$$

where v is a unit vector. By Proposition 17.6, the maximum of the above quantity is the largest eigenvalue of $U \frac{1}{(n-1)} D^2 U^\top$, namely $\frac{\sigma_1^2}{n-1}$, and it is achieved for u_1 , the first column of U . Now we get

$$Y_1 = (X - \mu)u_1 = V D U^\top u_1,$$

and since the columns of U form an orthonormal basis, $U^\top u_1 = e_1 = (1, 0, \dots, 0)$, and so Y_1 is indeed the first column of VD .

By the induction hypothesis, the centered points Y_1, \dots, Y_k , where $Y_h = (X - \mu)u_h$ and u_1, \dots, u_k are the first k columns of U , are k principal components of X . Because

$$\text{cov}(Y, Z) = v^\top U \frac{1}{(n-1)} D^2 U^\top w,$$

where $Y = (X - \mu)v$ and $Z = (X - \mu)w$, the condition $\text{cov}(Y_h, Z) = 0$ for $h = 1, \dots, k$ is equivalent to the fact that w belongs to the orthogonal complement of the subspace spanned by $\{u_1, \dots, u_k\}$, and maximizing $\text{var}(Z)$ subject to $\text{cov}(Y_h, Z) = 0$ for $h = 1, \dots, k$ is equivalent to maximizing

$$w^\top U \frac{1}{(n-1)} D^2 U^\top w,$$

where w is a unit vector orthogonal to the subspace spanned by $\{u_1, \dots, u_k\}$. By Proposition 17.6, the maximum of the above quantity is the $(k+1)$ th eigenvalue of $U \frac{1}{(n-1)} D^2 U^\top$, namely $\frac{\sigma_{k+1}^2}{n-1}$, and it is achieved for u_{k+1} , the $(k+1)$ th column of U . Now we get

$$Y_{k+1} = (X - \mu)u_{k+1} = V D U^\top u_{k+1},$$

and since the columns of U form an orthonormal basis, $U^\top u_{k+1} = e_{k+1}$, and Y_{k+1} is indeed the $(k+1)$ th column of VD , which completes the proof of the induction step. \square

The d columns u_1, \dots, u_d of U are usually called the *principal directions* of $X - \mu$ (and X). We note that not only do we have $\text{cov}(Y_h, Y_k) = 0$ whenever $h \neq k$, but the directions u_1, \dots, u_d along which the data are projected are mutually orthogonal.

We know from our study of SVD that $\sigma_1^2, \dots, \sigma_d^2$ are the eigenvalues of the symmetric positive semidefinite matrix $(X - \mu)^\top (X - \mu)$ and that u_1, \dots, u_d are corresponding eigenvectors. Numerically, it is preferable to use SVD on $X - \mu$ rather than to compute explicitly $(X - \mu)^\top (X - \mu)$ and then diagonalize it. Indeed, the explicit computation of $A^\top A$ from

a matrix A can be numerically quite unstable, and good SVD algorithms avoid computing $A^\top A$ explicitly.

In general, since an SVD of X is not unique, *the principal directions u_1, \dots, u_d are not unique*. This can happen when a data set has some *rotational symmetries*, and in such a case, PCA is not a very good method for analyzing the data set.

17.4 Best Affine Approximation

A problem very close to PCA (and based on least squares) is to *best approximate a data set of n points X_1, \dots, X_n , with $X_i \in \mathbb{R}^d$, by a p -dimensional affine subspace A of \mathbb{R}^d , with $1 \leq p \leq d-1$ (the terminology rank $d-p$ is also used).*

First, consider $p = d-1$. Then $A = A_1$ is an affine hyperplane (in \mathbb{R}^d), and it is given by an equation of the form

$$a_1 x_1 + \dots + a_d x_d + c = 0.$$

By *best approximation*, we mean that (a_1, \dots, a_d, c) solves the homogeneous linear system

$$\begin{pmatrix} x_{11} & \cdots & x_{1d} & 1 \\ \vdots & \vdots & \vdots & \vdots \\ x_{n1} & \cdots & x_{nd} & 1 \end{pmatrix} \begin{pmatrix} a_1 \\ \vdots \\ a_d \\ c \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix}$$

in the *least squares sense*, subject to the condition that $a = (a_1, \dots, a_d)$ is a unit vector, that is, $a^\top a = 1$, where $X_i = (x_{i1}, \dots, x_{id})$.

If we form the symmetric matrix

$$\begin{pmatrix} x_{11} & \cdots & x_{1d} & 1 \\ \vdots & \vdots & \vdots & \vdots \\ x_{n1} & \cdots & x_{nd} & 1 \end{pmatrix}^\top \begin{pmatrix} x_{11} & \cdots & x_{1d} & 1 \\ \vdots & \vdots & \vdots & \vdots \\ x_{n1} & \cdots & x_{nd} & 1 \end{pmatrix}$$

involved in the normal equations, we see that the bottom row (and last column) of that matrix is

$$n\mu_1 \quad \cdots \quad n\mu_d \quad n,$$

where $n\mu_j = \sum_{i=1}^n x_{ij}$ is n times the mean of the column C_j of X .

Therefore, if (a_1, \dots, a_d, c) is a least squares solution, that is, a solution of the normal equations, we must have

$$n\mu_1 a_1 + \cdots + n\mu_d a_d + nc = 0,$$

that is,

$$a_1 \mu_1 + \cdots + a_d \mu_d + c = 0,$$

which means that the *hyperplane* A_1 must pass through the centroid μ of the data points X_1, \dots, X_n . Then we can rewrite the original system with respect to the centered data $X_i - \mu$, and we find that the variable c drops out and we get the system

$$(X - \mu)a = 0,$$

where $a = (a_1, \dots, a_d)$.

Thus, we are looking for a unit vector a solving $(X - \mu)a = 0$ in the least squares sense, that is, some a such that $a^\top a = 1$ minimizing

$$a^\top (X - \mu)^\top (X - \mu) a.$$

Compute some SVD VDU^\top of $X - \mu$, where the main diagonal of D consists of the singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d$ of $X - \mu$ arranged in descending order. Then

$$a^\top (X - \mu)^\top (X - \mu) a = a^\top U D^2 U^\top a,$$

where $D^2 = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$ is a diagonal matrix, so pick a to be *the last column in* U (corresponding to the smallest eigenvalue σ_d^2 of $(X - \mu)^\top (X - \mu)$). This is a solution to our best fit problem.

Therefore, if U_{d-1} is the linear hyperplane defined by a , that is,

$$U_{d-1} = \{u \in \mathbb{R}^d \mid \langle u, a \rangle = 0\},$$

where a is the last column in U for some SVD VDU^\top of $X - \mu$, we have shown that the affine hyperplane $A_1 = \mu + U_{d-1}$ is a best approximation of the data set X_1, \dots, X_n in the least squares sense.

It is easy to show that this hyperplane $A_1 = \mu + U_{d-1}$ minimizes the sum of the square distances of each X_i to its orthogonal projection onto A_1 . Also, since U_{d-1} is the orthogonal complement of a , the last column of U , we see that U_{d-1} is spanned by the first $d-1$ columns of U , that is, the first $d-1$ principal directions of $X - \mu$.

All this can be generalized to a *best* $(d-k)$ -dimensional affine subspace A_k approximating X_1, \dots, X_n in the least squares sense ($1 \leq k \leq d-1$). Such an affine subspace A_k is cut out by k independent hyperplanes H_i (with $1 \leq i \leq k$), each given by some equation

$$a_{i1}x_1 + \dots + a_{id}x_d + c_i = 0.$$

If we write $a_i = (a_{i1}, \dots, a_{id})$, to say that the H_i are independent means that a_1, \dots, a_k are linearly independent. In fact, we may assume that a_1, \dots, a_k form an *orthonormal system*.

Then, finding a best $(d-k)$ -dimensional affine subspace A_k amounts to solving the homogeneous linear system

$$\begin{pmatrix} X & \mathbf{1} & 0 & \cdots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 0 & X & \mathbf{1} \end{pmatrix} \begin{pmatrix} a_1 \\ c_1 \\ \vdots \\ a_k \\ c_k \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix},$$

in the least squares sense, subject to the conditions $a_i^\top a_j = \delta_{ij}$, for all i, j with $1 \leq i, j \leq k$, where the matrix of the system is a block diagonal matrix consisting of k diagonal blocks $(X, \mathbf{1})$, where $\mathbf{1}$ denotes the column vector $(1, \dots, 1) \in \mathbb{R}^n$.

Again, it is easy to see that each hyperplane H_i must pass through the centroid μ of X_1, \dots, X_n , and by switching to the centered data $X_i - \mu$ we get the system

$$\begin{pmatrix} X - \mu & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & X - \mu \end{pmatrix} \begin{pmatrix} a_1 \\ \vdots \\ a_k \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix},$$

with $a_i^\top a_j = \delta_{ij}$ for all i, j with $1 \leq i, j \leq k$.

If $VDU^\top = X - \mu$ is an SVD decomposition, it is easy to see that a least squares solution of this system is given by the last k columns of U , assuming that the main diagonal of D consists of the singular values $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_d$ of $X - \mu$ arranged in descending order. But now the $(d - k)$ -dimensional subspace U_{d-k} cut out by the hyperplanes defined by a_1, \dots, a_k is simply the orthogonal complement of (a_1, \dots, a_k) , which is the subspace spanned by the first $d - k$ columns of U .

So the best $(d - k)$ -dimensional affine subspace A_k approximating X_1, \dots, X_n in the least squares sense is

$$A_k = \mu + U_{d-k},$$

where U_{d-k} is the linear subspace spanned by the first $d - k$ principal directions of $X - \mu$, that is, the first $d - k$ columns of U . Consequently, we get the following interesting interpretation of PCA (actually, principal directions):

Theorem 17.8. *Let X be an $n \times d$ matrix of data points X_1, \dots, X_n , and let μ be the centroid of the X_i 's. If $X - \mu = VDU^\top$ is an SVD decomposition of $X - \mu$ and if the main diagonal of D consists of the singular values $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_d$, then a best $(d - k)$ -dimensional affine approximation A_k of X_1, \dots, X_n in the least squares sense is given by*

$$A_k = \mu + U_{d-k},$$

where U_{d-k} is the linear subspace spanned by the first $d - k$ columns of U , the first $d - k$ principal directions of $X - \mu$ ($1 \leq k \leq d - 1$).

There are many applications of PCA to data compression, dimension reduction, and pattern analysis. The basic idea is that in many cases, given a data set X_1, \dots, X_n , with $X_i \in \mathbb{R}^d$, only a “small” subset of $m < d$ of the features is needed to describe the data set accurately.

If u_1, \dots, u_d are the principal directions of $X - \mu$, then the first m projections of the data (the first m principal components, i.e., the first m columns of VD) onto the first m principal directions represent the data without much loss of information. Thus, instead of using the

original data points X_1, \dots, X_n , with $X_i \in \mathbb{R}^d$, we can use their projections onto the first m principal directions Y_1, \dots, Y_m , where $Y_i \in \mathbb{R}^m$ and $m < d$, obtaining a compressed version of the original data set.

For example, PCA is used in computer vision for *face recognition*. Sirovitch and Kirby (1987) seem to be the first to have had the idea of using PCA to compress facial images. They introduced the term *eigenpicture* to refer to the principal directions, u_i . However, an explicit face recognition algorithm was given only later, by Turk and Pentland (1991). They renamed eigenpictures as *eigenfaces*.

For details on the topic of eigenfaces, see Forsyth and Ponce [39] (Chapter 22, Section 22.3.2), where you will also find exact references to Turk and Pentland's papers.

Another interesting application of PCA is to the *recognition of handwritten digits*. Such an application is described in Hastie, Tibshirani, and Friedman, [55] (Chapter 14, Section 14.5.1).

17.5 Summary

The main concepts and results of this chapter are listed below:

- *Least squares problems.*
- Existence of a least squares solution of smallest norm (Theorem 17.1).
- The *pseudo-inverse* A^+ of a matrix A .
- The least squares solution of smallest norm is given by the pseudo-inverse (Theorem 17.2)
- Projection properties of the pseudo-inverse.
- The pseudo-inverse of a normal matrix.
- The *Penrose characterization* of the pseudo-inverse.
- Data compression and SVD.
- Best approximation of rank $< r$ of a matrix.
- *Principal component analysis.*
- Review of basic statistical concepts: *mean, variance, covariance, covariance matrix.*
- Centered data, *centroid*.
- The *principal components (PCA)*.

- The *Rayleigh–Ritz theorem* (Theorem 17.6).
- The main theorem: *SVD yields PCA* (Theorem 17.7).
- Best affine approximation.
- SVD yields a best affine approximation (Theorem 17.8).
- Face recognition, eigenfaces.

Chapter 18

Quadratic Optimization Problems

18.1 Quadratic Optimization: The Positive Definite Case

In this chapter, we consider two classes of quadratic optimization problems that appear frequently in engineering and in computer science (especially in computer vision):

1. Minimizing

$$f(x) = \frac{1}{2}x^\top Ax + x^\top b$$

over all $x \in \mathbb{R}^n$, or subject to linear or affine constraints.

2. Minimizing

$$f(x) = \frac{1}{2}x^\top Ax + x^\top b$$

over the unit sphere.

In both cases, A is a symmetric matrix. We also seek necessary and sufficient conditions for f to have a global minimum.

Many problems in physics and engineering can be stated as the minimization of some energy function, with or without constraints. Indeed, it is a fundamental principle of mechanics that nature acts so as to minimize energy. Furthermore, if a physical system is in a stable state of equilibrium, then the energy in that state should be minimal. For example, a small ball placed on top of a sphere is in an unstable equilibrium position. A small motion causes the ball to roll down. On the other hand, a ball placed inside and at the bottom of a sphere is in a stable equilibrium position, because the potential energy is minimal.

The simplest kind of energy function is a quadratic function. Such functions can be conveniently defined in the form

$$P(x) = x^\top Ax - x^\top b,$$

where A is a symmetric $n \times n$ matrix, and x, b , are vectors in \mathbb{R}^n , viewed as column vectors. Actually, for reasons that will be clear shortly, it is preferable to put a factor $\frac{1}{2}$ in front of the quadratic term, so that

$$P(x) = \frac{1}{2}x^\top Ax - x^\top b.$$

The question is, under what conditions (on A) does $P(x)$ have a global minimum, preferably unique?

We give a complete answer to the above question in two stages:

1. In this section, we show that if A is symmetric positive definite, then $P(x)$ has a unique global minimum precisely when

$$Ax = b.$$

2. In Section 18.2, we give necessary and sufficient conditions in the general case, in terms of the pseudo-inverse of A .

We begin with the matrix version of Definition 16.2.

Definition 18.1. A symmetric *positive definite matrix* is a matrix whose eigenvalues are strictly positive, and a symmetric *positive semidefinite matrix* is a matrix whose eigenvalues are nonnegative.

Equivalent criteria are given in the following proposition.

Proposition 18.1. *Given any Euclidean space E of dimension n , the following properties hold:*

- (1) *Every self-adjoint linear map $f: E \rightarrow E$ is positive definite iff*

$$\langle x, f(x) \rangle > 0$$

for all $x \in E$ with $x \neq 0$.

- (2) *Every self-adjoint linear map $f: E \rightarrow E$ is positive semidefinite iff*

$$\langle x, f(x) \rangle \geq 0$$

for all $x \in E$.

Proof. (1) First, assume that f is positive definite. Recall that every self-adjoint linear map has an orthonormal basis (e_1, \dots, e_n) of eigenvectors, and let $\lambda_1, \dots, \lambda_n$ be the corresponding eigenvalues. With respect to this basis, for every $x = x_1e_1 + \dots + x_n e_n \neq 0$, we have

$$\langle x, f(x) \rangle = \left\langle \sum_{i=1}^n x_i e_i, f\left(\sum_{i=1}^n x_i e_i\right) \right\rangle = \left\langle \sum_{i=1}^n x_i e_i, \sum_{i=1}^n \lambda_i x_i e_i \right\rangle = \sum_{i=1}^n \lambda_i x_i^2,$$

which is strictly positive, since $\lambda_i > 0$ for $i = 1, \dots, n$, and $x_i^2 > 0$ for some i , since $x \neq 0$.

Conversely, assume that

$$\langle x, f(x) \rangle > 0$$

for all $x \neq 0$. Then for $x = e_i$, we get

$$\langle e_i, f(e_i) \rangle = \langle e_i, \lambda_i e_i \rangle = \lambda_i,$$

and thus $\lambda_i > 0$ for all $i = 1, \dots, n$.

(2) As in (1), we have

$$\langle x, f(x) \rangle = \sum_{i=1}^n \lambda_i x_i^2,$$

and since $\lambda_i \geq 0$ for $i = 1, \dots, n$ because f is positive semidefinite, we have $\langle x, f(x) \rangle \geq 0$, as claimed. The converse is as in (1) except that we get only $\lambda_i \geq 0$ since $\langle e_i, f(e_i) \rangle \geq 0$. \square

Some special notation is customary (especially in the field of convex optimization) to express that a symmetric matrix is positive definite or positive semidefinite.

Definition 18.2. Given any $n \times n$ symmetric matrix A we write $A \succeq 0$ if A is positive semidefinite and we write $A \succ 0$ if A is positive definite.

It should be noted that we can define the relation

$$A \succeq B$$

between any two $n \times n$ matrices (symmetric or not) iff $A - B$ is symmetric positive semidefinite. It is easy to check that this relation is actually a partial order on matrices, called the *positive semidefinite cone ordering*; for details, see Boyd and Vandenberghe [17], Section 2.4.

If A is symmetric positive definite, it is easily checked that A^{-1} is also symmetric positive definite. Also, if C is a symmetric positive definite $m \times m$ matrix and A is an $m \times n$ matrix of rank n (and so $m \geq n$), then $A^\top C A$ is symmetric positive definite.

We can now prove that

$$P(x) = \frac{1}{2} x^\top A x - x^\top b$$

has a global minimum when A is symmetric positive definite.

Proposition 18.2. *Given a quadratic function*

$$P(x) = \frac{1}{2} x^\top A x - x^\top b,$$

if A is symmetric positive definite, then $P(x)$ has a unique global minimum for the solution of the linear system $Ax = b$. The minimum value of $P(x)$ is

$$P(A^{-1}b) = -\frac{1}{2} b^\top A^{-1} b.$$

Proof. Since A is positive definite, it is invertible, since its eigenvalues are all strictly positive. Let $x = A^{-1}b$, and compute $P(y) - P(x)$ for any $y \in \mathbb{R}^n$. Since $Ax = b$, we get

$$\begin{aligned} P(y) - P(x) &= \frac{1}{2}y^\top Ay - y^\top b - \frac{1}{2}x^\top Ax + x^\top b \\ &= \frac{1}{2}y^\top Ay - y^\top Ax + \frac{1}{2}x^\top Ax \\ &= \frac{1}{2}(y - x)^\top A(y - x). \end{aligned}$$

Since A is positive definite, the last expression is nonnegative, and thus

$$P(y) \geq P(x)$$

for all $y \in \mathbb{R}^n$, which proves that $x = A^{-1}b$ is a global minimum of $P(x)$. A simple computation yields

$$P(A^{-1}b) = -\frac{1}{2}b^\top A^{-1}b.$$

□

Remarks:

- (1) The quadratic function $P(x)$ is also given by

$$P(x) = \frac{1}{2}x^\top Ax - b^\top x,$$

but the definition using $x^\top b$ is more convenient for the proof of Proposition 18.2.

- (2) If $P(x)$ contains a constant term $c \in \mathbb{R}$, so that

$$P(x) = \frac{1}{2}x^\top Ax - x^\top b + c,$$

the proof of Proposition 18.2 still shows that $P(x)$ has a unique global minimum for $x = A^{-1}b$, but the minimal value is

$$P(A^{-1}b) = -\frac{1}{2}b^\top A^{-1}b + c.$$

Thus, when the energy function $P(x)$ of a system is given by a quadratic function

$$P(x) = \frac{1}{2}x^\top Ax - x^\top b,$$

where A is symmetric positive definite, finding the global minimum of $P(x)$ is equivalent to solving the linear system $Ax = b$. Sometimes, it is useful to recast a linear problem $Ax = b$

as a variational problem (finding the minimum of some energy function). However, very often, a minimization problem comes with extra constraints that must be satisfied for all admissible solutions. For instance, we may want to minimize the quadratic function

$$Q(y_1, y_2) = \frac{1}{2}(y_1^2 + y_2^2)$$

subject to the constraint

$$2y_1 - y_2 = 5.$$

The solution for which $Q(y_1, y_2)$ is minimum is no longer $(y_1, y_2) = (0, 0)$, but instead, $(y_1, y_2) = (2, -1)$, as will be shown later.

Geometrically, the graph of the function defined by $z = Q(y_1, y_2)$ in \mathbb{R}^3 is a paraboloid of revolution P with axis of revolution Oz . The constraint

$$2y_1 - y_2 = 5$$

corresponds to the vertical plane H parallel to the z -axis and containing the line of equation $2y_1 - y_2 = 5$ in the xy -plane. Thus, the constrained minimum of Q is located on the parabola that is the intersection of the paraboloid P with the plane H .

A nice way to solve constrained minimization problems of the above kind is to use the method of *Lagrange multipliers*. But first, let us define precisely what kind of minimization problems we intend to solve.

Definition 18.3. The *quadratic constrained minimization problem* consists in minimizing a quadratic function

$$Q(y) = \frac{1}{2}y^\top C^{-1}y - b^\top y$$

subject to the linear constraints

$$A^\top y = f,$$

where C^{-1} is an $m \times m$ symmetric positive definite matrix, A is an $m \times n$ matrix of rank n (so that $m \geq n$), and where $b, y \in \mathbb{R}^m$ (viewed as column vectors), and $f \in \mathbb{R}^n$ (viewed as a column vector).

The reason for using C^{-1} instead of C is that the constrained minimization problem has an interpretation as a set of equilibrium equations in which the matrix that arises naturally is C (see Strang [104]). Since C and C^{-1} are both symmetric positive definite, this doesn't make any difference, but it seems preferable to stick to Strang's notation.

The method of Lagrange consists in incorporating the n constraints $A^\top y = f$ into the quadratic function $Q(y)$, by introducing extra variables $\lambda = (\lambda_1, \dots, \lambda_n)$ called *Lagrange multipliers*, one for each constraint. We form the *Lagrangian*

$$L(y, \lambda) = Q(y) + \lambda^\top (A^\top y - f) = \frac{1}{2}y^\top C^{-1}y - (b - A\lambda)^\top y - \lambda^\top f.$$

We shall prove that our constrained minimization problem has a unique solution given by the system of linear equations

$$\begin{aligned} C^{-1}y + A\lambda &= b, \\ A^\top y &= f, \end{aligned}$$

which can be written in matrix form as

$$\begin{pmatrix} C^{-1} & A \\ A^\top & 0 \end{pmatrix} \begin{pmatrix} y \\ \lambda \end{pmatrix} = \begin{pmatrix} b \\ f \end{pmatrix}.$$

Note that the matrix of this system is symmetric. Eliminating y from the first equation

$$C^{-1}y + A\lambda = b,$$

we get

$$y = C(b - A\lambda),$$

and substituting into the second equation, we get

$$A^\top C(b - A\lambda) = f,$$

that is,

$$A^\top CA\lambda = A^\top Cb - f.$$

However, by a previous remark, since C is symmetric positive definite and the columns of A are linearly independent, $A^\top CA$ is symmetric positive definite, and thus invertible. Note that this way of solving the system requires solving for the Lagrange multipliers first.

Letting $e = b - A\lambda$, we also note that the system

$$\begin{pmatrix} C^{-1} & A \\ A^\top & 0 \end{pmatrix} \begin{pmatrix} y \\ \lambda \end{pmatrix} = \begin{pmatrix} b \\ f \end{pmatrix}$$

is equivalent to the system

$$\begin{aligned} e &= b - A\lambda, \\ y &= Ce, \\ A^\top y &= f. \end{aligned}$$

The latter system is called the *equilibrium equations* by Strang [104]. Indeed, Strang shows that the equilibrium equations of many physical systems can be put in the above form. This includes spring-mass systems, electrical networks, and trusses, which are structures built from elastic bars. In each case, y , e , b , C , λ , f , and $K = A^\top CA$ have a physical

interpretation. The matrix $K = A^\top CA$ is usually called the *stiffness matrix*. Again, the reader is referred to Strang [104].

In order to prove that our constrained minimization problem has a unique solution, we proceed to prove that the constrained minimization of $Q(y)$ subject to $A^\top y = f$ is equivalent to the unconstrained maximization of another function $-P(\lambda)$. We get $P(\lambda)$ by minimizing the Lagrangian $L(y, \lambda)$ treated as a function of y alone. Since C^{-1} is symmetric positive definite and

$$L(y, \lambda) = \frac{1}{2}y^\top C^{-1}y - (b - A\lambda)^\top y - \lambda^\top f,$$

by Proposition 18.2 the global minimum (with respect to y) of $L(y, \lambda)$ is obtained for the solution y of

$$C^{-1}y = b - A\lambda,$$

that is, when

$$y = C(b - A\lambda),$$

and the minimum of $L(y, \lambda)$ is

$$\min_y L(y, \lambda) = -\frac{1}{2}(A\lambda - b)^\top C(A\lambda - b) - \lambda^\top f.$$

Letting

$$P(\lambda) = \frac{1}{2}(A\lambda - b)^\top C(A\lambda - b) + \lambda^\top f,$$

we claim that the solution of the constrained minimization of $Q(y)$ subject to $A^\top y = f$ is equivalent to the unconstrained maximization of $-P(\lambda)$. Of course, since we minimized $L(y, \lambda)$ with respect to y , we have

$$L(y, \lambda) \geq -P(\lambda)$$

for all y and all λ . However, when the constraint $A^\top y = f$ holds, $L(y, \lambda) = Q(y)$, and thus for any admissible y , which means that $A^\top y = f$, we have

$$\min_y Q(y) \geq \max_\lambda -P(\lambda).$$

In order to prove that the unique minimum of the constrained problem $Q(y)$ subject to $A^\top y = f$ is the unique maximum of $-P(\lambda)$, we compute $Q(y) + P(\lambda)$.

Proposition 18.3. *The quadratic constrained minimization problem of Definition 18.3 has a unique solution (y, λ) given by the system*

$$\begin{pmatrix} C^{-1} & A \\ A^\top & 0 \end{pmatrix} \begin{pmatrix} y \\ \lambda \end{pmatrix} = \begin{pmatrix} b \\ f \end{pmatrix}.$$

Furthermore, the component λ of the above solution is the unique value for which $-P(\lambda)$ is maximum.

Proof. As we suggested earlier, let us compute $Q(y) + P(\lambda)$, assuming that the constraint $A^\top y = f$ holds. Eliminating f , since $b^\top y = y^\top b$ and $\lambda^\top A^\top y = y^\top A\lambda$, we get

$$\begin{aligned} Q(y) + P(\lambda) &= \frac{1}{2}y^\top C^{-1}y - b^\top y + \frac{1}{2}(A\lambda - b)^\top C(A\lambda - b) + \lambda^\top f \\ &= \frac{1}{2}(C^{-1}y + A\lambda - b)^\top C(C^{-1}y + A\lambda - b). \end{aligned}$$

Since C is positive definite, the last expression is nonnegative. In fact, it is null iff

$$C^{-1}y + A\lambda - b = 0,$$

that is,

$$C^{-1}y + A\lambda = b.$$

But then the unique constrained minimum of $Q(y)$ subject to $A^\top y = f$ is equal to the unique maximum of $-P(\lambda)$ exactly when $A^\top y = f$ and $C^{-1}y + A\lambda = b$, which proves the proposition. \square

Remarks:

- (1) There is a form of duality going on in this situation. The constrained minimization of $Q(y)$ subject to $A^\top y = f$ is called the *primal problem*, and the unconstrained maximization of $-P(\lambda)$ is called the *dual problem*. Duality is the fact stated slightly loosely as

$$\min_y Q(y) = \max_\lambda -P(\lambda).$$

Recalling that $e = b - A\lambda$, since

$$P(\lambda) = \frac{1}{2}(A\lambda - b)^\top C(A\lambda - b) + \lambda^\top f,$$

we can also write

$$P(\lambda) = \frac{1}{2}e^\top C e + \lambda^\top f.$$

This expression often represents the total potential energy of a system. Again, the optimal solution is the one that minimizes the potential energy (and thus maximizes $-P(\lambda)$).

- (2) It is immediately verified that the equations of Proposition 18.3 are equivalent to the equations stating that the partial derivatives of the Lagrangian $L(y, \lambda)$ are null:

$$\begin{aligned} \frac{\partial L}{\partial y_i} &= 0, \quad i = 1, \dots, m, \\ \frac{\partial L}{\partial \lambda_j} &= 0, \quad j = 1, \dots, n. \end{aligned}$$

Thus, the constrained minimum of $Q(y)$ subject to $A^\top y = f$ is an extremum of the Lagrangian $L(y, \lambda)$. As we showed in Proposition 18.3, this extremum corresponds to simultaneously minimizing $L(y, \lambda)$ with respect to y and maximizing $L(y, \lambda)$ with respect to λ . Geometrically, such a point is a *saddle point* for $L(y, \lambda)$.

- (3) The Lagrange multipliers sometimes have a natural physical meaning. For example, in the spring-mass system they correspond to node displacements. In some general sense, Lagrange multipliers are correction terms needed to satisfy equilibrium equations and the price paid for the constraints. For more details, see Strang [104].

Going back to the constrained minimization of $Q(y_1, y_2) = \frac{1}{2}(y_1^2 + y_2^2)$ subject to

$$2y_1 - y_2 = 5,$$

the Lagrangian is

$$L(y_1, y_2, \lambda) = \frac{1}{2}(y_1^2 + y_2^2) + \lambda(2y_1 - y_2 - 5),$$

and the equations stating that the Lagrangian has a saddle point are

$$\begin{aligned} y_1 + 2\lambda &= 0, \\ y_2 - \lambda &= 0, \\ 2y_1 - y_2 - 5 &= 0. \end{aligned}$$

We obtain the solution $(y_1, y_2, \lambda) = (2, -1, -1)$.

Much more should be said about the use of Lagrange multipliers in optimization or variational problems. This is a vast topic. Least squares methods and Lagrange multipliers are used to tackle many problems in computer graphics and computer vision; see Trucco and Verri [111], Metaxas [79], Jain, Katsuri, and Schunck [60], Faugeras [37], and Foley, van Dam, Feiner, and Hughes [38]. For a lucid introduction to optimization methods, see Ciarlet [24].

18.2 Quadratic Optimization: The General Case

In this section, we complete the study initiated in Section 18.1 and give necessary and sufficient conditions for the quadratic function $\frac{1}{2}x^\top Ax + x^\top b$ to have a global minimum. We begin with the following simple fact:

Proposition 18.4. *If A is an invertible symmetric matrix, then the function*

$$f(x) = \frac{1}{2}x^\top Ax + x^\top b$$

has a minimum value iff $A \succeq 0$, in which case this optimal value is obtained for a unique value of x , namely $x^ = -A^{-1}b$, and with*

$$f(A^{-1}b) = -\frac{1}{2}b^\top A^{-1}b.$$

Proof. Observe that

$$\frac{1}{2}(x + A^{-1}b)^\top A(x + A^{-1}b) = \frac{1}{2}x^\top Ax + x^\top b + \frac{1}{2}b^\top A^{-1}b.$$

Thus,

$$f(x) = \frac{1}{2}x^\top Ax + x^\top b = \frac{1}{2}(x + A^{-1}b)^\top A(x + A^{-1}b) - \frac{1}{2}b^\top A^{-1}b.$$

If A has some negative eigenvalue, say $-\lambda$ (with $\lambda > 0$), if we pick any eigenvector u of A associated with λ , then for any $\alpha \in \mathbb{R}$ with $\alpha \neq 0$, if we let $x = \alpha u - A^{-1}b$, then since $Au = -\lambda u$, we get

$$\begin{aligned} f(x) &= \frac{1}{2}(x + A^{-1}b)^\top A(x + A^{-1}b) - \frac{1}{2}b^\top A^{-1}b \\ &= \frac{1}{2}\alpha u^\top A\alpha u - \frac{1}{2}b^\top A^{-1}b \\ &= -\frac{1}{2}\alpha^2\lambda \|u\|_2^2 - \frac{1}{2}b^\top A^{-1}b, \end{aligned}$$

and since α can be made as large as we want and $\lambda > 0$, we see that f has no minimum. Consequently, in order for f to have a minimum, we must have $A \succeq 0$. In this case, since $(x + A^{-1}b)^\top A(x + A^{-1}b) \geq 0$, it is clear that the minimum value of f is achieved when $x + A^{-1}b = 0$, that is, $x = -A^{-1}b$. \square

Let us now consider the case of an arbitrary symmetric matrix A .

Proposition 18.5. *If A is a symmetric matrix, then the function*

$$f(x) = \frac{1}{2}x^\top Ax + x^\top b$$

has a minimum value iff $A \succeq 0$ and $(I - AA^+)b = 0$, in which case this minimum value is

$$p^* = -\frac{1}{2}b^\top A^+b.$$

Furthermore, if $A = U^\top \Sigma U$ is an SVD of A , then the optimal value is achieved by all $x \in \mathbb{R}^n$ of the form

$$x = -A^+b + U^\top \begin{pmatrix} 0 \\ z \end{pmatrix},$$

for any $z \in \mathbb{R}^{n-r}$, where r is the rank of A .

Proof. The case that A is invertible is taken care of by Proposition 18.4, so we may assume that A is singular. If A has rank $r < n$, then we can diagonalize A as

$$A = U^\top \begin{pmatrix} \Sigma_r & 0 \\ 0 & 0 \end{pmatrix} U,$$

where U is an orthogonal matrix and where Σ_r is an $r \times r$ diagonal invertible matrix. Then we have

$$\begin{aligned} f(x) &= \frac{1}{2} x^\top U^\top \begin{pmatrix} \Sigma_r & 0 \\ 0 & 0 \end{pmatrix} Ux + x^\top U^\top Ub \\ &= \frac{1}{2} (Ux)^\top \begin{pmatrix} \Sigma_r & 0 \\ 0 & 0 \end{pmatrix} Ux + (Ux)^\top Ub. \end{aligned}$$

If we write

$$Ux = \begin{pmatrix} y \\ z \end{pmatrix} \quad \text{and} \quad Ub = \begin{pmatrix} c \\ d \end{pmatrix},$$

with $y, c \in \mathbb{R}^r$ and $z, d \in \mathbb{R}^{n-r}$, we get

$$\begin{aligned} f(x) &= \frac{1}{2} (Ux)^\top \begin{pmatrix} \Sigma_r & 0 \\ 0 & 0 \end{pmatrix} Ux + (Ux)^\top Ub \\ &= \frac{1}{2} (y^\top, z^\top) \begin{pmatrix} \Sigma_r & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} y \\ z \end{pmatrix} + (y^\top, z^\top) \begin{pmatrix} c \\ d \end{pmatrix} \\ &= \frac{1}{2} y^\top \Sigma_r y + y^\top c + z^\top d. \end{aligned}$$

For $y = 0$, we get

$$f(x) = z^\top d,$$

so if $d \neq 0$, the function f has no minimum. Therefore, if f has a minimum, then $d = 0$. However, $d = 0$ means that

$$Ub = \begin{pmatrix} c \\ 0 \end{pmatrix},$$

and we know from Section 17.1 that b is in the range of A (here, U is U^\top), which is equivalent to $(I - AA^+)b = 0$. If $d = 0$, then

$$f(x) = \frac{1}{2} y^\top \Sigma_r y + y^\top c,$$

and since Σ_r is invertible, by Proposition 18.4, the function f has a minimum iff $\Sigma_r \succeq 0$, which is equivalent to $A \succeq 0$.

Therefore, we have proved that if f has a minimum, then $(I - AA^+)b = 0$ and $A \succeq 0$. Conversely, if $(I - AA^+)b = 0$ and $A \succeq 0$, what we just did proves that f does have a minimum.

When the above conditions hold, the minimum is achieved if $y = -\Sigma_r^{-1}c$, $z = 0$ and $d = 0$, that is, for x^* given by

$$Ux^* = \begin{pmatrix} -\Sigma_r^{-1}c \\ 0 \end{pmatrix} \quad \text{and} \quad Ub = \begin{pmatrix} c \\ 0 \end{pmatrix},$$

from which we deduce that

$$x^* = -U^\top \begin{pmatrix} \Sigma_r^{-1}c \\ 0 \end{pmatrix} = -U^\top \begin{pmatrix} \Sigma_r^{-1}c & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} c \\ 0 \end{pmatrix} = -U^\top \begin{pmatrix} \Sigma_r^{-1}c & 0 \\ 0 & 0 \end{pmatrix} Ub = -A^+b$$

and the minimum value of f is

$$f(x^*) = -\frac{1}{2}b^\top A^+b.$$

For any $x \in \mathbb{R}^n$ of the form

$$x = -A^+b + U^\top \begin{pmatrix} 0 \\ z \end{pmatrix},$$

for any $z \in \mathbb{R}^{n-r}$, our previous calculations show that $f(x) = -\frac{1}{2}b^\top A^+b$. \square

The case in which we add either linear constraints of the form $C^\top x = 0$ or affine constraints of the form $C^\top x = t$ (where $t \neq 0$) can be reduced to the unconstrained case using a QR -decomposition of C or N . Let us show how to do this for linear constraints of the form $C^\top x = 0$.

If we use a QR decomposition of C , by permuting the columns, we may assume that

$$C = Q^\top \begin{pmatrix} R & S \\ 0 & 0 \end{pmatrix} \Pi,$$

where R is an $r \times r$ invertible upper triangular matrix and S is an $r \times (m-r)$ matrix (C has rank r). Then, if we let

$$x = Q^\top \begin{pmatrix} y \\ z \end{pmatrix},$$

where $y \in \mathbb{R}^r$ and $z \in \mathbb{R}^{n-r}$, then $C^\top x = 0$ becomes

$$\Pi^\top \begin{pmatrix} R & 0 \\ S & 0 \end{pmatrix} Qx = \Pi^\top \begin{pmatrix} R^\top & 0 \\ S^\top & 0 \end{pmatrix} \begin{pmatrix} y \\ z \end{pmatrix} = 0,$$

which implies $y = 0$, and every solution of $C^\top x = 0$ is of the form

$$x = Q^\top \begin{pmatrix} 0 \\ z \end{pmatrix}.$$

Our original problem becomes

$$\begin{aligned} &\text{minimize} && \frac{1}{2}(y^\top, z^\top)QAQ^\top \begin{pmatrix} y \\ z \end{pmatrix} + (y^\top, z^\top)Qb \\ &\text{subject to} && y = 0, y \in \mathbb{R}^r, z \in \mathbb{R}^{n-r}. \end{aligned}$$

Thus, the constraint $C^\top x = 0$ has been eliminated, and if we write

$$QAQ^\top = \begin{pmatrix} G_{11} & G_{12} \\ G_{21} & G_{22} \end{pmatrix}$$

and

$$Qb = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}, \quad b_1 \in \mathbb{R}^r, \quad b_2 \in \mathbb{R}^{n-r},$$

our problem becomes

$$\text{minimize } \frac{1}{2} z^\top G_{22} z + z^\top b_2, \quad z \in \mathbb{R}^{n-r},$$

the problem solved in Proposition 18.5.

Constraints of the form $C^\top x = t$ (where $t \neq 0$) can be handled in a similar fashion. In this case, we may assume that C is an $n \times m$ matrix with full rank (so that $m \leq n$) and $t \in \mathbb{R}^m$. Then we use a QR -decomposition of the form

$$C = P \begin{pmatrix} R \\ 0 \end{pmatrix},$$

where P is an orthogonal matrix and R is an $m \times m$ invertible upper triangular matrix. If we write

$$x = P \begin{pmatrix} y \\ z \end{pmatrix},$$

where $y \in \mathbb{R}^m$ and $z \in \mathbb{R}^{n-m}$, the equation $C^\top x = t$ becomes

$$(R^\top, 0)P^\top x = t,$$

that is,

$$(R^\top, 0) \begin{pmatrix} y \\ z \end{pmatrix} = t,$$

which yields

$$R^\top y = t.$$

Since R is invertible, we get $y = (R^\top)^{-1}t$, and then it is easy to see that our original problem reduces to an unconstrained problem in terms of the matrix $P^\top AP$; the details are left as an exercise.

18.3 Maximizing a Quadratic Function on the Unit Sphere

In this section we discuss various quadratic optimization problems mostly arising from computer vision (image segmentation and contour grouping). These problems can be reduced to the following basic optimization problem: Given an $n \times n$ real symmetric matrix A

$$\begin{aligned} &\text{maximize} && x^\top A x \\ &\text{subject to} && x^\top x = 1, \quad x \in \mathbb{R}^n. \end{aligned}$$

In view of Proposition 17.6, the maximum value of $x^\top Ax$ on the unit sphere is equal to the largest eigenvalue λ_1 of the matrix A , and it is achieved for any unit eigenvector u_1 associated with λ_1 .

A variant of the above problem often encountered in computer vision consists in minimizing $x^\top Ax$ on the ellipsoid given by an equation of the form

$$x^\top Bx = 1,$$

where B is a symmetric positive definite matrix. Since B is positive definite, it can be diagonalized as

$$B = QDQ^\top,$$

where Q is an orthogonal matrix and D is a diagonal matrix,

$$D = \text{diag}(d_1, \dots, d_n),$$

with $d_i > 0$, for $i = 1, \dots, n$. If we define the matrices $B^{1/2}$ and $B^{-1/2}$ by

$$B^{1/2} = Q \text{diag}(\sqrt{d_1}, \dots, \sqrt{d_n}) Q^\top$$

and

$$B^{-1/2} = Q \text{diag}(1/\sqrt{d_1}, \dots, 1/\sqrt{d_n}) Q^\top,$$

it is clear that these matrices are symmetric, that $B^{-1/2}BB^{-1/2} = I$, and that $B^{1/2}$ and $B^{-1/2}$ are mutual inverses. Then, if we make the change of variable

$$x = B^{-1/2}y,$$

the equation $x^\top Bx = 1$ becomes $y^\top y = 1$, and the optimization problem

$$\begin{array}{ll} \text{maximize} & x^\top Ax \\ \text{subject to} & x^\top Bx = 1, x \in \mathbb{R}^n, \end{array}$$

is equivalent to the problem

$$\begin{array}{ll} \text{maximize} & y^\top B^{-1/2}AB^{-1/2}y \\ \text{subject to} & y^\top y = 1, y \in \mathbb{R}^n, \end{array}$$

where $y = B^{1/2}x$ and where $B^{-1/2}AB^{-1/2}$ is symmetric.

The complex version of our basic optimization problem in which A is a Hermitian matrix also arises in computer vision. Namely, given an $n \times n$ complex Hermitian matrix A ,

$$\begin{array}{ll} \text{maximize} & x^* Ax \\ \text{subject to} & x^* x = 1, x \in \mathbb{C}^n. \end{array}$$

Again by Proposition 17.6, the maximum value of x^*Ax on the unit sphere is equal to the largest eigenvalue λ_1 of the matrix A and it is achieved for any unit eigenvector u_1 associated with λ_1 .

It is worth pointing out that if A is a *skew-Hermitian* matrix, that is, if $A^* = -A$, then x^*Ax is *pure imaginary or zero*.

Indeed, since $z = x^*Ax$ is a scalar, we have $z^* = \bar{z}$ (the conjugate of z), so we have

$$\overline{x^*Ax} = (x^*Ax)^* = x^*A^*x = -x^*Ax,$$

so $\overline{x^*Ax} + x^*Ax = 2\operatorname{Re}(x^*Ax) = 0$, which means that x^*Ax is pure imaginary or zero.

In particular, if A is a real matrix and if A is *skew-symmetric*, then

$$x^\top Ax = 0.$$

Thus, for any real matrix (symmetric or not),

$$x^\top Ax = x^\top H(A)x,$$

where $H(A) = (A + A^\top)/2$, the symmetric part of A .

There are situations in which it is necessary to add linear constraints to the problem of maximizing a quadratic function on the sphere. This problem was completely solved by Golub [48] (1973). The problem is the following: Given an $n \times n$ real symmetric matrix A and an $n \times p$ matrix C ,

$$\begin{aligned} &\text{minimize} && x^\top Ax \\ &\text{subject to} && x^\top x = 1, \quad C^\top x = 0, \quad x \in \mathbb{R}^n. \end{aligned}$$

Golub shows that the linear constraint $C^\top x = 0$ can be eliminated as follows: If we use a QR decomposition of C , by permuting the columns, we may assume that

$$C = Q^\top \begin{pmatrix} R & S \\ 0 & 0 \end{pmatrix} \Pi,$$

where R is an $r \times r$ invertible upper triangular matrix and S is an $r \times (p-r)$ matrix (assuming C has rank r). Then if we let

$$x = Q^\top \begin{pmatrix} y \\ z \end{pmatrix},$$

where $y \in \mathbb{R}^r$ and $z \in \mathbb{R}^{n-r}$, then $C^\top x = 0$ becomes

$$\Pi^\top \begin{pmatrix} R^\top & 0 \\ S^\top & 0 \end{pmatrix} Qx = \Pi^\top \begin{pmatrix} R^\top & 0 \\ S^\top & 0 \end{pmatrix} \begin{pmatrix} y \\ z \end{pmatrix} = 0,$$

which implies $y = 0$, and every solution of $C^\top x = 0$ is of the form

$$x = Q^\top \begin{pmatrix} 0 \\ z \end{pmatrix}.$$

Our original problem becomes

$$\begin{aligned} &\text{minimize} && (y^\top, z^\top) Q A Q^\top \begin{pmatrix} y \\ z \end{pmatrix} \\ &\text{subject to} && z^\top z = 1, \ z \in \mathbb{R}^{n-r}, \\ &&& y = 0, \ y \in \mathbb{R}^r. \end{aligned}$$

Thus, the constraint $C^\top x = 0$ has been eliminated, and if we write

$$Q A Q^\top = \begin{pmatrix} G_{11} & G_{12} \\ G_{12}^\top & G_{22} \end{pmatrix},$$

our problem becomes

$$\begin{aligned} &\text{minimize} && z^\top G_{22} z \\ &\text{subject to} && z^\top z = 1, \ z \in \mathbb{R}^{n-r}, \end{aligned}$$

a standard eigenvalue problem. Observe that if we let

$$J = \begin{pmatrix} 0 & 0 \\ 0 & I_{n-r} \end{pmatrix},$$

then

$$J Q A Q^\top J = \begin{pmatrix} 0 & 0 \\ 0 & G_{22} \end{pmatrix},$$

and if we set

$$P = Q^\top J Q,$$

then

$$P A P = Q^\top J Q A Q^\top J Q.$$

Now, $Q^\top J Q A Q^\top J Q$ and $J Q A Q^\top J$ have the same eigenvalues, so $P A P$ and $J Q A Q^\top J$ also have the same eigenvalues. It follows that the solutions of our optimization problem are among the eigenvalues of $K = P A P$, and at least r of those are 0. Using the fact that $C C^+$ is the projection onto the range of C , where C^+ is the pseudo-inverse of C , it can also be shown that

$$P = I - C C^+,$$

the projection onto the kernel of C^\top . In particular, when $n \geq p$ and C has full rank (the columns of C are linearly independent), then we know that $C^+ = (C^\top C)^{-1} C^\top$ and

$$P = I - C (C^\top C)^{-1} C^\top.$$

This fact is used by Cour and Shi [25] and implicitly by Yu and Shi [116].

The problem of adding affine constraints of the form $N^\top x = t$, where $t \neq 0$, also comes up in practice. At first glance, this problem may not seem harder than the linear problem in which $t = 0$, but it is. This problem was extensively studied in a paper by Gander, Golub, and von Matt [45] (1989).

Gander, Golub, and von Matt consider the following problem: Given an $(n+m) \times (n+m)$ real symmetric matrix A (with $n > 0$), an $(n+m) \times m$ matrix N with full rank, and a nonzero vector $t \in \mathbb{R}^m$ with $\|(N^\top)^\dagger t\| < 1$ (where $(N^\top)^\dagger$ denotes the pseudo-inverse of N^\top),

$$\begin{aligned} & \text{minimize} && x^\top A x \\ & \text{subject to} && x^\top x = 1, \quad N^\top x = t, \quad x \in \mathbb{R}^{n+m}. \end{aligned}$$

The condition $\|(N^\top)^\dagger t\| < 1$ ensures that the problem has a solution and is not trivial. The authors begin by proving that the affine constraint $N^\top x = t$ can be eliminated. One way to do so is to use a QR decomposition of N . If

$$N = P \begin{pmatrix} R \\ 0 \end{pmatrix},$$

where P is an orthogonal matrix and R is an $m \times m$ invertible upper triangular matrix, then if we observe that

$$\begin{aligned} x^\top A x &= x^\top P P^\top A P P^\top x, \\ N^\top x &= (R^\top, 0) P^\top x = t, \\ x^\top x &= x^\top P P^\top x = 1, \end{aligned}$$

and if we write

$$P^\top A P = \begin{pmatrix} B & \Gamma^\top \\ \Gamma & C \end{pmatrix}$$

and

$$P^\top x = \begin{pmatrix} y \\ z \end{pmatrix},$$

then we get

$$\begin{aligned} x^\top A x &= y^\top B y + 2z^\top \Gamma y + z^\top C z, \\ R^\top y &= t, \\ y^\top y + z^\top z &= 1. \end{aligned}$$

Thus

$$y = (R^\top)^{-1} t,$$

and if we write

$$s^2 = 1 - y^\top y > 0$$

and

$$b = \Gamma y,$$

we get the simplified problem

$$\begin{array}{ll} \text{minimize} & z^\top C z + 2z^\top b \\ \text{subject to} & z^\top z = s^2, \quad z \in \mathbb{R}^m. \end{array}$$

Unfortunately, if $b \neq 0$, Proposition 17.6 is no longer applicable. It is still possible to find the minimum of the function $z^\top C z + 2z^\top b$ using Lagrange multipliers, but such a solution is too involved to be presented here. Interested readers will find a thorough discussion in Gander, Golub, and von Matt [45].

18.4 Summary

The main concepts and results of this chapter are listed below:

- Quadratic optimization problems; *quadratic functions*.
- Symmetric *positive definite* and *positive semidefinite* matrices.
- The *positive semidefinite cone ordering*.
- Existence of a global minimum when A is symmetric positive definite.
- Constrained quadratic optimization problems.
- *Lagrange multipliers*; *Lagrangian*.
- *Primal* and *dual* problems.
- Quadratic optimization problems: the case of a symmetric invertible matrix A .
- Quadratic optimization problems: the general case of a symmetric matrix A .
- Adding linear constraints of the form $C^\top x = 0$.
- Adding affine constraints of the form $C^\top x = t$, with $t \neq 0$.
- Maximizing a quadratic function over the unit sphere.
- Maximizing a quadratic function over an ellipsoid.
- Maximizing a Hermitian quadratic form.
- Adding linear constraints of the form $C^\top x = 0$.
- Adding affine constraints of the form $N^\top x = t$, with $t \neq 0$.

Chapter 19

Basics of Affine Geometry

L'algèbre n'est qu'une géométrie écrite; la géométrie n'est qu'une algèbre figurée.
—Sophie Germain

19.1 Affine Spaces

Geometrically, curves and surfaces are usually considered to be sets of points with some special properties, living in a space consisting of “points.” Typically, one is also interested in geometric properties invariant under certain transformations, for example, translations, rotations, projections, etc. One could model the space of points as a vector space, but this is not very satisfactory for a number of reasons. One reason is that the point corresponding to the zero vector (0), called the origin, plays a special role, when there is really no reason to have a privileged origin. Another reason is that certain notions, such as parallelism, are handled in an awkward manner. But the deeper reason is that vector spaces and affine spaces really have different geometries. The geometric properties of a vector space are invariant under the group of bijective linear maps, whereas the geometric properties of an affine space are invariant under the group of bijective affine maps, and these two groups are not isomorphic. Roughly speaking, there are more affine maps than linear maps.

Affine spaces provide a better framework for doing geometry. In particular, it is possible to deal with points, curves, surfaces, etc., in an **intrinsic manner**, that is, independently of any specific choice of a coordinate system. As in physics, this is highly desirable to really understand what is going on. Of course, coordinate systems have to be chosen to finally carry out computations, but one should learn to resist the temptation to resort to coordinate systems until it is really necessary.

Affine spaces are the right framework for dealing with motions, trajectories, and physical forces, among other things. Thus, affine geometry is crucial to a clean presentation of kinematics, dynamics, and other parts of physics (for example, elasticity). After all, a rigid motion is an affine map, but not a linear map in general. Also, given an $m \times n$ matrix A and a vector $b \in \mathbb{R}^m$, the set $U = \{x \in \mathbb{R}^n \mid Ax = b\}$ of solutions of the system $Ax = b$ is an

affine space, but not a vector space (linear space) in general.

Use coordinate systems only when needed!

This chapter proceeds as follows. We take advantage of the fact that almost every affine concept is the counterpart of some concept in linear algebra. We begin by defining affine spaces, stressing the physical interpretation of the definition in terms of points (particles) and vectors (forces). Corresponding to linear combinations of vectors, we define affine combinations of points (barycenters), realizing that we are forced to restrict our attention to families of scalars adding up to 1. Corresponding to linear subspaces, we introduce affine subspaces as subsets closed under affine combinations. Then, we characterize affine subspaces in terms of certain vector spaces called their directions. This allows us to define a clean notion of parallelism. Next, corresponding to linear independence and bases, we define affine independence and affine frames. We also define convexity. Corresponding to linear maps, we define affine maps as maps preserving affine combinations. We show that every affine map is completely defined by the image of one point and a linear map. Then, we investigate briefly some simple affine maps, the translations and the central dilatations. At this point, we give a glimpse of affine geometry. We prove the theorems of Thales, Pappus, and Desargues. After this, the definition of affine hyperplanes in terms of affine forms is reviewed. The section ends with a closer look at the intersection of affine subspaces.

Our presentation of affine geometry is far from being comprehensive, and it is biased toward the algorithmic geometry of curves and surfaces. For more details, the reader is referred to Pedoe [88], Snapper and Troyer [99], Berger [8, 9], Coxeter [26], Samuel [90], Tisseron [109], and Hilbert and Cohn-Vossen [56].

Suppose we have a particle moving in 3D space and that we want to describe the trajectory of this particle. If one looks up a good textbook on dynamics, such as Greenwood [51], one finds out that the particle is modeled as a point, and that the position of this point x is determined with respect to a “frame” in \mathbb{R}^3 by a vector. Curiously, the notion of a frame is rarely defined precisely, but it is easy to infer that a frame is a pair $(O, (e_1, e_2, e_3))$ consisting of an origin O (which is a point) together with a basis of three vectors (e_1, e_2, e_3) . For example, the standard frame in \mathbb{R}^3 has origin $O = (0, 0, 0)$ and the basis of three vectors $e_1 = (1, 0, 0)$, $e_2 = (0, 1, 0)$, and $e_3 = (0, 0, 1)$. The position of a point x is then defined by the “unique vector” from O to x .

But wait a minute, this definition seems to be defining frames and the position of a point without defining what a point is! Well, let us identify points with elements of \mathbb{R}^3 . If so, given any two points $a = (a_1, a_2, a_3)$ and $b = (b_1, b_2, b_3)$, there is a unique *free vector*, denoted by \overrightarrow{ab} , from a to b , the vector $\overrightarrow{ab} = (b_1 - a_1, b_2 - a_2, b_3 - a_3)$. Note that

$$b = a + \overrightarrow{ab},$$

addition being understood as addition in \mathbb{R}^3 . Then, in the standard frame, given a point $x = (x_1, x_2, x_3)$, the position of x is the vector $\overrightarrow{Ox} = (x_1, x_2, x_3)$, which coincides with the point itself. In the standard frame, points and vectors are identified. Points and free vectors are illustrated in Figure 19.1.

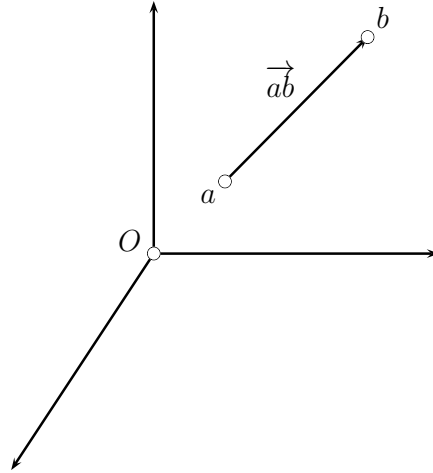


Figure 19.1: Points and free vectors

What if we pick a frame with a different origin, say $\Omega = (\omega_1, \omega_2, \omega_3)$, but the same basis vectors (e_1, e_2, e_3) ? This time, the point $x = (x_1, x_2, x_3)$ is defined by two position vectors:

$$\overrightarrow{Ox} = (x_1, x_2, x_3)$$

in the frame $(O, (e_1, e_2, e_3))$ and

$$\overrightarrow{\Omega x} = (x_1 - \omega_1, x_2 - \omega_2, x_3 - \omega_3)$$

in the frame $(\Omega, (e_1, e_2, e_3))$.

This is because

$$\overrightarrow{Ox} = \overrightarrow{O\Omega} + \overrightarrow{\Omega x} \quad \text{and} \quad \overrightarrow{O\Omega} = (\omega_1, \omega_2, \omega_3).$$

We note that in the second frame $(\Omega, (e_1, e_2, e_3))$, points and position vectors are no longer identified. This gives us evidence that points are not vectors. It may be computationally convenient to deal with points using position vectors, but such a treatment is not frame invariant, which has undesirable effects.

Inspired by physics, we deem it important to define points and properties of points that are frame invariant. An undesirable side effect of the present approach shows up if we attempt to define linear combinations of points. First, let us review the notion of linear combination of vectors. Given two vectors u and v of coordinates (u_1, u_2, u_3) and (v_1, v_2, v_3) with respect to the basis (e_1, e_2, e_3) , for any two scalars λ, μ , we can define the linear combination $\lambda u + \mu v$ as the vector of coordinates

$$(\lambda u_1 + \mu v_1, \lambda u_2 + \mu v_2, \lambda u_3 + \mu v_3).$$

If we choose a different basis (e'_1, e'_2, e'_3) and if the matrix P expressing the vectors (e'_1, e'_2, e'_3) over the basis (e_1, e_2, e_3) is

$$P = \begin{pmatrix} a_1 & b_1 & c_1 \\ a_2 & b_2 & c_2 \\ a_3 & b_3 & c_3 \end{pmatrix},$$

which means that the columns of P are the coordinates of the e'_j over the basis (e_1, e_2, e_3) , since

$$u_1 e_1 + u_2 e_2 + u_3 e_3 = u'_1 e'_1 + u'_2 e'_2 + u'_3 e'_3$$

and

$$v_1 e_1 + v_2 e_2 + v_3 e_3 = v'_1 e'_1 + v'_2 e'_2 + v'_3 e'_3,$$

it is easy to see that the coordinates (u_1, u_2, u_3) and (v_1, v_2, v_3) of u and v with respect to the basis (e_1, e_2, e_3) are given in terms of the coordinates (u'_1, u'_2, u'_3) and (v'_1, v'_2, v'_3) of u and v with respect to the basis (e'_1, e'_2, e'_3) by the matrix equations

$$\begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} = P \begin{pmatrix} u'_1 \\ u'_2 \\ u'_3 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix} = P \begin{pmatrix} v'_1 \\ v'_2 \\ v'_3 \end{pmatrix}.$$

From the above, we get

$$\begin{pmatrix} u'_1 \\ u'_2 \\ u'_3 \end{pmatrix} = P^{-1} \begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} v'_1 \\ v'_2 \\ v'_3 \end{pmatrix} = P^{-1} \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix},$$

and by linearity, the coordinates

$$(\lambda u'_1 + \mu v'_1, \lambda u'_2 + \mu v'_2, \lambda u'_3 + \mu v'_3)$$

of $\lambda u + \mu v$ with respect to the basis (e'_1, e'_2, e'_3) are given by

$$\begin{pmatrix} \lambda u'_1 + \mu v'_1 \\ \lambda u'_2 + \mu v'_2 \\ \lambda u'_3 + \mu v'_3 \end{pmatrix} = \lambda P^{-1} \begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} + \mu P^{-1} \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix} = P^{-1} \begin{pmatrix} \lambda u_1 + \mu v_1 \\ \lambda u_2 + \mu v_2 \\ \lambda u_3 + \mu v_3 \end{pmatrix}.$$

Everything worked out because the change of basis does not involve a change of origin. On the other hand, if we consider the change of frame from the frame $(O, (e_1, e_2, e_3))$ to the frame $(\Omega, (e_1, e_2, e_3))$, where $\overrightarrow{O\Omega} = (\omega_1, \omega_2, \omega_3)$, given two points a, b of coordinates (a_1, a_2, a_3) and (b_1, b_2, b_3) with respect to the frame $(O, (e_1, e_2, e_3))$ and of coordinates (a'_1, a'_2, a'_3) and (b'_1, b'_2, b'_3) with respect to the frame $(\Omega, (e_1, e_2, e_3))$, since

$$(a'_1, a'_2, a'_3) = (a_1 - \omega_1, a_2 - \omega_2, a_3 - \omega_3)$$

and

$$(b'_1, b'_2, b'_3) = (b_1 - \omega_1, b_2 - \omega_2, b_3 - \omega_3),$$

the coordinates of $\lambda a + \mu b$ with respect to the frame $(O, (e_1, e_2, e_3))$ are

$$(\lambda a_1 + \mu b_1, \lambda a_2 + \mu b_2, \lambda a_3 + \mu b_3),$$

but the coordinates

$$(\lambda a'_1 + \mu b'_1, \lambda a'_2 + \mu b'_2, \lambda a'_3 + \mu b'_3)$$

of $\lambda a + \mu b$ with respect to the frame $(\Omega, (e_1, e_2, e_3))$ are

$$(\lambda a_1 + \mu b_1 - (\lambda + \mu)\omega_1, \lambda a_2 + \mu b_2 - (\lambda + \mu)\omega_2, \lambda a_3 + \mu b_3 - (\lambda + \mu)\omega_3),$$

which are different from

$$(\lambda a_1 + \mu b_1 - \omega_1, \lambda a_2 + \mu b_2 - \omega_2, \lambda a_3 + \mu b_3 - \omega_3),$$

unless $\lambda + \mu = 1$.

Thus, we have discovered a major difference between vectors and points: The notion of linear combination of vectors is basis independent, but the notion of linear combination of points is frame dependent. In order to salvage the notion of linear combination of points, some restriction is needed: The scalar coefficients must add up to 1.

A clean way to handle the problem of frame invariance and to deal with points in a more intrinsic manner is to make a clearer distinction between points and vectors. We duplicate \mathbb{R}^3 into two copies, the first copy corresponding to points, where we forget the vector space structure, and the second copy corresponding to free vectors, where the vector space structure is important. Furthermore, we make explicit the important fact that the vector space \mathbb{R}^3 acts on the set of points \mathbb{R}^3 : Given any **point** $a = (a_1, a_2, a_3)$ and any **vector** $v = (v_1, v_2, v_3)$, we obtain the **point**

$$a + v = (a_1 + v_1, a_2 + v_2, a_3 + v_3),$$

which can be thought of as the result of translating a to b using the vector v . We can imagine that v is placed such that its origin coincides with a and that its tip coincides with b . This action $+: \mathbb{R}^3 \times \mathbb{R}^3 \rightarrow \mathbb{R}^3$ satisfies some crucial properties. For example,

$$\begin{aligned} a + 0 &= a, \\ (a + u) + v &= a + (u + v), \end{aligned}$$

and for any two points a, b , there is a unique free vector \overrightarrow{ab} such that

$$b = a + \overrightarrow{ab}.$$

It turns out that the above properties, although trivial in the case of \mathbb{R}^3 , are all that is needed to define the abstract notion of affine space (or affine structure). The basic idea is to consider two (distinct) sets E and \overrightarrow{E} , where E is a set of points (with no structure) and \overrightarrow{E} is a vector space (of free vectors) acting on the set E .

Did you say “A fine space”?

Intuitively, we can think of the elements of \vec{E} as forces moving the points in E , considered as physical particles. The effect of applying a force (free vector) $u \in \vec{E}$ to a point $a \in E$ is a translation. By this, we mean that for every force $u \in \vec{E}$, the action of the force u is to “move” every point $a \in E$ to the point $a + u \in E$ obtained by the translation corresponding to u viewed as a vector. Since translations can be composed, it is natural that \vec{E} is a vector space.

For simplicity, it is assumed that all vector spaces under consideration are defined over the field \mathbb{R} of real numbers. Most of the definitions and results also hold for an arbitrary field K , although some care is needed when dealing with fields of characteristic different from zero (see the problems). It is also assumed that all families $(\lambda_i)_{i \in I}$ of scalars have finite support. Recall that a family $(\lambda_i)_{i \in I}$ of scalars has *finite support* if $\lambda_i = 0$ for all $i \in I - J$, where J is a finite subset of I . Obviously, finite families of scalars have finite support, and for simplicity, the reader may assume that all families of scalars are finite. The formal definition of an affine space is as follows.

Definition 19.1. An *affine space* is either the degenerate space reduced to the empty set, or a triple $\langle E, \vec{E}, + \rangle$ consisting of a nonempty set E (of *points*), a vector space \vec{E} (of *translations*, or *free vectors*), and an action $+: E \times \vec{E} \rightarrow E$, satisfying the following conditions.

(A1) $a + 0 = a$, for every $a \in E$.

(A2) $(a + u) + v = a + (u + v)$, for every $a \in E$, and every $u, v \in \vec{E}$.

(A3) For any two points $a, b \in E$, there is a unique $u \in \vec{E}$ such that $a + u = b$.

The unique vector $u \in \vec{E}$ such that $a + u = b$ is denoted by \overrightarrow{ab} , or sometimes by \mathbf{ab} , or even by $b - a$. Thus, we also write

$$b = a + \overrightarrow{ab}$$

(or $b = a + \mathbf{ab}$, or even $b = a + (b - a)$).

The *dimension of the affine space* $\langle E, \vec{E}, + \rangle$ is the dimension $\dim(\vec{E})$ of the vector space \vec{E} . For simplicity, it is denoted by $\dim(E)$.

Conditions (A1) and (A2) say that the (abelian) group \vec{E} acts on E , and condition (A3) says that \vec{E} acts transitively and faithfully on E . Note that

$$\overrightarrow{a(a+v)} = v$$

for all $a \in E$ and all $v \in \vec{E}$, since $\overrightarrow{a(a+v)}$ is the unique vector such that $a + v = a + \overrightarrow{a(a+v)}$. Thus, $b = a + v$ is equivalent to $\overrightarrow{ab} = v$. Figure 19.2 gives an intuitive picture of an affine space. It is natural to think of all vectors as having the same origin, the null vector.

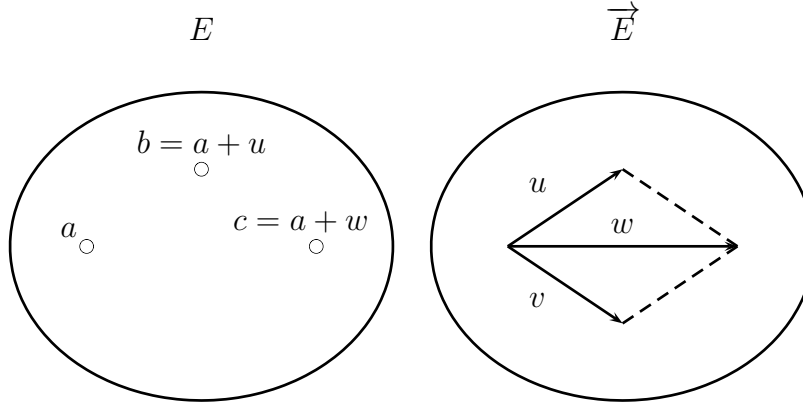


Figure 19.2: Intuitive picture of an affine space

The axioms defining an affine space $\langle E, \vec{E}, + \rangle$ can be interpreted intuitively as saying that E and \vec{E} are two different ways of looking at the same object, but wearing different sets of glasses, the second set of glasses depending on the choice of an “origin” in E . Indeed, we can choose to look at the points in E , forgetting that every pair (a, b) of points defines a unique vector \vec{ab} in \vec{E} , or we can choose to look at the vectors u in \vec{E} , forgetting the points in E . Furthermore, if we also pick any point a in E , a point that can be viewed as an *origin* in E , then we can recover all the points in E as the translated points $a + u$ for all $u \in \vec{E}$. This can be formalized by defining two maps between E and \vec{E} .

For every $a \in E$, consider the mapping from \vec{E} to E given by

$$u \mapsto a + u,$$

where $u \in \vec{E}$, and consider the mapping from E to \vec{E} given by

$$b \mapsto \vec{ab},$$

where $b \in E$. The composition of the first mapping with the second is

$$u \mapsto a + u \mapsto \overrightarrow{a(a + u)},$$

which, in view of (A3), yields u . The composition of the second with the first mapping is

$$b \mapsto \vec{ab} \mapsto a + \vec{ab},$$

which, in view of (A3), yields b . Thus, these compositions are the identity from \vec{E} to \vec{E} and the identity from E to E , and the mappings are both bijections.

When we identify E with \vec{E} via the mapping $b \mapsto \vec{ab}$, we say that we consider E as the vector space obtained by taking a as the origin in E , and we denote it by E_a . Because E_a is

a vector space, to be consistent with our notational conventions we should use the notation \vec{E}_a (using an arrow), instead of E_a . However, for simplicity, we stick to the notation E_a .

Thus, an affine space $\langle E, \vec{E}, + \rangle$ is a way of defining a vector space structure on a set of points E , without making a commitment to a **fixed** origin in E . Nevertheless, as soon as we commit to an origin a in E , we can view E as the vector space E_a . However, we urge the reader to think of E as a physical set of points and of \vec{E} as a set of forces acting on E , rather than reducing E to some isomorphic copy of \mathbb{R}^n . After all, points are points, and not vectors! For notational simplicity, we will often denote an affine space $\langle E, \vec{E}, + \rangle$ by (E, \vec{E}) , or even by E . The vector space \vec{E} is called the *vector space associated with E* .



One should be careful about the overloading of the addition symbol $+$. Addition is well-defined on vectors, as in $u + v$; the translate $a + u$ of a point $a \in E$ by a vector $u \in \vec{E}$ is also well-defined, but addition of points $a + b$ **does not make sense**. In this respect, the notation $b - a$ for the unique vector u such that $b = a + u$ is somewhat confusing, since it suggests that points can be subtracted (but not added!).

Any vector space \vec{E} has an affine space structure specified by choosing $E = \vec{E}$, and letting $+$ be addition in the vector space \vec{E} . We will refer to the affine structure $\langle \vec{E}, \vec{E}, + \rangle$ on a vector space \vec{E} as the *canonical (or natural) affine structure on \vec{E}* . In particular, the vector space \mathbb{R}^n can be viewed as the affine space $\langle \mathbb{R}^n, \mathbb{R}^n, + \rangle$, denoted by \mathbb{A}^n . In general, if K is any field, the affine space $\langle K^n, K^n, + \rangle$ is denoted by \mathbb{A}_K^n . In order to distinguish between the double role played by members of \mathbb{R}^n , points and vectors, we will denote points by row vectors, and vectors by column vectors. Thus, the action of the vector space \mathbb{R}^n over the set \mathbb{R}^n simply viewed as a set of points is given by

$$(a_1, \dots, a_n) + \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix} = (a_1 + u_1, \dots, a_n + u_n).$$

We will also use the convention that if $x = (x_1, \dots, x_n) \in \mathbb{R}^n$, then the column vector associated with x is denoted by \mathbf{x} (in boldface notation). Abusing the notation slightly, if $a \in \mathbb{R}^n$ is a point, we also write $a \in \mathbb{A}^n$. The affine space \mathbb{A}^n is called the *real affine space of dimension n* . In most cases, we will consider $n = 1, 2, 3$.

19.2 Examples of Affine Spaces

Let us now give an example of an affine space that is not given as a vector space (at least, not in an obvious fashion). Consider the subset L of \mathbb{A}^2 consisting of all points (x, y) satisfying the equation

$$x + y - 1 = 0.$$

The set L is the line of slope -1 passing through the points $(1, 0)$ and $(0, 1)$ shown in Figure 19.3.

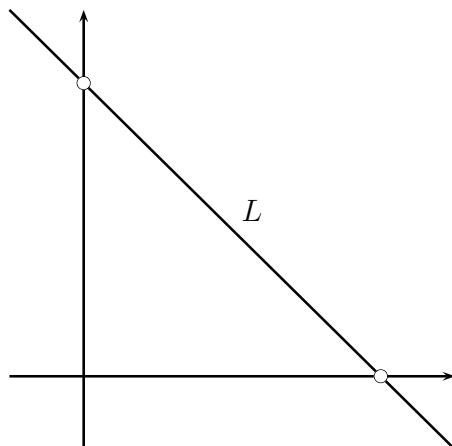


Figure 19.3: An affine space: the line of equation $x + y - 1 = 0$

The line L can be made into an official affine space by defining the action $+: L \times \mathbb{R} \rightarrow L$ of \mathbb{R} on L defined such that for every point $(x, 1 - x)$ on L and any $u \in \mathbb{R}$,

$$(x, 1 - x) + u = (x + u, 1 - x - u).$$

It is immediately verified that this action makes L into an affine space. For example, for any two points $a = (a_1, 1 - a_1)$ and $b = (b_1, 1 - b_1)$ on L , the unique (vector) $u \in \mathbb{R}$ such that $b = a + u$ is $u = b_1 - a_1$. Note that the vector space \mathbb{R} is isomorphic to the line of equation $x + y = 0$ passing through the origin.

Similarly, consider the subset H of \mathbb{A}^3 consisting of all points (x, y, z) satisfying the equation

$$x + y + z - 1 = 0.$$

The set H is the plane passing through the points $(1, 0, 0)$, $(0, 1, 0)$, and $(0, 0, 1)$. The plane H can be made into an official affine space by defining the action $+: H \times \mathbb{R}^2 \rightarrow H$ of \mathbb{R}^2 on H defined such that for every point $(x, y, 1 - x - y)$ on H and any $\begin{pmatrix} u \\ v \end{pmatrix} \in \mathbb{R}^2$,

$$(x, y, 1 - x - y) + \begin{pmatrix} u \\ v \end{pmatrix} = (x + u, y + v, 1 - x - u - y - v).$$

For a slightly wilder example, consider the subset P of \mathbb{A}^3 consisting of all points (x, y, z) satisfying the equation

$$x^2 + y^2 - z = 0.$$

The set P is a paraboloid of revolution, with axis Oz . The surface P can be made into an official affine space by defining the action $+: P \times \mathbb{R}^2 \rightarrow P$ of \mathbb{R}^2 on P defined such that for

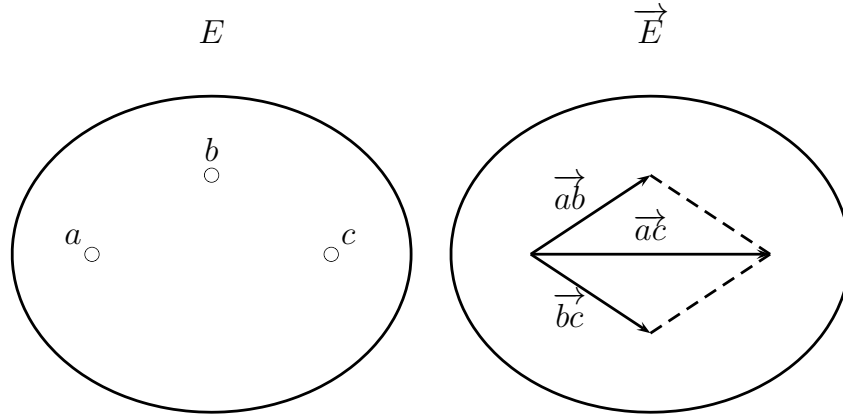


Figure 19.4: Points and corresponding vectors in affine geometry

every point $(x, y, x^2 + y^2)$ on P and any $\begin{pmatrix} u \\ v \end{pmatrix} \in \mathbb{R}^2$,

$$(x, y, x^2 + y^2) + \begin{pmatrix} u \\ v \end{pmatrix} = (x + u, y + v, (x + u)^2 + (y + v)^2).$$

This should dispell any idea that affine spaces are dull. Affine spaces not already equipped with an obvious vector space structure arise in projective geometry.

19.3 Chasles's Identity

Given any three points $a, b, c \in E$, since $c = a + \vec{ac}$, $b = a + \vec{ab}$, and $c = b + \vec{bc}$, we get

$$c = b + \vec{bc} = (a + \vec{ab}) + \vec{bc} = a + (\vec{ab} + \vec{bc})$$

by (A2), and thus, by (A3),

$$\vec{ab} + \vec{bc} = \vec{ac},$$

which is known as *Chasles's identity*, and illustrated in Figure 19.4.

Since $a = a + \vec{aa}$ and by (A1) $a = a + 0$, by (A3) we get

$$\vec{aa} = 0.$$

Thus, letting $a = c$ in Chasles's identity, we get

$$\vec{ba} = -\vec{ab}.$$

Given any four points $a, b, c, d \in E$, since by Chasles's identity

$$\vec{ab} + \vec{bc} = \vec{ad} + \vec{dc} = \vec{ac},$$

we have the *parallelogram law*

$$\vec{ab} = \vec{dc} \quad \text{iff} \quad \vec{bc} = \vec{ad}.$$

19.4 Affine Combinations, Barycenters

A fundamental concept in linear algebra is that of a linear combination. The corresponding concept in affine geometry is that of an *affine combination*, also called a *barycenter*. However, there is a problem with the naive approach involving a coordinate system, as we saw in Section 19.1. Since this problem is the reason for introducing affine combinations, at the risk of boring certain readers, we give another example showing what goes wrong if we are not careful in defining linear combinations of points.

Consider \mathbb{R}^2 as an affine space, under its natural coordinate system with origin $O = (0, 0)$ and basis vectors $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$ and $\begin{pmatrix} 0 \\ 1 \end{pmatrix}$. Given any two points $a = (a_1, a_2)$ and $b = (b_1, b_2)$, it is natural to define the affine combination $\lambda a + \mu b$ as the point of coordinates

$$(\lambda a_1 + \mu b_1, \lambda a_2 + \mu b_2).$$

Thus, when $a = (-1, -1)$ and $b = (2, 2)$, the point $a + b$ is the point $c = (1, 1)$.

Let us now consider the new coordinate system with respect to the origin $c = (1, 1)$ (and the same basis vectors). This time, the coordinates of a are $(-2, -2)$, the coordinates of b are $(1, 1)$, and the point $a + b$ is the point d of coordinates $(-1, -1)$. However, it is clear that the point d is identical to the origin $O = (0, 0)$ of the first coordinate system.

Thus, $a + b$ corresponds to two different points depending on which coordinate system is used for its computation!

This shows that some extra condition is needed in order for affine combinations to make sense. It turns out that if the scalars sum up to 1, the definition is intrinsic, as the following lemma shows.

Lemma 19.1. *Given an affine space E , let $(a_i)_{i \in I}$ be a family of points in E , and let $(\lambda_i)_{i \in I}$ be a family of scalars. For any two points $a, b \in E$, the following properties hold:*

(1) *If $\sum_{i \in I} \lambda_i = 1$, then*

$$a + \sum_{i \in I} \lambda_i \vec{aa_i} = b + \sum_{i \in I} \lambda_i \vec{ba_i}.$$

(2) *If $\sum_{i \in I} \lambda_i = 0$, then*

$$\sum_{i \in I} \lambda_i \vec{aa_i} = \sum_{i \in I} \lambda_i \vec{ba_i}.$$

Proof. (1) By Chasles's identity (see Section 19.3), we have

$$\begin{aligned}
 a + \sum_{i \in I} \lambda_i \overrightarrow{aa_i} &= a + \sum_{i \in I} \lambda_i (\overrightarrow{ab} + \overrightarrow{ba_i}) \\
 &= a + \left(\sum_{i \in I} \lambda_i \right) \overrightarrow{ab} + \sum_{i \in I} \lambda_i \overrightarrow{ba_i} \\
 &= a + \overrightarrow{ab} + \sum_{i \in I} \lambda_i \overrightarrow{ba_i} && \text{since } \sum_{i \in I} \lambda_i = 1 \\
 &= b + \sum_{i \in I} \lambda_i \overrightarrow{ba_i} && \text{since } b = a + \overrightarrow{ab}.
 \end{aligned}$$

(2) We also have

$$\begin{aligned}
 \sum_{i \in I} \lambda_i \overrightarrow{aa_i} &= \sum_{i \in I} \lambda_i (\overrightarrow{ab} + \overrightarrow{ba_i}) \\
 &= \left(\sum_{i \in I} \lambda_i \right) \overrightarrow{ab} + \sum_{i \in I} \lambda_i \overrightarrow{ba_i} \\
 &= \sum_{i \in I} \lambda_i \overrightarrow{ba_i},
 \end{aligned}$$

since $\sum_{i \in I} \lambda_i = 0$. □

Thus, by Lemma 19.1, for any family of points $(a_i)_{i \in I}$ in E , for any family $(\lambda_i)_{i \in I}$ of scalars such that $\sum_{i \in I} \lambda_i = 1$, the point

$$x = a + \sum_{i \in I} \lambda_i \overrightarrow{aa_i}$$

is independent of the choice of the origin $a \in E$. This property motivates the following definition.

Definition 19.2. For any family of points $(a_i)_{i \in I}$ in E , for any family $(\lambda_i)_{i \in I}$ of scalars such that $\sum_{i \in I} \lambda_i = 1$, and for any $a \in E$, the point

$$a + \sum_{i \in I} \lambda_i \overrightarrow{aa_i}$$

(which is independent of $a \in E$, by Lemma 19.1) is called the *barycenter* (or *barycentric combination*, or *affine combination*) of the points a_i assigned the weights λ_i , and it is denoted by

$$\sum_{i \in I} \lambda_i a_i.$$

In dealing with barycenters, it is convenient to introduce the notion of a *weighted point*, which is just a pair (a, λ) , where $a \in E$ is a point, and $\lambda \in \mathbb{R}$ is a scalar. Then, given a family of weighted points $((a_i, \lambda_i))_{i \in I}$, where $\sum_{i \in I} \lambda_i = 1$, we also say that the point $\sum_{i \in I} \lambda_i a_i$ is the *barycenter of the family of weighted points* $((a_i, \lambda_i))_{i \in I}$.

Note that the barycenter x of the family of weighted points $((a_i, \lambda_i))_{i \in I}$ is the unique point such that

$$\overrightarrow{ax} = \sum_{i \in I} \lambda_i \overrightarrow{aa_i} \quad \text{for every } a \in E,$$

and setting $a = x$, the point x is the unique point such that

$$\sum_{i \in I} \lambda_i \overrightarrow{xa_i} = 0.$$

In physical terms, the barycenter is the *center of mass* of the family of weighted points $((a_i, \lambda_i))_{i \in I}$ (where the masses have been normalized, so that $\sum_{i \in I} \lambda_i = 1$, and negative masses are allowed).

Remarks:

- (1) Since the barycenter of a family $((a_i, \lambda_i))_{i \in I}$ of weighted points is defined for families $(\lambda_i)_{i \in I}$ of scalars with finite support (and such that $\sum_{i \in I} \lambda_i = 1$), we might as well assume that I is finite. Then, for all $m \geq 2$, it is easy to prove that the barycenter of m weighted points can be obtained by repeated computations of barycenters of two weighted points.
- (2) This result still holds, provided that the field K has at least three distinct elements, but the proof is trickier!
- (3) When $\sum_{i \in I} \lambda_i = 0$, the vector $\sum_{i \in I} \lambda_i \overrightarrow{aa_i}$ does not depend on the point a , and we may denote it by $\sum_{i \in I} \lambda_i a_i$. This observation will be used to define a vector space in which linear combinations of both points and vectors make sense, regardless of the value of $\sum_{i \in I} \lambda_i$.

Figure 19.5 illustrates the geometric construction of the barycenters g_1 and g_2 of the weighted points $(a, \frac{1}{4})$, $(b, \frac{1}{4})$, and $(c, \frac{1}{2})$, and $(a, -1)$, $(b, 1)$, and $(c, 1)$.

The point g_1 can be constructed geometrically as the middle of the segment joining c to the middle $\frac{1}{2}a + \frac{1}{2}b$ of the segment (a, b) , since

$$g_1 = \frac{1}{2} \left(\frac{1}{2}a + \frac{1}{2}b \right) + \frac{1}{2}c.$$

The point g_2 can be constructed geometrically as the point such that the middle $\frac{1}{2}b + \frac{1}{2}c$ of the segment (b, c) is the middle of the segment (a, g_2) , since

$$g_2 = -a + 2 \left(\frac{1}{2}b + \frac{1}{2}c \right).$$

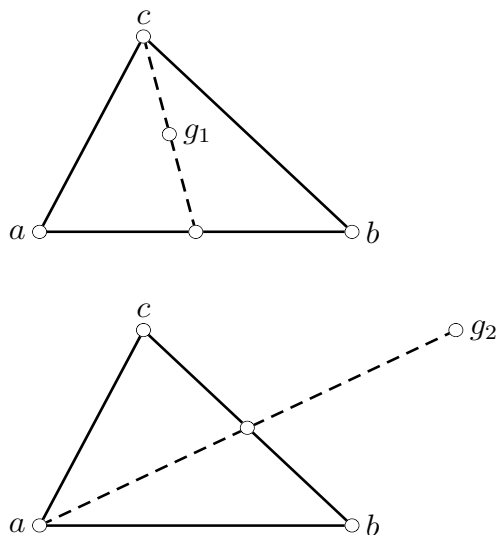


Figure 19.5: Barycenters, $g_1 = \frac{1}{4}a + \frac{1}{4}b + \frac{1}{2}c$, $g_2 = -a + b + c$

Later on, we will see that a polynomial curve can be defined as a set of barycenters of a fixed number of points. For example, let (a, b, c, d) be a sequence of points in \mathbb{A}^2 . Observe that

$$(1-t)^3 + 3t(1-t)^2 + 3t^2(1-t) + t^3 = 1,$$

since the sum on the left-hand side is obtained by expanding $(t + (1-t))^3 = 1$ using the binomial formula. Thus,

$$(1-t)^3 a + 3t(1-t)^2 b + 3t^2(1-t) c + t^3 d$$

is a well-defined affine combination. Then, we can define the curve $F: \mathbb{A} \rightarrow \mathbb{A}^2$ such that

$$F(t) = (1-t)^3 a + 3t(1-t)^2 b + 3t^2(1-t) c + t^3 d.$$

Such a curve is called a *Bézier curve*, and (a, b, c, d) are called its *control points*. Note that the curve passes through a and d , but generally not through b and c . It can be shown that any point $F(t)$ on the curve can be constructed using an algorithm performing affine interpolation steps (the *de Casteljau algorithm*).

19.5 Affine Subspaces

In linear algebra, a (linear) subspace can be characterized as a nonempty subset of a vector space closed under linear combinations. In affine spaces, the notion corresponding to the notion of (linear) subspace is the notion of affine subspace. It is natural to define an affine subspace as a subset of an affine space closed under affine combinations.

Definition 19.3. Given an affine space $\langle E, \vec{E}, + \rangle$, a subset V of E is an *affine subspace* (of $\langle E, \vec{E}, + \rangle$) if for every family of weighted points $((a_i, \lambda_i))_{i \in I}$ in V such that $\sum_{i \in I} \lambda_i = 1$, the barycenter $\sum_{i \in I} \lambda_i a_i$ belongs to V .

An affine subspace is also called a *flat* by some authors. According to Definition 19.3, the empty set is trivially an affine subspace, and every intersection of affine subspaces is an affine subspace.

As an example, consider the subset U of \mathbb{R}^2 defined by

$$U = \{(x, y) \in \mathbb{R}^2 \mid ax + by = c\},$$

i.e., the set of solutions of the equation

$$ax + by = c,$$

where it is assumed that $a \neq 0$ or $b \neq 0$. Given any m points $(x_i, y_i) \in U$ and any m scalars λ_i such that $\lambda_1 + \cdots + \lambda_m = 1$, we claim that

$$\sum_{i=1}^m \lambda_i (x_i, y_i) \in U.$$

Indeed, $(x_i, y_i) \in U$ means that

$$ax_i + by_i = c,$$

and if we multiply both sides of this equation by λ_i and add up the resulting m equations, we get

$$\sum_{i=1}^m (\lambda_i ax_i + \lambda_i by_i) = \sum_{i=1}^m \lambda_i c,$$

and since $\lambda_1 + \cdots + \lambda_m = 1$, we get

$$a \left(\sum_{i=1}^m \lambda_i x_i \right) + b \left(\sum_{i=1}^m \lambda_i y_i \right) = \left(\sum_{i=1}^m \lambda_i \right) c = c,$$

which shows that

$$\left(\sum_{i=1}^m \lambda_i x_i, \sum_{i=1}^m \lambda_i y_i \right) = \sum_{i=1}^m \lambda_i (x_i, y_i) \in U.$$

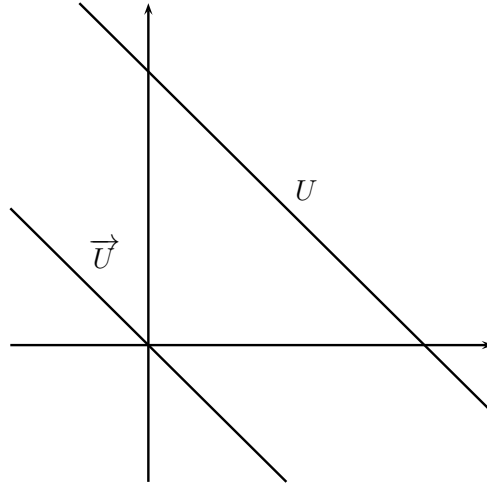
Thus, U is an affine subspace of \mathbb{A}^2 . In fact, it is just a usual line in \mathbb{A}^2 .

It turns out that U is closely related to the subset of \mathbb{R}^2 defined by

$$\vec{U} = \{(x, y) \in \mathbb{R}^2 \mid ax + by = 0\},$$

i.e., the set of solutions of the homogeneous equation

$$ax + by = 0$$

Figure 19.6: An affine line U and its direction

obtained by setting the right-hand side of $ax + by = c$ to zero. Indeed, for any m scalars λ_i , the same calculation as above yields that

$$\sum_{i=1}^m \lambda_i(x_i, y_i) \in \vec{U},$$

this time **without any restriction on the** λ_i , since the right-hand side of the equation is null. Thus, \vec{U} is a subspace of \mathbb{R}^2 . In fact, \vec{U} is one-dimensional, and it is just a usual line in \mathbb{R}^2 . This line can be identified with a line passing through the origin of \mathbb{A}^2 , a line that is parallel to the line U of equation $ax + by = c$, as illustrated in Figure 19.6.

Now, if (x_0, y_0) is any point in U , we claim that

$$U = (x_0, y_0) + \vec{U},$$

where

$$(x_0, y_0) + \vec{U} = \{(x_0 + u_1, y_0 + u_2) \mid (u_1, u_2) \in \vec{U}\}.$$

First, $(x_0, y_0) + \vec{U} \subseteq U$, since $ax_0 + by_0 = c$ and $au_1 + bu_2 = 0$ for all $(u_1, u_2) \in \vec{U}$. Second, if $(x, y) \in U$, then $ax + by = c$, and since we also have $ax_0 + by_0 = c$, by subtraction, we get

$$a(x - x_0) + b(y - y_0) = 0,$$

which shows that $(x - x_0, y - y_0) \in \vec{U}$, and thus $(x, y) \in (x_0, y_0) + \vec{U}$. Hence, we also have $U \subseteq (x_0, y_0) + \vec{U}$, and $U = (x_0, y_0) + \vec{U}$.

The above example shows that the affine line U defined by the equation

$$ax + by = c$$

is obtained by “translating” the parallel line \vec{U} of equation

$$ax + by = 0$$

passing through the origin. In fact, given any point $(x_0, y_0) \in U$,

$$U = (x_0, y_0) + \vec{U}.$$

More generally, it is easy to prove the following fact. Given any $m \times n$ matrix A and any vector $b \in \mathbb{R}^m$, the subset U of \mathbb{R}^n defined by

$$U = \{x \in \mathbb{R}^n \mid Ax = b\}$$

is an affine subspace of \mathbb{A}^n .

Actually, observe that $Ax = b$ should really be written as $Ax^\top = b$, to be consistent with our convention that points are represented by row vectors. We can also use the boldface notation for column vectors, in which case the equation is written as $A\mathbf{x} = b$. For the sake of minimizing the amount of notation, we stick to the simpler (yet incorrect) notation $Ax = b$. If we consider the corresponding homogeneous equation $Ax = 0$, the set

$$\vec{U} = \{x \in \mathbb{R}^n \mid Ax = 0\}$$

is a subspace of \mathbb{R}^n , and for any $x_0 \in U$, we have

$$U = x_0 + \vec{U}.$$

This is a general situation. Affine subspaces can be characterized in terms of subspaces of \vec{E} . Let V be a nonempty subset of E . For every family (a_1, \dots, a_n) in V , for any family $(\lambda_1, \dots, \lambda_n)$ of scalars, and for every point $a \in V$, observe that for every $x \in E$,

$$x = a + \sum_{i=1}^n \lambda_i \vec{aa_i}$$

is the barycenter of the family of weighted points

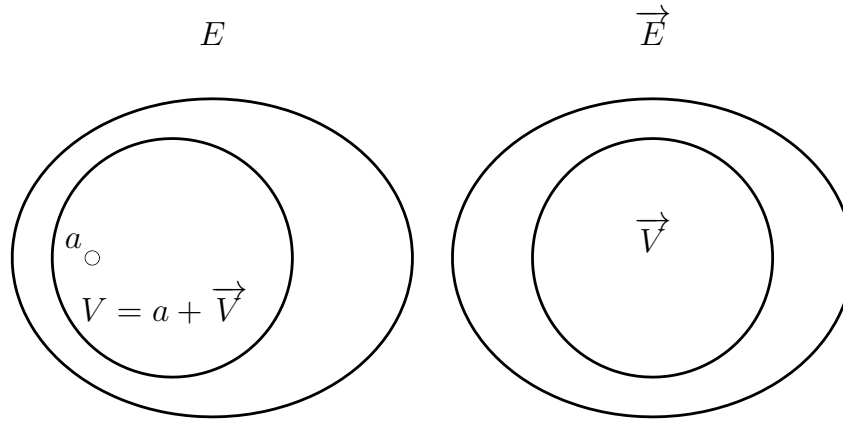
$$\left((a_1, \lambda_1), \dots, (a_n, \lambda_n), \left(a, 1 - \sum_{i=1}^n \lambda_i \right) \right),$$

since

$$\sum_{i=1}^n \lambda_i + \left(1 - \sum_{i=1}^n \lambda_i \right) = 1.$$

Given any point $a \in E$ and any subset \vec{V} of \vec{E} , let $a + \vec{V}$ denote the following subset of E :

$$a + \vec{V} = \{a + v \mid v \in \vec{V}\}.$$

Figure 19.7: An affine subspace V and its direction \vec{V}

Lemma 19.2. Let $\langle E, \vec{E}, + \rangle$ be an affine space.

(1) A nonempty subset V of E is an affine subspace iff for every point $a \in V$, the set

$$\vec{V}_a = \{\vec{ax} \mid x \in V\}$$

is a subspace of \vec{E} . Consequently, $V = a + \vec{V}_a$. Furthermore,

$$\vec{V} = \{\vec{xy} \mid x, y \in V\}$$

is a subspace of \vec{E} and $\vec{V}_a = \vec{V}$ for all $a \in E$. Thus, $V = a + \vec{V}$.

(2) For any subspace \vec{V} of \vec{E} and for any $a \in E$, the set $V = a + \vec{V}$ is an affine subspace.

Proof. The proof is straightforward, and is omitted. It is also given in Gallier [43]. \square

In particular, when E is the natural affine space associated with a vector space \vec{E} , Lemma 19.2 shows that every affine subspace of E is of the form $u + \vec{U}$, for a subspace \vec{U} of \vec{E} . The subspaces of \vec{E} are the affine subspaces of E that contain 0.

The subspace \vec{V} associated with an affine subspace V is called the *direction of V* . It is also clear that the map $+: V \times \vec{V} \rightarrow V$ induced by $+: E \times \vec{E} \rightarrow E$ confers to $\langle V, \vec{V}, + \rangle$ an affine structure. Figure 19.7 illustrates the notion of affine subspace.

By the dimension of the subspace V , we mean the dimension of \vec{V} .

An affine subspace of dimension 1 is called a *line*, and an affine subspace of dimension 2 is called a *plane*.

An affine subspace of codimension 1 is called a *hyperplane* (recall that a subspace F of a vector space E has codimension 1 iff there is some subspace G of dimension 1 such that $E = F \oplus G$, the direct sum of F and G , see Strang [105] or Lang [67]).

We say that two affine subspaces U and V are *parallel* if their directions are identical. Equivalently, since $\vec{U} = \vec{V}$, we have $U = a + \vec{U}$ and $V = b + \vec{U}$ for any $a \in U$ and any $b \in V$, and thus V is obtained from U by the translation \vec{ab} .

In general, when we talk about n points a_1, \dots, a_n , we mean the sequence (a_1, \dots, a_n) , and not the set $\{a_1, \dots, a_n\}$ (the a_i 's need not be distinct).

By Lemma 19.2, a line is specified by a point $a \in E$ and a nonzero vector $v \in \vec{E}$, i.e., a line is the set of all points of the form $a + \lambda v$, for $\lambda \in \mathbb{R}$.

We say that three points a, b, c are *collinear* if the vectors \vec{ab} and \vec{ac} are linearly dependent. If two of the points a, b, c are distinct, say $a \neq b$, then there is a unique $\lambda \in \mathbb{R}$ such that $\vec{ac} = \lambda \vec{ab}$, and we define the ratio $\frac{\vec{ac}}{\vec{ab}} = \lambda$.

A plane is specified by a point $a \in E$ and two linearly independent vectors $u, v \in \vec{E}$, i.e., a plane is the set of all points of the form $a + \lambda u + \mu v$, for $\lambda, \mu \in \mathbb{R}$.

We say that four points a, b, c, d are *coplanar* if the vectors \vec{ab}, \vec{ac} , and \vec{ad} are linearly dependent. Hyperplanes will be characterized a little later.

Lemma 19.3. *Given an affine space $\langle E, \vec{E}, + \rangle$, for any family $(a_i)_{i \in I}$ of points in E , the set V of barycenters $\sum_{i \in I} \lambda_i a_i$ (where $\sum_{i \in I} \lambda_i = 1$) is the smallest affine subspace containing $(a_i)_{i \in I}$.*

Proof. If $(a_i)_{i \in I}$ is empty, then $V = \emptyset$, because of the condition $\sum_{i \in I} \lambda_i = 1$. If $(a_i)_{i \in I}$ is nonempty, then the smallest affine subspace containing $(a_i)_{i \in I}$ must contain the set V of barycenters $\sum_{i \in I} \lambda_i a_i$, and thus, it is enough to show that V is closed under affine combinations, which is immediately verified. \square

Given a nonempty subset S of E , the smallest affine subspace of E generated by S is often denoted by $\langle S \rangle$. For example, a line specified by two distinct points a and b is denoted by $\langle a, b \rangle$, or even (a, b) , and similarly for planes, etc.

Remarks:

- (1) Since it can be shown that the barycenter of n weighted points can be obtained by repeated computations of barycenters of two weighted points, a nonempty subset V of E is an affine subspace iff for every two points $a, b \in V$, the set V contains all barycentric combinations of a and b . If V contains at least two points, then V is an affine subspace iff for any two distinct points $a, b \in V$, the set V contains the line determined by a and b , that is, the set of all points $(1 - \lambda)a + \lambda b$, $\lambda \in \mathbb{R}$.
- (2) This result still holds if the field K has at least three distinct elements, but the proof is trickier!

19.6 Affine Independence and Affine Frames

Corresponding to the notion of linear independence in vector spaces, we have the notion of affine independence. Given a family $(a_i)_{i \in I}$ of points in an affine space E , we will reduce the

notion of (affine) independence of these points to the (linear) independence of the families $(\overrightarrow{a_i a_j})_{j \in (I - \{i\})}$ of vectors obtained by choosing any a_i as an origin. First, the following lemma shows that it is sufficient to consider only one of these families.

Lemma 19.4. *Given an affine space $\langle E, \overrightarrow{E}, + \rangle$, let $(a_i)_{i \in I}$ be a family of points in E . If the family $(\overrightarrow{a_i a_j})_{j \in (I - \{i\})}$ is linearly independent for some $i \in I$, then $(\overrightarrow{a_i a_j})_{j \in (I - \{i\})}$ is linearly independent for every $i \in I$.*

Proof. Assume that the family $(\overrightarrow{a_i a_j})_{j \in (I - \{i\})}$ is linearly independent for some specific $i \in I$. Let $k \in I$ with $k \neq i$, and assume that there are some scalars $(\lambda_j)_{j \in (I - \{k\})}$ such that

$$\sum_{j \in (I - \{k\})} \lambda_j \overrightarrow{a_k a_j} = 0.$$

Since

$$\overrightarrow{a_k a_j} = \overrightarrow{a_k a_i} + \overrightarrow{a_i a_j},$$

we have

$$\begin{aligned} \sum_{j \in (I - \{k\})} \lambda_j \overrightarrow{a_k a_j} &= \sum_{j \in (I - \{k\})} \lambda_j \overrightarrow{a_k a_i} + \sum_{j \in (I - \{k\})} \lambda_j \overrightarrow{a_i a_j}, \\ &= \sum_{j \in (I - \{k\})} \lambda_j \overrightarrow{a_k a_i} + \sum_{j \in (I - \{i, k\})} \lambda_j \overrightarrow{a_i a_j}, \\ &= \sum_{j \in (I - \{i, k\})} \lambda_j \overrightarrow{a_i a_j} - \left(\sum_{j \in (I - \{k\})} \lambda_j \right) \overrightarrow{a_i a_k}, \end{aligned}$$

and thus

$$\sum_{j \in (I - \{i, k\})} \lambda_j \overrightarrow{a_i a_j} - \left(\sum_{j \in (I - \{k\})} \lambda_j \right) \overrightarrow{a_i a_k} = 0.$$

Since the family $(\overrightarrow{a_i a_j})_{j \in (I - \{i\})}$ is linearly independent, we must have $\lambda_j = 0$ for all $j \in (I - \{i, k\})$ and $\sum_{j \in (I - \{k\})} \lambda_j = 0$, which implies that $\lambda_j = 0$ for all $j \in (I - \{k\})$. \square

We define affine independence as follows.

Definition 19.4. Given an affine space $\langle E, \overrightarrow{E}, + \rangle$, a family $(a_i)_{i \in I}$ of points in E is *affinely independent* if the family $(\overrightarrow{a_i a_j})_{j \in (I - \{i\})}$ is linearly independent for some $i \in I$.

Definition 19.4 is reasonable, because by Lemma 19.4, the independence of the family $(\overrightarrow{a_i a_j})_{j \in (I - \{i\})}$ does not depend on the choice of a_i . A crucial property of linearly independent vectors (u_1, \dots, u_m) is that if a vector v is a linear combination

$$v = \sum_{i=1}^m \lambda_i u_i$$

of the u_i , then the λ_i are unique. A similar result holds for affinely independent points.

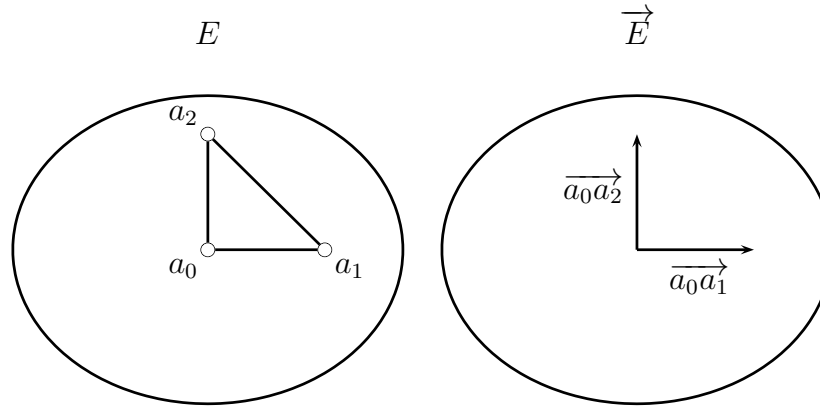


Figure 19.8: Affine independence and linear independence

Lemma 19.5. *Given an affine space $\langle E, \vec{E}, + \rangle$, let (a_0, \dots, a_m) be a family of $m+1$ points in E . Let $x \in E$, and assume that $x = \sum_{i=0}^m \lambda_i a_i$, where $\sum_{i=0}^m \lambda_i = 1$. Then, the family $(\lambda_0, \dots, \lambda_m)$ such that $x = \sum_{i=0}^m \lambda_i a_i$ is unique iff the family $(\overrightarrow{a_0a_1}, \dots, \overrightarrow{a_0a_m})$ is linearly independent.*

Proof. The proof is straightforward and is omitted. It is also given in Gallier [43]. \square

Lemma 19.5 suggests the notion of affine frame. Affine frames are the affine analogues of bases in vector spaces. Let $\langle E, \vec{E}, + \rangle$ be a nonempty affine space, and let (a_0, \dots, a_m) be a family of $m+1$ points in E . The family (a_0, \dots, a_m) determines the family of m vectors $(\overrightarrow{a_0a_1}, \dots, \overrightarrow{a_0a_m})$ in \vec{E} . Conversely, given a point a_0 in E and a family of m vectors (u_1, \dots, u_m) in \vec{E} , we obtain the family of $m+1$ points (a_0, \dots, a_m) in E , where $a_i = a_0 + u_i$, $1 \leq i \leq m$.

Thus, for any $m \geq 1$, it is equivalent to consider a family of $m+1$ points (a_0, \dots, a_m) in E , and a pair $(a_0, (u_1, \dots, u_m))$, where the u_i are vectors in \vec{E} . Figure 19.8 illustrates the notion of affine independence.

Remark: The above observation also applies to infinite families $(a_i)_{i \in I}$ of points in E and families $(u_i)_{i \in I - \{0\}}$ of vectors in \vec{E} , provided that the index set I contains 0.

When $(\overrightarrow{a_0a_1}, \dots, \overrightarrow{a_0a_m})$ is a basis of \vec{E} then, for every $x \in E$, since $x = a_0 + \overrightarrow{a_0x}$, there is a unique family (x_1, \dots, x_m) of scalars such that

$$x = a_0 + x_1 \overrightarrow{a_0a_1} + \dots + x_m \overrightarrow{a_0a_m}.$$

The scalars (x_1, \dots, x_m) may be considered as coordinates with respect to $(a_0, (\overrightarrow{a_0a_1}, \dots, \overrightarrow{a_0a_m}))$. Since

$$x = a_0 + \sum_{i=1}^m x_i \overrightarrow{a_0a_i} \quad \text{iff} \quad x = \left(1 - \sum_{i=1}^m x_i\right) a_0 + \sum_{i=1}^m x_i a_i,$$

$x \in E$ can also be expressed uniquely as

$$x = \sum_{i=0}^m \lambda_i a_i$$

with $\sum_{i=0}^m \lambda_i = 1$, and where $\lambda_0 = 1 - \sum_{i=1}^m \lambda_i$, and $\lambda_i = x_i$ for $1 \leq i \leq m$. The scalars $(\lambda_0, \dots, \lambda_m)$ are also certain kinds of coordinates with respect to (a_0, \dots, a_m) . All this is summarized in the following definition.

Definition 19.5. Given an affine space $\langle E, \vec{E}, + \rangle$, an *affine frame with origin* a_0 is a family (a_0, \dots, a_m) of $m+1$ points in E such that the list of vectors $(\overrightarrow{a_0 a_1}, \dots, \overrightarrow{a_0 a_m})$ is a basis of \vec{E} . The pair $(a_0, (\overrightarrow{a_0 a_1}, \dots, \overrightarrow{a_0 a_m}))$ is also called an *affine frame with origin* a_0 . Then, every $x \in E$ can be expressed as

$$x = a_0 + x_1 \overrightarrow{a_0 a_1} + \dots + x_m \overrightarrow{a_0 a_m}$$

for a unique family (x_1, \dots, x_m) of scalars, called the *coordinates of x w.r.t. the affine frame* $(a_0, (\overrightarrow{a_0 a_1}, \dots, \overrightarrow{a_0 a_m}))$. Furthermore, every $x \in E$ can be written as

$$x = \lambda_0 a_0 + \dots + \lambda_m a_m$$

for some unique family $(\lambda_0, \dots, \lambda_m)$ of scalars such that $\lambda_0 + \dots + \lambda_m = 1$ called the *barycentric coordinates of x with respect to the affine frame* (a_0, \dots, a_m) .

The coordinates (x_1, \dots, x_m) and the barycentric coordinates $(\lambda_0, \dots, \lambda_m)$ are related by the equations $\lambda_0 = 1 - \sum_{i=1}^m x_i$ and $\lambda_i = x_i$, for $1 \leq i \leq m$. An affine frame is called an *affine basis* by some authors. A family $(a_i)_{i \in I}$ of points in E is *affinely dependent* if it is not affinely independent. We can also characterize affinely dependent families as follows.

Lemma 19.6. Given an affine space $\langle E, \vec{E}, + \rangle$, let $(a_i)_{i \in I}$ be a family of points in E . The family $(a_i)_{i \in I}$ is affinely dependent iff there is a family $(\lambda_i)_{i \in I}$ such that $\lambda_j \neq 0$ for some $j \in I$, $\sum_{i \in I} \lambda_i = 0$, and $\sum_{i \in I} \lambda_i \overrightarrow{x a_i} = 0$ for every $x \in E$.

Proof. By Lemma 19.5, the family $(a_i)_{i \in I}$ is affinely dependent iff the family of vectors $(\overrightarrow{a_i a_j})_{j \in (I - \{i\})}$ is linearly dependent for some $i \in I$. For any $i \in I$, the family $(\overrightarrow{a_i a_j})_{j \in (I - \{i\})}$ is linearly dependent iff there is a family $(\lambda_j)_{j \in (I - \{i\})}$ such that $\lambda_j \neq 0$ for some j , and such that

$$\sum_{j \in (I - \{i\})} \lambda_j \overrightarrow{a_i a_j} = 0.$$

Then, for any $x \in E$, we have

$$\begin{aligned} \sum_{j \in (I - \{i\})} \lambda_j \overrightarrow{a_i a_j} &= \sum_{j \in (I - \{i\})} \lambda_j (\overrightarrow{x a_j} - \overrightarrow{x a_i}) \\ &= \sum_{j \in (I - \{i\})} \lambda_j \overrightarrow{x a_j} - \left(\sum_{j \in (I - \{i\})} \lambda_j \right) \overrightarrow{x a_i}, \end{aligned}$$

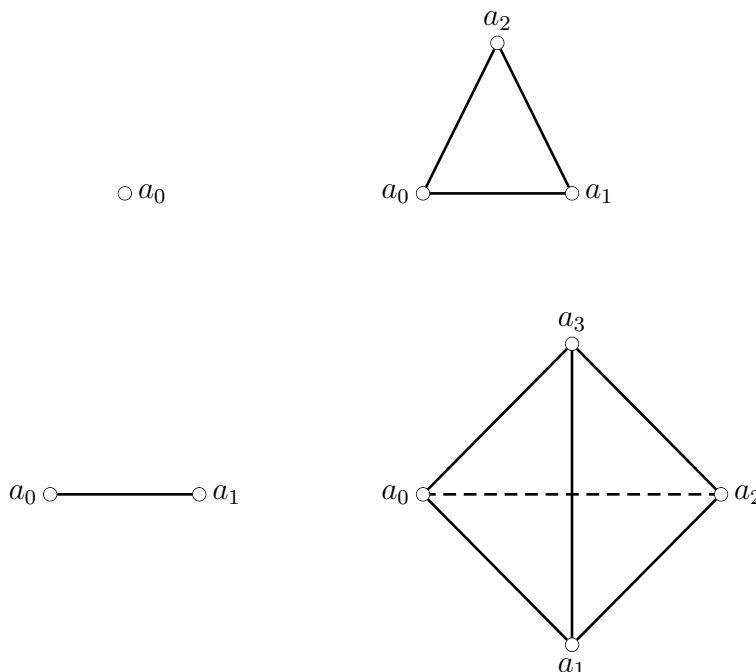


Figure 19.9: Examples of affine frames and their convex hulls

and letting $\lambda_i = -\left(\sum_{j \in (I - \{i\})} \lambda_j\right)$, we get $\sum_{i \in I} \lambda_i \vec{x} \vec{a}_i = 0$, with $\sum_{i \in I} \lambda_i = 0$ and $\lambda_j \neq 0$ for some $j \in I$. The converse is obvious by setting $x = a_i$ for some i such that $\lambda_i \neq 0$, since $\sum_{i \in I} \lambda_i = 0$ implies that $\lambda_j \neq 0$, for some $j \neq i$. \square

Even though Lemma 19.6 is rather dull, it is one of the key ingredients in the proof of beautiful and deep theorems about convex sets, such as Carathéodory's theorem, Radon's theorem, and Helly's theorem.

A family of two points (a, b) in E is affinely independent iff $\vec{ab} \neq 0$, iff $a \neq b$. If $a \neq b$, the affine subspace generated by a and b is the set of all points $(1 - \lambda)a + \lambda b$, which is the unique line passing through a and b . A family of three points (a, b, c) in E is affinely independent iff \vec{ab} and \vec{ac} are linearly independent, which means that a , b , and c are not on the same line (they are not collinear). In this case, the affine subspace generated by (a, b, c) is the set of all points $(1 - \lambda - \mu)a + \lambda b + \mu c$, which is the unique plane containing a , b , and c . A family of four points (a, b, c, d) in E is affinely independent iff \vec{ab} , \vec{ac} , and \vec{ad} are linearly independent, which means that a , b , c , and d are not in the same plane (they are not coplanar). In this case, a , b , c , and d are the vertices of a tetrahedron. Figure 19.9 shows affine frames and their convex hulls for $|I| = 0, 1, 2, 3$.

Given $n+1$ affinely independent points (a_0, \dots, a_n) in E , we can consider the set of points $\lambda_0 a_0 + \dots + \lambda_n a_n$, where $\lambda_0 + \dots + \lambda_n = 1$ and $\lambda_i \geq 0$ ($\lambda_i \in \mathbb{R}$). Such affine combinations are called *convex combinations*. This set is called the *convex hull* of (a_0, \dots, a_n) (or *n-simplex*

spanned by (a_0, \dots, a_n) . When $n = 1$, we get the segment between a_0 and a_1 , including a_0 and a_1 . When $n = 2$, we get the interior of the triangle whose vertices are a_0, a_1, a_2 , including boundary points (the edges). When $n = 3$, we get the interior of the tetrahedron whose vertices are a_0, a_1, a_2, a_3 , including boundary points (faces and edges). The set

$$\{a_0 + \lambda_1 \overrightarrow{a_0 a_1} + \dots + \lambda_n \overrightarrow{a_0 a_n} \mid \text{where } 0 \leq \lambda_i \leq 1 \ (\lambda_i \in \mathbb{R})\}$$

is called the *parallelotope spanned by (a_0, \dots, a_n)* . When E has dimension 2, a parallelotope is also called a *parallelogram*, and when E has dimension 3, a *parallelepiped*.

More generally, we say that a subset V of E is *convex* if for any two points $a, b \in V$, we have $c \in V$ for every point $c = (1 - \lambda)a + \lambda b$, with $0 \leq \lambda \leq 1$ ($\lambda \in \mathbb{R}$).



Points are not vectors! The following example illustrates why treating points as vectors may cause problems. Let a, b, c be three affinely independent points in \mathbb{A}^3 . Any point x in the plane (a, b, c) can be expressed as

$$x = \lambda_0 a + \lambda_1 b + \lambda_2 c,$$

where $\lambda_0 + \lambda_1 + \lambda_2 = 1$. How can we compute $\lambda_0, \lambda_1, \lambda_2$? Letting $a = (a_1, a_2, a_3)$, $b = (b_1, b_2, b_3)$, $c = (c_1, c_2, c_3)$, and $x = (x_1, x_2, x_3)$ be the coordinates of a, b, c, x in the standard frame of \mathbb{A}^3 , it is tempting to solve the system of equations

$$\begin{pmatrix} a_1 & b_1 & c_1 \\ a_2 & b_2 & c_2 \\ a_3 & b_3 & c_3 \end{pmatrix} \begin{pmatrix} \lambda_0 \\ \lambda_1 \\ \lambda_2 \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}.$$

However, there is a problem when the origin of the coordinate system belongs to the plane (a, b, c) , since in this case, the matrix is not invertible! What we should really be doing is to solve the system

$$\lambda_0 \overrightarrow{Oa} + \lambda_1 \overrightarrow{Ob} + \lambda_2 \overrightarrow{Oc} = \overrightarrow{Ox},$$

where O is any point **not** in the plane (a, b, c) . An alternative is to use certain well-chosen cross products.

It can be shown that barycentric coordinates correspond to various ratios of areas and volumes; see the problems.

19.7 Affine Maps

Corresponding to linear maps we have the notion of an affine map. An affine map is defined as a map preserving affine combinations.

Definition 19.6. Given two affine spaces $\langle E, \overrightarrow{E}, + \rangle$ and $\langle E', \overrightarrow{E'}, +' \rangle$, a function $f: E \rightarrow E'$ is an *affine map* iff for every family $((a_i, \lambda_i))_{i \in I}$ of weighted points in E such that $\sum_{i \in I} \lambda_i = 1$, we have

$$f\left(\sum_{i \in I} \lambda_i a_i\right) = \sum_{i \in I} \lambda_i f(a_i).$$

In other words, f preserves barycenters.

Affine maps can be obtained from linear maps as follows. For simplicity of notation, the same symbol $+$ is used for both affine spaces (instead of using both $+$ and $+'$).

Given any point $a \in E$, any point $b \in E'$, and any linear map $h: \vec{E} \rightarrow \vec{E}'$, we claim that the map $f: E \rightarrow E'$ defined such that

$$f(a + v) = b + h(v)$$

is an affine map. Indeed, for any family $(\lambda_i)_{i \in I}$ of scalars with $\sum_{i \in I} \lambda_i = 1$ and any family $(v_i)_{i \in I}$, since

$$\sum_{i \in I} \lambda_i(a + v_i) = a + \sum_{i \in I} \lambda_i \overrightarrow{a(a + v_i)} = a + \sum_{i \in I} \lambda_i v_i$$

and

$$\sum_{i \in I} \lambda_i(b + h(v_i)) = b + \sum_{i \in I} \lambda_i \overrightarrow{b(b + h(v_i))} = b + \sum_{i \in I} \lambda_i h(v_i),$$

we have

$$\begin{aligned} f\left(\sum_{i \in I} \lambda_i(a + v_i)\right) &= f\left(a + \sum_{i \in I} \lambda_i v_i\right) \\ &= b + h\left(\sum_{i \in I} \lambda_i v_i\right) \\ &= b + \sum_{i \in I} \lambda_i h(v_i) \\ &= \sum_{i \in I} \lambda_i(b + h(v_i)) \\ &= \sum_{i \in I} \lambda_i f(a + v_i). \end{aligned}$$

Note that the condition $\sum_{i \in I} \lambda_i = 1$ was implicitly used (in a hidden call to Lemma 19.1) in deriving that

$$\sum_{i \in I} \lambda_i(a + v_i) = a + \sum_{i \in I} \lambda_i v_i$$

and

$$\sum_{i \in I} \lambda_i(b + h(v_i)) = b + \sum_{i \in I} \lambda_i h(v_i).$$

As a more concrete example, the map

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} 3 \\ 1 \end{pmatrix}$$

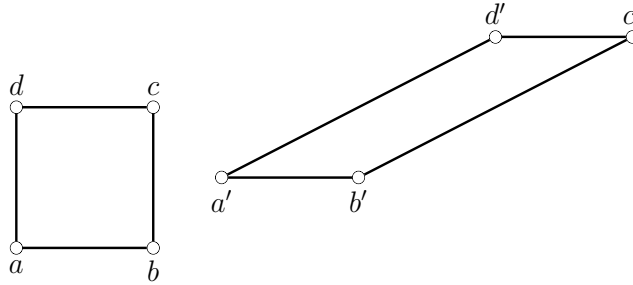


Figure 19.10: The effect of a shear

defines an affine map in \mathbb{A}^2 . It is a “shear” followed by a translation. The effect of this shear on the square (a, b, c, d) is shown in Figure 19.10. The image of the square (a, b, c, d) is the parallelogram (a', b', c', d') .

Let us consider one more example. The map

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} 1 & 1 \\ 1 & 3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} 3 \\ 0 \end{pmatrix}$$

is an affine map. Since we can write

$$\begin{pmatrix} 1 & 1 \\ 1 & 3 \end{pmatrix} = \sqrt{2} \begin{pmatrix} \sqrt{2}/2 & -\sqrt{2}/2 \\ 2/2 & \sqrt{2}/2 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix},$$

this affine map is the composition of a shear, followed by a rotation of angle $\pi/4$, followed by a magnification of ratio $\sqrt{2}$, followed by a translation. The effect of this map on the square (a, b, c, d) is shown in Figure 19.11. The image of the square (a, b, c, d) is the parallelogram (a', b', c', d') .

The following lemma shows the converse of what we just showed. Every affine map is determined by the image of any point and a linear map.

Lemma 19.7. *Given an affine map $f: E \rightarrow E'$, there is a unique linear map $\vec{f}: \vec{E} \rightarrow \vec{E}'$ such that*

$$f(a + v) = f(a) + \vec{f}(v),$$

for every $a \in E$ and every $v \in \vec{E}$.

Proof. Let $a \in E$ be any point in E . We claim that the map defined such that

$$\vec{f}(v) = \overrightarrow{f(a)f(a+v)}$$

for every $v \in \vec{E}$ is a linear map $\vec{f}: \vec{E} \rightarrow \vec{E}'$. Indeed, we can write

$$a + \lambda v = \lambda(a + v) + (1 - \lambda)a,$$

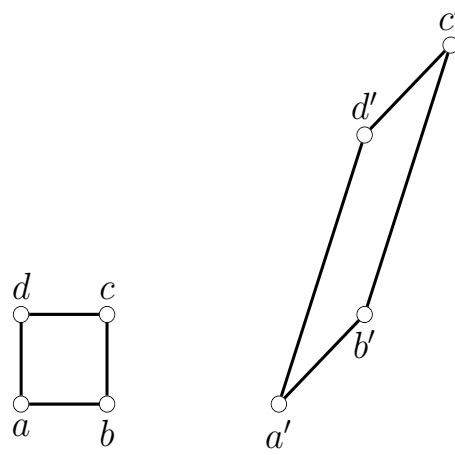


Figure 19.11: The effect of an affine map

since $a + \lambda v = a + \lambda \overrightarrow{a(a+v)} + (1-\lambda)\overrightarrow{a\vec{a}}$, and also

$$a + u + v = (a + u) + (a + v) - a,$$

since $a + u + v = a + \overrightarrow{a(a+u)} + \overrightarrow{a(a+v)} - \overrightarrow{a\vec{a}}$. Since f preserves barycenters, we get

$$f(a + \lambda v) = \lambda f(a + v) + (1 - \lambda)f(a).$$

If we recall that $x = \sum_{i \in I} \lambda_i a_i$ is the barycenter of a family $((a_i, \lambda_i))_{i \in I}$ of weighted points (with $\sum_{i \in I} \lambda_i = 1$) iff

$$\overrightarrow{bx} = \sum_{i \in I} \lambda_i \overrightarrow{ba_i} \quad \text{for every } b \in E,$$

we get

$$\overrightarrow{f(a)f(a+\lambda v)} = \lambda \overrightarrow{f(a)f(a+v)} + (1-\lambda)\overrightarrow{f(a)f(a)} = \lambda \overrightarrow{f(a)f(a+v)},$$

showing that $\overrightarrow{f}(\lambda v) = \lambda \overrightarrow{f}(v)$. We also have

$$f(a + u + v) = f(a + u) + f(a + v) - f(a),$$

from which we get

$$\overrightarrow{f(a)f(a+u+v)} = \overrightarrow{f(a)f(a+u)} + \overrightarrow{f(a)f(a+v)},$$

showing that $\overrightarrow{f}(u + v) = \overrightarrow{f}(u) + \overrightarrow{f}(v)$. Consequently, \overrightarrow{f} is a linear map. For any other point $b \in E$, since

$$b + v = a + \overrightarrow{ab} + v = a + \overrightarrow{a(a+v)} - \overrightarrow{a\vec{a}} + \overrightarrow{ab},$$

$b + v = (a + v) - a + b$, and since f preserves barycenters, we get

$$f(b + v) = f(a + v) - f(a) + f(b),$$

which implies that

$$\begin{aligned} \overrightarrow{f(b)f(b+v)} &= \overrightarrow{f(b)f(a+v)} - \overrightarrow{f(b)f(a)} + \overrightarrow{f(b)f(b)}, \\ &= \overrightarrow{f(a)f(b)} + \overrightarrow{f(b)f(a+v)}, \\ &= \overrightarrow{f(a)f(a+v)}. \end{aligned}$$

Thus, $\overrightarrow{f(b)f(b+v)} = \overrightarrow{f(a)f(a+v)}$, which shows that the definition of \vec{f} does not depend on the choice of $a \in E$. The fact that \vec{f} is unique is obvious: We must have $\vec{f}(v) = \overrightarrow{f(a)f(a+v)}$. \square

The unique linear map $\vec{f}: \vec{E} \rightarrow \vec{E}'$ given by Lemma 19.7 is called the *linear map associated with the affine map f* .

Note that the condition

$$f(a + v) = f(a) + \vec{f}(v),$$

for every $a \in E$ and every $v \in \vec{E}$, can be stated equivalently as

$$f(x) = f(a) + \vec{f}(\overrightarrow{a\vec{x}}), \quad \text{or} \quad \overrightarrow{f(a)f(x)} = \vec{f}(\overrightarrow{a\vec{x}}),$$

for all $a, x \in E$. Lemma 19.7 shows that for any affine map $f: E \rightarrow E'$, there are points $a \in E$, $b \in E'$, and a unique linear map $\vec{f}: \vec{E} \rightarrow \vec{E}'$, such that

$$f(a + v) = b + \vec{f}(v),$$

for all $v \in \vec{E}$ (just let $b = f(a)$, for any $a \in E$). Affine maps for which \vec{f} is the identity map are called *translations*. Indeed, if $\vec{f} = \text{id}$,

$$\begin{aligned} f(x) &= f(a) + \vec{f}(\overrightarrow{a\vec{x}}) = f(a) + \overrightarrow{a\vec{x}} = x + \overrightarrow{x\vec{a}} + \overrightarrow{af(a)} + \overrightarrow{a\vec{x}} \\ &= x + \overrightarrow{x\vec{a}} + \overrightarrow{af(a)} - \overrightarrow{x\vec{a}} = x + \overrightarrow{af(a)}, \end{aligned}$$

and so

$$\overrightarrow{xf(x)} = \overrightarrow{af(a)},$$

which shows that f is the translation induced by the vector $\overrightarrow{af(a)}$ (which does not depend on a).

Since an affine map preserves barycenters, and since an affine subspace V is closed under barycentric combinations, the image $f(V)$ of V is an affine subspace in E' . So, for example, the image of a line is a point or a line, and the image of a plane is either a point, a line, or a plane.

It is easily verified that the composition of two affine maps is an affine map. Also, given affine maps $f: E \rightarrow E'$ and $g: E' \rightarrow E''$, we have

$$g(f(a + v)) = g\left(f(a) + \overrightarrow{f}(v)\right) = g(f(a)) + \overrightarrow{g}\left(\overrightarrow{f}(v)\right),$$

which shows that $\overrightarrow{g \circ f} = \overrightarrow{g} \circ \overrightarrow{f}$. It is easy to show that an affine map $f: E \rightarrow E'$ is injective iff $\overrightarrow{f}: \overrightarrow{E} \rightarrow \overrightarrow{E'}$ is injective, and that $f: E \rightarrow E'$ is surjective iff $\overrightarrow{f}: \overrightarrow{E} \rightarrow \overrightarrow{E'}$ is surjective. An affine map $f: E \rightarrow E'$ is constant iff $\overrightarrow{f}: \overrightarrow{E} \rightarrow \overrightarrow{E'}$ is the null (constant) linear map equal to 0 for all $v \in \overrightarrow{E}$.

If E is an affine space of dimension m and (a_0, a_1, \dots, a_m) is an affine frame for E , then for any other affine space F and for any sequence (b_0, b_1, \dots, b_m) of $m+1$ points in F , there is a unique affine map $f: E \rightarrow F$ such that $f(a_i) = b_i$, for $0 \leq i \leq m$. Indeed, f must be such that

$$f(\lambda_0 a_0 + \dots + \lambda_m a_m) = \lambda_0 b_0 + \dots + \lambda_m b_m,$$

where $\lambda_0 + \dots + \lambda_m = 1$, and this defines a unique affine map on all of E , since (a_0, a_1, \dots, a_m) is an affine frame for E .

Using affine frames, affine maps can be represented in terms of matrices. We explain how an affine map $f: E \rightarrow E$ is represented with respect to a frame (a_0, \dots, a_n) in E , the more general case where an affine map $f: E \rightarrow F$ is represented with respect to two affine frames (a_0, \dots, a_n) in E and (b_0, \dots, b_m) in F being analogous. Since

$$f(a_0 + x) = f(a_0) + \overrightarrow{f}(x)$$

for all $x \in \overrightarrow{E}$, we have

$$\overrightarrow{a_0 f(a_0 + x)} = \overrightarrow{a_0 f(a_0)} + \overrightarrow{f}(x).$$

Since x , $\overrightarrow{a_0 f(a_0)}$, and $\overrightarrow{a_0 f(a_0 + x)}$, can be expressed as

$$\begin{aligned} x &= x_1 \overrightarrow{a_0 a_1} + \dots + x_n \overrightarrow{a_0 a_n}, \\ \overrightarrow{a_0 f(a_0)} &= b_1 \overrightarrow{a_0 a_1} + \dots + b_n \overrightarrow{a_0 a_n}, \\ \overrightarrow{a_0 f(a_0 + x)} &= y_1 \overrightarrow{a_0 a_1} + \dots + y_n \overrightarrow{a_0 a_n}, \end{aligned}$$

if $A = (a_{ij})$ is the $n \times n$ matrix of the linear map \overrightarrow{f} over the basis $(\overrightarrow{a_0 a_1}, \dots, \overrightarrow{a_0 a_n})$, letting x , y , and b denote the column vectors of components (x_1, \dots, x_n) , (y_1, \dots, y_n) , and (b_1, \dots, b_n) ,

$$\overrightarrow{a_0 f(a_0 + x)} = \overrightarrow{a_0 f(a_0)} + \overrightarrow{f}(x)$$

is equivalent to

$$y = Ax + b.$$

Note that $b \neq 0$ unless $f(a_0) = a_0$. Thus, f is generally not a linear transformation, unless it has a *fixed point*, i.e., there is a point a_0 such that $f(a_0) = a_0$. The vector b is the “translation

part" of the affine map. Affine maps do not always have a fixed point. Obviously, nonnull translations have no fixed point. A less trivial example is given by the affine map

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

This map is a reflection about the x -axis followed by a translation along the x -axis. The affine map

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} 1 & -\sqrt{3} \\ \sqrt{3}/4 & 1/4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

can also be written as

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} 2 & 0 \\ 0 & 1/2 \end{pmatrix} \begin{pmatrix} 1/2 & -\sqrt{3}/2 \\ \sqrt{3}/2 & 1/2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

which shows that it is the composition of a rotation of angle $\pi/3$, followed by a stretch (by a factor of 2 along the x -axis, and by a factor of $\frac{1}{2}$ along the y -axis), followed by a translation. It is easy to show that this affine map has a unique fixed point. On the other hand, the affine map

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} 8/5 & -6/5 \\ 3/10 & 2/5 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

has no fixed point, even though

$$\begin{pmatrix} 8/5 & -6/5 \\ 3/10 & 2/5 \end{pmatrix} = \begin{pmatrix} 2 & 0 \\ 0 & 1/2 \end{pmatrix} \begin{pmatrix} 4/5 & -3/5 \\ 3/5 & 4/5 \end{pmatrix},$$

and the second matrix is a rotation of angle θ such that $\cos \theta = \frac{4}{5}$ and $\sin \theta = \frac{3}{5}$. For more on fixed points of affine maps, see the problems.

There is a useful trick to convert the equation $y = Ax + b$ into what looks like a linear equation. The trick is to consider an $(n+1) \times (n+1)$ matrix. We add 1 as the $(n+1)$ th component to the vectors x , y , and b , and form the $(n+1) \times (n+1)$ matrix

$$\begin{pmatrix} A & b \\ 0 & 1 \end{pmatrix}$$

so that $y = Ax + b$ is equivalent to

$$\begin{pmatrix} y \\ 1 \end{pmatrix} = \begin{pmatrix} A & b \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ 1 \end{pmatrix}.$$

This trick is very useful in kinematics and dynamics, where A is a rotation matrix. Such affine maps are called *rigid motions*.

If $f: E \rightarrow E'$ is a bijective affine map, given any three collinear points a, b, c in E , with $a \neq b$, where, say, $c = (1 - \lambda)a + \lambda b$, since f preserves barycenters, we have $f(c) =$

$(1 - \lambda)f(a) + \lambda f(b)$, which shows that $f(a), f(b), f(c)$ are collinear in E' . There is a converse to this property, which is simpler to state when the ground field is $K = \mathbb{R}$. The converse states that given any bijective function $f: E \rightarrow E'$ between two real affine spaces of the same dimension $n \geq 2$, if f maps any three collinear points to collinear points, then f is affine. The proof is rather long (see Berger [8] or Samuel [90]).

Given three collinear points a, b, c , where $a \neq c$, we have $b = (1 - \beta)a + \beta c$ for some unique β , and we define the *ratio of the sequence* a, b, c , as

$$\text{ratio}(a, b, c) = \frac{\beta}{(1 - \beta)} = \frac{\overrightarrow{ab}}{\overrightarrow{bc}},$$

provided that $\beta \neq 1$, i.e., $b \neq c$. When $b = c$, we agree that $\text{ratio}(a, b, c) = \infty$. We warn our readers that other authors define the ratio of a, b, c as $-\text{ratio}(a, b, c) = \frac{\overrightarrow{ba}}{\overrightarrow{bc}}$. Since affine maps preserve barycenters, it is clear that affine maps preserve the ratio of three points.

19.8 Affine Groups

We now take a quick look at the bijective affine maps. Given an affine space E , the set of affine bijections $f: E \rightarrow E$ is clearly a group, called the *affine group of E* , and denoted by $\mathbf{GA}(E)$. Recall that the group of bijective linear maps of the vector space \vec{E} is denoted by $\mathbf{GL}(\vec{E})$. Then, the map $f \mapsto \vec{f}$ defines a group homomorphism $L: \mathbf{GA}(E) \rightarrow \mathbf{GL}(\vec{E})$. The kernel of this map is the set of translations on E .

The subset of all linear maps of the form $\lambda \text{id}_{\vec{E}}$, where $\lambda \in \mathbb{R} - \{0\}$, is a subgroup of $\mathbf{GL}(\vec{E})$, and is denoted by $\mathbb{R}^* \text{id}_{\vec{E}}$ (where $\lambda \text{id}_{\vec{E}}(u) = \lambda u$, and $\mathbb{R}^* = \mathbb{R} - \{0\}$). The subgroup $\mathbf{DIL}(E) = L^{-1}(\mathbb{R}^* \text{id}_{\vec{E}})$ of $\mathbf{GA}(E)$ is particularly interesting. It turns out that it is the disjoint union of the translations and of the dilatations of ratio $\lambda \neq 1$. The elements of $\mathbf{DIL}(E)$ are called *affine dilatations*.

Given any point $a \in E$, and any scalar $\lambda \in \mathbb{R}$, a *dilatation or central dilatation (or homothety) of center a and ratio λ* is a map $H_{a,\lambda}$ defined such that

$$H_{a,\lambda}(x) = a + \lambda \overrightarrow{ax},$$

for every $x \in E$.

Remark: The terminology does not seem to be universally agreed upon. The terms *affine dilatation* and *central dilatation* are used by Pedoe [88]. Snapper and Troyer use the term *dilation* for an affine dilatation and *magnification* for a central dilatation [99]. Samuel uses *homothety* for a central dilatation, a direct translation of the French “homothétie” [90]. Since dilation is shorter than dilatation and somewhat easier to pronounce, perhaps we should use that!

Observe that $H_{a,\lambda}(a) = a$, and when $\lambda \neq 0$ and $x \neq a$, $H_{a,\lambda}(x)$ is on the line defined by a and x , and is obtained by “scaling” \overrightarrow{ax} by λ .

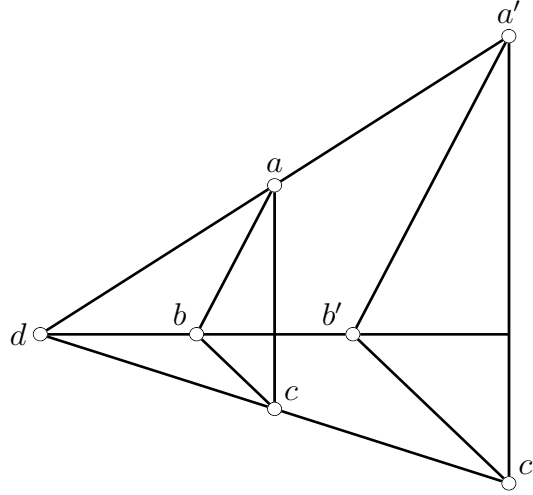


Figure 19.12: The effect of a central dilatation

Figure 19.12 shows the effect of a central dilatation of center d . The triangle (a, b, c) is magnified to the triangle (a', b', c') . Note how every line is mapped to a parallel line.

When $\lambda = 1$, $H_{a,1}$ is the identity. Note that $\overrightarrow{H_{a,\lambda}} = \lambda \text{id}_{\overrightarrow{E}}$. When $\lambda \neq 0$, it is clear that $H_{a,\lambda}$ is an affine bijection. It is immediately verified that

$$H_{a,\lambda} \circ H_{a,\mu} = H_{a,\lambda\mu}.$$

We have the following useful result.

Lemma 19.8. *Given any affine space E , for any affine bijection $f \in \mathbf{GA}(E)$, if $\overrightarrow{f} = \lambda \text{id}_{\overrightarrow{E}}$, for some $\lambda \in \mathbb{R}^*$ with $\lambda \neq 1$, then there is a unique point $c \in E$ such that $f = H_{c,\lambda}$.*

Proof. The proof is straightforward, and is omitted. It is also given in Gallier [43]. □

Clearly, if $\overrightarrow{f} = \text{id}_{\overrightarrow{E}}$, the affine map f is a translation. Thus, the group of affine dilatations $\mathbf{DIL}(E)$ is the disjoint union of the translations and of the dilatations of ratio $\lambda \neq 0, 1$. Affine dilatations can be given a purely geometric characterization.

Another point worth mentioning is that affine bijections preserve the ratio of volumes of parallelotopes. Indeed, given any basis $B = (u_1, \dots, u_m)$ of the vector space \overrightarrow{E} associated with the affine space E , given any $m+1$ affinely independent points (a_0, \dots, a_m) , we can compute the determinant $\det_B(\overrightarrow{a_0a_1}, \dots, \overrightarrow{a_0a_m})$ w.r.t. the basis B . For any bijective affine map $f: E \rightarrow E$, since

$$\det_B\left(\overrightarrow{f(a_0a_1)}, \dots, \overrightarrow{f(a_0a_m)}\right) = \det(\overrightarrow{f}) \det_B(\overrightarrow{a_0a_1}, \dots, \overrightarrow{a_0a_m})$$

and the determinant of a linear map is intrinsic (i.e., depends only on \vec{f} , and not on the particular basis B), we conclude that the ratio

$$\frac{\det_B(\vec{f}(\overrightarrow{a_0a_1}), \dots, \vec{f}(\overrightarrow{a_0a_m}))}{\det_B(\overrightarrow{a_0a_1}, \dots, \overrightarrow{a_0a_m})} = \det(\vec{f})$$

is independent of the basis B . Since $\det_B(\overrightarrow{a_0a_1}, \dots, \overrightarrow{a_0a_m})$ is the volume of the parallelotope spanned by (a_0, \dots, a_m) , where the parallelotope spanned by any point a and the vectors (u_1, \dots, u_m) has unit volume (see Berger [8], Section 9.12), we see that affine bijections preserve the ratio of volumes of parallelotopes. In fact, this ratio is independent of the choice of the parallelotopes of unit volume. In particular, the affine bijections $f \in \mathbf{GA}(E)$ such that $\det(\vec{f}) = 1$ preserve volumes. These affine maps form a subgroup $\mathbf{SA}(E)$ of $\mathbf{GA}(E)$ called the *special affine group of E* . We now take a glimpse at affine geometry.

19.9 Affine Geometry: A Glimpse

In this section we state and prove three fundamental results of affine geometry. Roughly speaking, affine geometry is the study of properties invariant under affine bijections. We now prove one of the oldest and most basic results of affine geometry, the theorem of Thales.

Lemma 19.9. *Given any affine space E , if H_1, H_2, H_3 are any three distinct parallel hyperplanes, and A and B are any two lines not parallel to H_i , letting $a_i = H_i \cap A$ and $b_i = H_i \cap B$, then the following ratios are equal:*

$$\frac{\overrightarrow{a_1a_3}}{\overrightarrow{a_1a_2}} = \frac{\overrightarrow{b_1b_3}}{\overrightarrow{b_1b_2}} = \rho.$$

Conversely, for any point d on the line A , if $\frac{\overrightarrow{a_1d}}{\overrightarrow{a_1a_2}} = \rho$, then $d = a_3$.

Proof. Figure 19.13 illustrates the theorem of Thales. We sketch a proof, leaving the details as an exercise. Since H_1, H_2, H_3 are parallel, they have the same direction \vec{H} , a hyperplane in \vec{E} . Let $u \in \vec{E} - \vec{H}$ be any nonnull vector such that $A = a_1 + \mathbb{R}u$. Since A is not parallel to H , we have $\vec{E} = \vec{H} \oplus \mathbb{R}u$, and thus we can define the linear map $p: \vec{E} \rightarrow \mathbb{R}u$, the projection on $\mathbb{R}u$ parallel to \vec{H} . This linear map induces an affine map $f: E \rightarrow A$, by defining f such that

$$f(b_1 + w) = a_1 + p(w),$$

for all $w \in \vec{E}$. Clearly, $f(b_1) = a_1$, and since H_1, H_2, H_3 all have direction \vec{H} , we also have $f(b_2) = a_2$ and $f(b_3) = a_3$. Since f is affine, it preserves ratios, and thus

$$\frac{\overrightarrow{a_1a_3}}{\overrightarrow{a_1a_2}} = \frac{\overrightarrow{b_1b_3}}{\overrightarrow{b_1b_2}}.$$

The converse is immediate. □

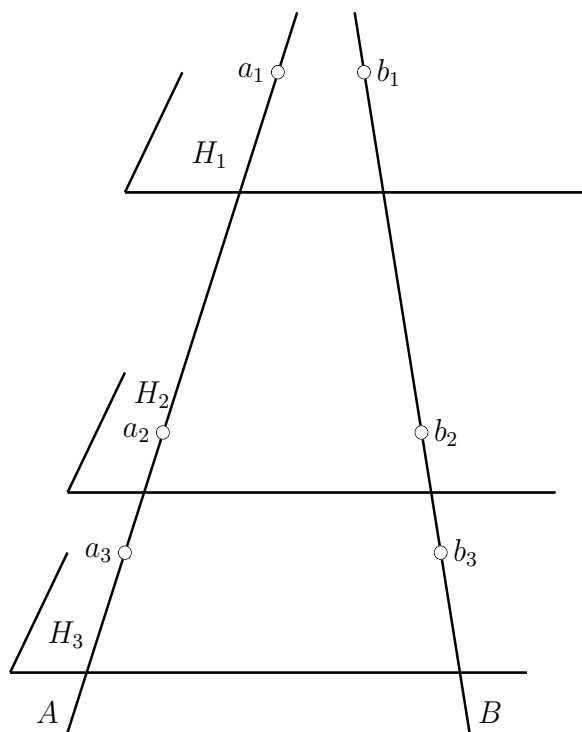


Figure 19.13: The theorem of Thales

We also have the following simple lemma, whose proof is left as an easy exercise.

Lemma 19.10. *Given any affine space E , given any two distinct points $a, b \in E$, and for any affine dilatation f different from the identity, if $a' = f(a)$, $D = \langle a, b \rangle$ is the line passing through a and b , and D' is the line parallel to D and passing through a' , the following are equivalent:*

- (i) $b' = f(b)$;
- (ii) *If f is a translation, then b' is the intersection of D' with the line parallel to $\langle a, a' \rangle$ passing through b ;*
If f is a dilatation of center c , then $b' = D' \cap \langle c, b \rangle$.

The first case is the parallelogram law, and the second case follows easily from Thales' theorem.

We are now ready to prove two classical results of affine geometry, Pappus's theorem and Desargues's theorem. Actually, these results are theorems of projective geometry, and we are stating affine versions of these important results. There are stronger versions that are best proved using projective geometry.

Lemma 19.11. *Given any affine plane E , any two distinct lines D and D' , then for any distinct points a, b, c on D and a', b', c' on D' , if a, b, c, a', b', c' are distinct from the intersection of D and D' (if D and D' intersect) and if the lines $\langle a, b' \rangle$ and $\langle a', b \rangle$ are parallel, and the lines $\langle b, c' \rangle$ and $\langle b', c \rangle$ are parallel, then the lines $\langle a, c' \rangle$ and $\langle a', c \rangle$ are parallel.*

Proof. Pappus's theorem is illustrated in Figure 19.14. If D and D' are not parallel, let d be their intersection. Let f be the dilatation of center d such that $f(a) = b$, and let g be the dilatation of center d such that $g(b) = c$. Since the lines $\langle a, b' \rangle$ and $\langle a', b \rangle$ are parallel, and the lines $\langle b, c' \rangle$ and $\langle b', c \rangle$ are parallel, by Lemma 19.10 we have $a' = f(b')$ and $b' = g(c')$. However, we observed that dilatations with the same center commute, and thus $f \circ g = g \circ f$, and thus, letting $h = g \circ f$, we get $c = h(a)$ and $a' = h(c')$. Again, by Lemma 19.10, the lines $\langle a, c' \rangle$ and $\langle a', c \rangle$ are parallel. If D and D' are parallel, we use translations instead of dilatations. \square

There is a converse to Pappus's theorem, which yields a fancier version of Pappus's theorem, but it is easier to prove it using projective geometry. It should be noted that in axiomatic presentations of projective geometry, Pappus's theorem is equivalent to the commutativity of the ground field K (in the present case, $K = \mathbb{R}$). We now prove an affine version of Desargues's theorem.

Lemma 19.12. *Given any affine space E , and given any two triangles (a, b, c) and (a', b', c') , where a, b, c, a', b', c' are all distinct, if $\langle a, b \rangle$ and $\langle a', b' \rangle$ are parallel and $\langle b, c \rangle$ and $\langle b', c' \rangle$ are parallel, then $\langle a, c \rangle$ and $\langle a', c' \rangle$ are parallel iff the lines $\langle a, a' \rangle$, $\langle b, b' \rangle$, and $\langle c, c' \rangle$ are either parallel or concurrent (i.e., intersect in a common point).*

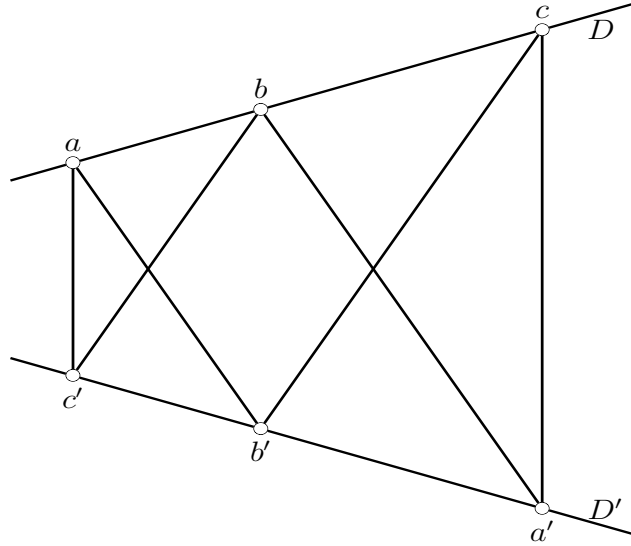


Figure 19.14: Pappus's theorem (affine version)

Proof. We prove half of the lemma, the direction in which it is assumed that $\langle a, c \rangle$ and $\langle a', c' \rangle$ are parallel, leaving the converse as an exercise. Since the lines $\langle a, b \rangle$ and $\langle a', b' \rangle$ are parallel, the points a, b, a', b' are coplanar. Thus, either $\langle a, a' \rangle$ and $\langle b, b' \rangle$ are parallel, or they have some intersection d . We consider the second case where they intersect, leaving the other case as an easy exercise. Let f be the dilatation of center d such that $f(a) = a'$. By Lemma 19.10, we get $f(b) = b'$. If $f(c) = c''$, again by Lemma 19.10 twice, the lines $\langle b, c \rangle$ and $\langle b', c'' \rangle$ are parallel, and the lines $\langle a, c \rangle$ and $\langle a', c'' \rangle$ are parallel. From this it follows that $c'' = c'$. Indeed, recall that $\langle b, c \rangle$ and $\langle b', c' \rangle$ are parallel, and similarly $\langle a, c \rangle$ and $\langle a', c' \rangle$ are parallel. Thus, the lines $\langle b', c'' \rangle$ and $\langle b', c' \rangle$ are identical, and similarly the lines $\langle a', c'' \rangle$ and $\langle a', c' \rangle$ are identical. Since $\overrightarrow{a'c'}$ and $\overrightarrow{b'c'}$ are linearly independent, these lines have a unique intersection, which must be $c'' = c'$.

The direction where it is assumed that the lines $\langle a, a' \rangle$, $\langle b, b' \rangle$ and $\langle c, c' \rangle$, are either parallel or concurrent is left as an exercise (in fact, the proof is quite similar). \square

Desargues's theorem is illustrated in Figure 19.15.

There is a fancier version of Desargues's theorem, but it is easier to prove it using projective geometry. It should be noted that in axiomatic presentations of projective geometry, Desargues's theorem is related to the associativity of the ground field K (in the present case, $K = \mathbb{R}$). Also, Desargues's theorem yields a geometric characterization of the affine dilatations. An affine dilatation f on an affine space E is a bijection that maps every line D to a line $f(D)$ parallel to D . We leave the proof as an exercise.

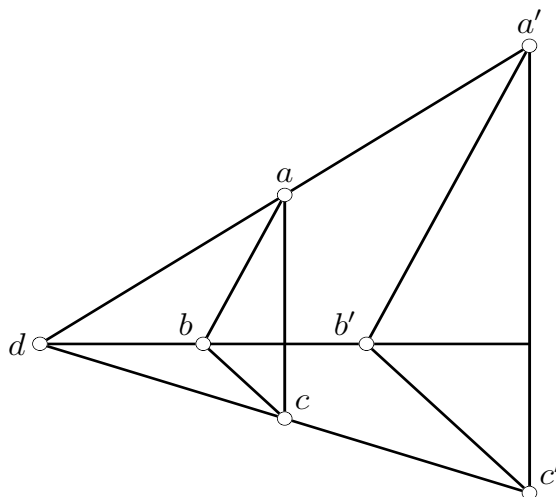


Figure 19.15: Desargues's theorem (affine version)

19.10 Affine Hyperplanes

We now consider affine forms and affine hyperplanes. In Section 19.5 we observed that the set L of solutions of an equation

$$ax + by = c$$

is an affine subspace of \mathbb{A}^2 of dimension 1, in fact, a line (provided that a and b are not both null). It would be equally easy to show that the set P of solutions of an equation

$$ax + by + cz = d$$

is an affine subspace of \mathbb{A}^3 of dimension 2, in fact, a plane (provided that a, b, c are not all null). More generally, the set H of solutions of an equation

$$\lambda_1 x_1 + \cdots + \lambda_m x_m = \mu$$

is an affine subspace of \mathbb{A}^m , and if $\lambda_1, \dots, \lambda_m$ are not all null, it turns out that it is a subspace of dimension $m - 1$ called a *hyperplane*.

We can interpret the equation

$$\lambda_1 x_1 + \cdots + \lambda_m x_m = \mu$$

in terms of the map $f: \mathbb{R}^m \rightarrow \mathbb{R}$ defined such that

$$f(x_1, \dots, x_m) = \lambda_1 x_1 + \cdots + \lambda_m x_m - \mu$$

for all $(x_1, \dots, x_m) \in \mathbb{R}^m$. It is immediately verified that this map is affine, and the set H of solutions of the equation

$$\lambda_1 x_1 + \cdots + \lambda_m x_m = \mu$$

is the *null set*, or *kernel*, of the affine map $f: \mathbb{A}^m \rightarrow \mathbb{R}$, in the sense that

$$H = f^{-1}(0) = \{x \in \mathbb{A}^m \mid f(x) = 0\},$$

where $x = (x_1, \dots, x_m)$.

Thus, it is interesting to consider *affine forms*, which are just affine maps $f: E \rightarrow \mathbb{R}$ from an affine space to \mathbb{R} . Unlike linear forms f^* , for which $\text{Ker } f^*$ is never empty (since it always contains the vector 0), it is possible that $f^{-1}(0) = \emptyset$ for an affine form f . Given an affine map $f: E \rightarrow \mathbb{R}$, we also denote $f^{-1}(0)$ by $\text{Ker } f$, and we call it the *kernel* of f . Recall that an (affine) hyperplane is an affine subspace of codimension 1. The relationship between affine hyperplanes and affine forms is given by the following lemma.

Lemma 19.13. *Let E be an affine space. The following properties hold:*

- (a) *Given any nonconstant affine form $f: E \rightarrow \mathbb{R}$, its kernel $H = \text{Ker } f$ is a hyperplane.*
- (b) *For any hyperplane H in E , there is a nonconstant affine form $f: E \rightarrow \mathbb{R}$ such that $H = \text{Ker } f$. For any other affine form $g: E \rightarrow \mathbb{R}$ such that $H = \text{Ker } g$, there is some $\lambda \in \mathbb{R}$ such that $g = \lambda f$ (with $\lambda \neq 0$).*
- (c) *Given any hyperplane H in E and any (nonconstant) affine form $f: E \rightarrow \mathbb{R}$ such that $H = \text{Ker } f$, every hyperplane H' parallel to H is defined by a nonconstant affine form g such that $g(a) = f(a) - \lambda$, for all $a \in E$ and some $\lambda \in \mathbb{R}$.*

Proof. The proof is straightforward, and is omitted. It is also given in Gallier [43]. □

When E is of dimension n , given an affine frame $(a_0, (u_1, \dots, u_n))$ of E with origin a_0 , recall from Definition 19.5 that every point of E can be expressed uniquely as $x = a_0 + x_1u_1 + \dots + x_nu_n$, where (x_1, \dots, x_n) are the *coordinates* of x with respect to the affine frame $(a_0, (u_1, \dots, u_n))$.

Also recall that every linear form f^* is such that $f^*(x) = \lambda_1x_1 + \dots + \lambda_nx_n$, for every $x = x_1u_1 + \dots + x_nu_n$ and some $\lambda_1, \dots, \lambda_n \in \mathbb{R}$. Since an affine form $f: E \rightarrow \mathbb{R}$ satisfies the property $f(a_0 + x) = f(a_0) + \overrightarrow{f}(x)$, denoting $f(a_0 + x)$ by $f(x_1, \dots, x_n)$, we see that we have

$$f(x_1, \dots, x_n) = \lambda_1x_1 + \dots + \lambda_nx_n + \mu,$$

where $\mu = f(a_0) \in \mathbb{R}$ and $\lambda_1, \dots, \lambda_n \in \mathbb{R}$. Thus, a hyperplane is the set of points whose coordinates (x_1, \dots, x_n) satisfy the (affine) equation

$$\lambda_1x_1 + \dots + \lambda_nx_n + \mu = 0.$$

19.11 Intersection of Affine Spaces

In this section we take a closer look at the intersection of affine subspaces. This subsection can be omitted at first reading.

First, we need a result of linear algebra. Given a vector space E and any two subspaces M and N , there are several interesting linear maps. We have the canonical injections $i: M \rightarrow M+N$ and $j: N \rightarrow M+N$, the canonical injections $in_1: M \rightarrow M \oplus N$ and $in_2: N \rightarrow M \oplus N$, and thus, injections $f: M \cap N \rightarrow M \oplus N$ and $g: M \cap N \rightarrow M \oplus N$, where f is the composition of the inclusion map from $M \cap N$ to M with in_1 , and g is the composition of the inclusion map from $M \cap N$ to N with in_2 . Then, we have the maps $f+g: M \cap N \rightarrow M \oplus N$, and $i-j: M \oplus N \rightarrow M+N$.

Lemma 19.14. *Given a vector space E and any two subspaces M and N , with the definitions above,*

$$0 \longrightarrow M \cap N \xrightarrow{f+g} M \oplus N \xrightarrow{i-j} M+N \longrightarrow 0$$

is a short exact sequence, which means that $f+g$ is injective, $i-j$ is surjective, and that $\text{Im}(f+g) = \text{Ker}(i-j)$. As a consequence, we have the Grassmann relation

$$\dim(M) + \dim(N) = \dim(M+N) + \dim(M \cap N).$$

Proof. It is obvious that $i-j$ is surjective and that $f+g$ is injective. Assume that $(i-j)(u+v) = 0$, where $u \in M$, and $v \in N$. Then, $i(u) = j(v)$, and thus, by definition of i and j , there is some $w \in M \cap N$, such that $i(u) = j(v) = w \in M \cap N$. By definition of f and g , $u = f(w)$ and $v = g(w)$, and thus $\text{Im}(f+g) = \text{Ker}(i-j)$, as desired. The second part of the lemma follows from standard results of linear algebra (see Artin [4], Strang [105], or Lang [67]). \square

We now prove a simple lemma about the intersection of affine subspaces.

Lemma 19.15. *Given any affine space E , for any two nonempty affine subspaces M and N , the following facts hold:*

- (1) $M \cap N \neq \emptyset$ iff $\overrightarrow{ab} \in \overrightarrow{M} + \overrightarrow{N}$ for some $a \in M$ and some $b \in N$.
- (2) $M \cap N$ consists of a single point iff $\overrightarrow{ab} \in \overrightarrow{M} + \overrightarrow{N}$ for some $a \in M$ and some $b \in N$, and $\overrightarrow{M} \cap \overrightarrow{N} = \{0\}$.
- (3) If S is the least affine subspace containing M and N , then $\overrightarrow{S} = \overrightarrow{M} + \overrightarrow{N} + K\overrightarrow{ab}$ (the vector space \overrightarrow{E} is defined over the field K).

Proof. (1) Pick any $a \in M$ and any $b \in N$, which is possible, since M and N are nonempty. Since $\overrightarrow{M} = \{\overrightarrow{ax} \mid x \in M\}$ and $\overrightarrow{N} = \{\overrightarrow{by} \mid y \in N\}$, if $M \cap N \neq \emptyset$, for any $c \in M \cap N$ we have $\overrightarrow{ab} = \overrightarrow{ac} - \overrightarrow{bc}$, with $\overrightarrow{ac} \in \overrightarrow{M}$ and $\overrightarrow{bc} \in \overrightarrow{N}$, and thus, $\overrightarrow{ab} \in \overrightarrow{M} + \overrightarrow{N}$. Conversely, assume that

$\vec{ab} \in \vec{M} + \vec{N}$ for some $a \in M$ and some $b \in N$. Then $\vec{ab} = \vec{ax} + \vec{by}$, for some $x \in M$ and some $y \in N$. But we also have

$$\vec{ab} = \vec{ax} + \vec{xy} + \vec{yb},$$

and thus we get $0 = \vec{xy} + \vec{yb} - \vec{by}$, that is, $\vec{xy} = 2\vec{by}$. Thus, b is the middle of the segment $[x, y]$, and since $\vec{yx} = 2\vec{yb}$, $x = 2b - y$ is the barycenter of the weighted points $(b, 2)$ and $(y, -1)$. Thus x also belongs to N , since N being an affine subspace, it is closed under barycenters. Thus, $x \in M \cap N$, and $M \cap N \neq \emptyset$.

(2) Note that in general, if $M \cap N \neq \emptyset$, then

$$\overrightarrow{M \cap N} = \vec{M} \cap \vec{N},$$

because

$$\overrightarrow{M \cap N} = \{\vec{ab} \mid a, b \in M \cap N\} = \{\vec{ab} \mid a, b \in M\} \cap \{\vec{ab} \mid a, b \in N\} = \vec{M} \cap \vec{N}.$$

Since $M \cap N = c + \overrightarrow{M \cap N}$ for any $c \in M \cap N$, we have

$$M \cap N = c + \vec{M} \cap \vec{N} \quad \text{for any } c \in M \cap N.$$

From this it follows that if $M \cap N \neq \emptyset$, then $M \cap N$ consists of a single point iff $\vec{M} \cap \vec{N} = \{0\}$. This fact together with what we proved in (1) proves (2).

(3) This is left as an easy exercise. □

Remarks:

- (1) The proof of Lemma 19.15 shows that if $M \cap N \neq \emptyset$, then $\vec{ab} \in \vec{M} + \vec{N}$ for all $a \in M$ and all $b \in N$.
- (2) Lemma 19.15 implies that for any two nonempty affine subspaces M and N , if $\vec{E} = \vec{M} \oplus \vec{N}$, then $M \cap N$ consists of a single point. Indeed, if $\vec{E} = \vec{M} \oplus \vec{N}$, then $\vec{ab} \in \vec{E}$ for all $a \in M$ and all $b \in N$, and since $\vec{M} \cap \vec{N} = \{0\}$, the result follows from part (2) of the lemma.

We can now state the following lemma.

Lemma 19.16. *Given an affine space E and any two nonempty affine subspaces M and N , if S is the least affine subspace containing M and N , then the following properties hold:*

- (1) *If $M \cap N = \emptyset$, then*

$$\dim(M) + \dim(N) < \dim(E) + \dim(\vec{M} + \vec{N})$$

and

$$\dim(S) = \dim(M) + \dim(N) + 1 - \dim(\vec{M} \cap \vec{N}).$$

- (2) *If $M \cap N \neq \emptyset$, then*

$$\dim(S) = \dim(M) + \dim(N) - \dim(M \cap N).$$

Proof. The proof is not difficult, using Lemma 19.15 and Lemma 19.14, but we leave it as an exercise. □

19.12 Problems

Problem 19.1. Given a triangle (a, b, c) , give a geometric construction of the barycenter of the weighted points $(a, \frac{1}{4})$, $(b, \frac{1}{4})$, and $(c, \frac{1}{2})$. Give a geometric construction of the barycenter of the weighted points $(a, \frac{3}{2})$, $(b, \frac{3}{2})$, and $(c, -2)$.

Problem 19.2. Given a tetrahedron (a, b, c, d) and any two distinct points $x, y \in \{a, b, c, d\}$, let $m_{x,y}$ be the middle of the edge (x, y) . Prove that the barycenter g of the weighted points $(a, \frac{1}{4})$, $(b, \frac{1}{4})$, $(c, \frac{1}{4})$, and $(d, \frac{1}{4})$ is the common intersection of the line segments $(m_{a,b}, m_{c,d})$, $(m_{a,c}, m_{b,d})$, and $(m_{a,d}, m_{b,c})$. Show that if g_d is the barycenter of the weighted points $(a, \frac{1}{3})$, $(b, \frac{1}{3})$, $(c, \frac{1}{3})$, then g is the barycenter of $(d, \frac{1}{4})$ and $(g_d, \frac{3}{4})$.

Problem 19.3. Let E be a nonempty set, and \vec{E} a vector space and assume that there is a function $\Phi: E \times E \rightarrow \vec{E}$, such that if we denote $\Phi(a, b)$ by \vec{ab} , the following properties hold:

- (1) $\vec{ab} + \vec{bc} = \vec{ac}$, for all $a, b, c \in E$;
- (2) For every $a \in E$, the map $\Phi_a: E \rightarrow \vec{E}$ defined such that for every $b \in E$, $\Phi_a(b) = \vec{ab}$, is a bijection.

Let $\Psi_a: \vec{E} \rightarrow E$ be the inverse of $\Phi_a: E \rightarrow \vec{E}$.

Prove that the function $+: E \times \vec{E} \rightarrow E$ defined such that

$$a + u = \Psi_a(u)$$

for all $a \in E$ and all $u \in \vec{E}$ makes $(E, \vec{E}, +)$ into an affine space.

Note. We showed in the text that an affine space $(E, \vec{E}, +)$ satisfies the properties stated above. Thus, we obtain an equivalent characterization of affine spaces.

Problem 19.4. Given any three points a, b, c in the affine plane \mathbb{A}^2 , letting (a_1, a_2) , (b_1, b_2) , and (c_1, c_2) be the coordinates of a, b, c , with respect to the standard affine frame for \mathbb{A}^2 , prove that a, b, c are collinear iff

$$\begin{vmatrix} a_1 & b_1 & c_1 \\ a_2 & b_2 & c_2 \\ 1 & 1 & 1 \end{vmatrix} = 0,$$

i.e., the determinant is null.

Letting (a_0, a_1, a_2) , (b_0, b_1, b_2) , and (c_0, c_1, c_2) be the barycentric coordinates of a, b, c with respect to the standard affine frame for \mathbb{A}^2 , prove that a, b, c are collinear iff

$$\begin{vmatrix} a_0 & b_0 & c_0 \\ a_1 & b_1 & c_1 \\ a_2 & b_2 & c_2 \end{vmatrix} = 0.$$

Given any four points a, b, c, d in the affine space \mathbb{A}^3 , letting (a_1, a_2, a_3) , (b_1, b_2, b_3) , (c_1, c_2, c_3) , and (d_1, d_2, d_3) be the coordinates of a, b, c, d , with respect to the standard affine frame for \mathbb{A}^3 , prove that a, b, c, d are coplanar iff

$$\begin{vmatrix} a_1 & b_1 & c_1 & d_1 \\ a_2 & b_2 & c_2 & d_2 \\ a_3 & b_3 & c_3 & d_3 \\ 1 & 1 & 1 & 1 \end{vmatrix} = 0,$$

i.e., the determinant is null.

Letting (a_0, a_1, a_2, a_3) , (b_0, b_1, b_2, b_3) , (c_0, c_1, c_2, c_3) , and (d_0, d_1, d_2, d_3) be the barycentric coordinates of a, b, c, d , with respect to the standard affine frame for \mathbb{A}^3 , prove that a, b, c, d are coplanar iff

$$\begin{vmatrix} a_0 & b_0 & c_0 & d_0 \\ a_1 & b_1 & c_1 & d_1 \\ a_2 & b_2 & c_2 & d_2 \\ a_3 & b_3 & c_3 & d_3 \end{vmatrix} = 0.$$

Problem 19.5. The function $f : \mathbb{A} \rightarrow \mathbb{A}^3$ given by

$$t \mapsto (t, t^2, t^3)$$

defines what is called a *twisted cubic* curve. Given any four pairwise distinct values t_1, t_2, t_3, t_4 , prove that the points $f(t_1), f(t_2), f(t_3)$, and $f(t_4)$ are not coplanar.

Hint. Have you heard of the Vandermonde determinant?

Problem 19.6. For any two distinct points $a, b \in \mathbb{A}^2$ of barycentric coordinates (a_0, a_1, a_2) and (b_0, b_1, b_2) with respect to any given affine frame (O, i, j) , show that the equation of the line $\langle a, b \rangle$ determined by a and b is

$$\begin{vmatrix} a_0 & b_0 & x \\ a_1 & b_1 & y \\ a_2 & b_2 & z \end{vmatrix} = 0,$$

or, equivalently,

$$(a_1b_2 - a_2b_1)x + (a_2b_0 - a_0b_2)y + (a_0b_1 - a_1b_0)z = 0,$$

where (x, y, z) are the barycentric coordinates of the generic point on the line $\langle a, b \rangle$.

Prove that the equation of a line in barycentric coordinates is of the form

$$ux + vy + wz = 0,$$

where $u \neq v$ or $v \neq w$ or $u \neq w$. Show that two equations

$$ux + vy + wz = 0 \quad \text{and} \quad u'x + v'y + w'z = 0$$

represent the same line in barycentric coordinates iff $(u', v', w') = \lambda(u, v, w)$ for some $\lambda \in \mathbb{R}$ (with $\lambda \neq 0$).

A triple (u, v, w) where $u \neq v$ or $v \neq w$ or $u \neq w$ is called a system of *tangential coordinates* of the line defined by the equation

$$ux + vy + wz = 0.$$

Problem 19.7. Given two lines D and D' in \mathbb{A}^2 defined by tangential coordinates (u, v, w) and (u', v', w') (as defined in Problem 19.6), let

$$d = \begin{vmatrix} u & v & w \\ u' & v' & w' \\ 1 & 1 & 1 \end{vmatrix} = vw' - wv' + wu' - uw' + uv' - vu'.$$

(a) Prove that D and D' have a unique intersection point iff $d \neq 0$, and that when it exists, the barycentric coordinates of this intersection point are

$$\frac{1}{d}(vw' - wv', wu' - uw', uv' - vu').$$

(b) Letting (O, i, j) be any affine frame for \mathbb{A}^2 , recall that when $x + y + z = 0$, for any point a , the vector

$$xa\vec{O} + ya\vec{i} + za\vec{j}$$

is independent of a and equal to

$$y\vec{O}i + z\vec{O}j = (y, z).$$

The triple (x, y, z) such that $x + y + z = 0$ is called the *barycentric coordinates* of the vector $y\vec{O}i + z\vec{O}j$ w.r.t. the affine frame (O, i, j) .

Given any affine frame (O, i, j) , prove that for $u \neq v$ or $v \neq w$ or $u \neq w$, the line of equation

$$ux + vy + wz = 0$$

in barycentric coordinates (x, y, z) (where $x + y + z = 1$) has for direction the set of vectors of barycentric coordinates (x, y, z) such that

$$ux + vy + wz = 0$$

(where $x + y + z = 0$).

Prove that D and D' are parallel iff $d = 0$. In this case, if $D \neq D'$, show that the common direction of D and D' is defined by the vector of barycentric coordinates

$$(vw' - wv', wu' - uw', uv' - vu').$$

(c) Given three lines D , D' , and D'' , at least two of which are distinct and defined by tangential coordinates (u, v, w) , (u', v', w') , and (u'', v'', w'') , prove that D , D' , and D'' are parallel or have a unique intersection point iff

$$\begin{vmatrix} u & v & w \\ u' & v' & w' \\ u'' & v'' & w'' \end{vmatrix} = 0.$$

Problem 19.8. Let (A, B, C) be a triangle in \mathbb{A}^2 . Let M, N, P be three points respectively on the lines BC, CA , and AB , of barycentric coordinates $(0, m', m'')$, $(n, 0, n'')$, and $(p, p', 0)$, w.r.t. the affine frame (A, B, C) .

(a) Assuming that $M \neq C$, $N \neq A$, and $P \neq B$, i.e., $m'n''p \neq 0$, show that

$$\frac{\overrightarrow{MB}}{\overrightarrow{MC}} \frac{\overrightarrow{NC}}{\overrightarrow{NA}} \frac{\overrightarrow{PA}}{\overrightarrow{PB}} = -\frac{m''np'}{m'n''p}.$$

(b) Prove *Menelaus's theorem*: The points M, N, P are collinear iff

$$m''np' + m'n''p = 0.$$

When $M \neq C$, $N \neq A$, and $P \neq B$, this is equivalent to

$$\frac{\overrightarrow{MB}}{\overrightarrow{MC}} \frac{\overrightarrow{NC}}{\overrightarrow{NA}} \frac{\overrightarrow{PA}}{\overrightarrow{PB}} = 1.$$

(c) Prove *Ceva's theorem*: The lines AM, BN, CP have a unique intersection point or are parallel iff

$$m''np' - m'n''p = 0.$$

When $M \neq C$, $N \neq A$, and $P \neq B$, this is equivalent to

$$\frac{\overrightarrow{MB}}{\overrightarrow{MC}} \frac{\overrightarrow{NC}}{\overrightarrow{NA}} \frac{\overrightarrow{PA}}{\overrightarrow{PB}} = -1.$$

Problem 19.9. This problem uses notions and results from Problems 19.6 and 19.7. In view of (a) and (b) of Problem 19.7, it is natural to extend the notion of barycentric coordinates of a point in \mathbb{A}^2 as follows. Given any affine frame (a, b, c) in \mathbb{A}^2 , we will say that the barycentric coordinates (x, y, z) of a point M , where $x + y + z = 1$, are the *normalized barycentric coordinates* of M . Then, any triple (x, y, z) such that $x + y + z \neq 0$ is also called a system of barycentric coordinates for the point of normalized barycentric coordinates

$$\frac{1}{x + y + z} (x, y, z).$$

With this convention, the intersection of the two lines D and D' is either a point or a vector, in both cases of barycentric coordinates

$$(vw' - wv', wu' - uw', uv' - vu').$$

When the above is a vector, we can think of it as a point at infinity (in the direction of the line defined by that vector).

Let (D_0, D'_0) , (D_1, D'_1) , and (D_2, D'_2) be three pairs of six distinct lines, such that the four lines belonging to any union of two of the above pairs are neither parallel nor concurrent (have a common intersection point). If D_0 and D'_0 have a unique intersection point, let M be this point, and if D_0 and D'_0 are parallel, let M denote a nonnull vector defining the common direction of D_0 and D'_0 . In either case, let (m, m', m'') be the barycentric coordinates of M , as explained at the beginning of the problem. We call M the *intersection* of D_0 and D'_0 . Similarly, define $N = (n, n', n'')$ as the intersection of D_1 and D'_1 , and $P = (p, p', p'')$ as the intersection of D_2 and D'_2 .

Prove that

$$\begin{vmatrix} m & n & p \\ m' & n' & p' \\ m'' & n'' & p'' \end{vmatrix} = 0$$

iff either

- (i) (D_0, D'_0) , (D_1, D'_1) , and (D_2, D'_2) are pairs of parallel lines; or
- (ii) the lines of some pair (D_i, D'_i) are parallel, each pair (D_j, D'_j) (with $j \neq i$) has a unique intersection point, and these two intersection points are distinct and determine a line parallel to the lines of the pair (D_i, D'_i) ; or
- (iii) each pair (D_i, D'_i) ($i = 0, 1, 2$) has a unique intersection point, and these points M, N, P are distinct and collinear.

Problem 19.10. Prove the following version of *Desargues's theorem*. Let A, B, C, A', B', C' be six distinct points of \mathbb{A}^2 . If no three of these points are collinear, then the lines $AA', BB',$ and CC' are parallel or collinear iff the intersection points M, N, P (in the sense of Problem 19.7) of the pairs of lines $(BC, B'C'), (CA, C'A'),$ and $(AB, A'B')$ are collinear in the sense of Problem 19.9.

Problem 19.11. Prove the following version of *Pappus's theorem*. Let D and D' be distinct lines, and let A, B, C and A', B', C' be distinct points respectively on D and D' . If these points are all distinct from the intersection of D and D' (if it exists), then the intersection points (in the sense of Problem 19.7) of the pairs of lines $(BC', CB'), (CA', AC'),$ and (AB', BA') are collinear in the sense of Problem 19.9.

Problem 19.12. The purpose of this problem is to prove *Pascal's theorem* for the nondegenerate conics. In the affine plane \mathbb{A}^2 , a *conic* is the set of points of coordinates (x, y) such that

$$\alpha x^2 + \beta y^2 + 2\gamma xy + 2\delta x + 2\lambda y + \mu = 0,$$

where $\alpha \neq 0$ or $\beta \neq 0$ or $\gamma \neq 0$. We can write the equation of the conic as

$$(x, y, 1) \begin{pmatrix} \alpha & \gamma & \delta \\ \gamma & \beta & \lambda \\ \delta & \lambda & \mu \end{pmatrix} \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = 0.$$

If we now use barycentric coordinates (x, y, z) (where $x + y + z = 1$), we can write

$$\begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix}.$$

Let

$$B = \begin{pmatrix} \alpha & \gamma & \delta \\ \gamma & \beta & \lambda \\ \delta & \lambda & \mu \end{pmatrix}, \quad C = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}, \quad X = \begin{pmatrix} x \\ y \\ z \end{pmatrix}.$$

(a) Letting $A = C^\top BC$, prove that the equation of the conic becomes

$$X^\top AX = 0.$$

Prove that A is symmetric, that $\det(A) = \det(B)$, and that $X^\top AX$ is homogeneous of degree 2. The equation $X^\top AX = 0$ is called the *homogeneous equation* of the conic.

We say that a conic of homogeneous equation $X^\top AX = 0$ is *nondegenerate* if $\det(A) \neq 0$, and *degenerate* if $\det(A) = 0$. Show that this condition does not depend on the choice of the affine frame.

(b) Given an affine frame (A, B, C) , prove that any conic passing through A, B, C has an equation of the form

$$ayz + bxz + cxy = 0.$$

Prove that a conic containing more than one point is degenerate iff it contains three distinct collinear points. In this case, the conic is the union of two lines.

(c) Prove *Pascal's theorem*. Given any six distinct points A, B, C, A', B', C' , if no three of the above points are collinear, then a nondegenerate conic passes through these six points iff the intersection points M, N, P (in the sense of Problem 19.7) of the pairs of lines (BC', CB') , (CA', AC') and (AB', BA') are collinear in the sense of Problem 19.9.

Hint. Use the affine frame (A, B, C) , and let (a, a', a'') , (b, b', b'') , and (c, c', c'') be the barycentric coordinates of A', B', C' respectively, and show that M, N, P have barycentric coordinates

$$(bc, cb', c''b), \quad (c'a, c'a', c''a'), \quad (ab'', a''b', a''b'').$$

Problem 19.13. The *centroid* of a triangle (a, b, c) is the barycenter of $(a, \frac{1}{3})$, $(b, \frac{1}{3})$, $(c, \frac{1}{3})$. If an affine map takes the vertices of triangle $\Delta_1 = \{(0, 0), (6, 0), (0, 9)\}$ to the vertices of triangle $\Delta_2 = \{(1, 1), (5, 4), (3, 1)\}$, does it also take the centroid of Δ_1 to the centroid of Δ_2 ? Justify your answer.

Problem 19.14. Let E be an affine space over \mathbb{R} , and let (a_1, \dots, a_n) be any $n \geq 3$ points in E . Let $(\lambda_1, \dots, \lambda_n)$ be any n scalars in \mathbb{R} , with $\lambda_1 + \dots + \lambda_n = 1$. Show that there must be some i , $1 \leq i \leq n$, such that $\lambda_i \neq 1$. To simplify the notation, assume that $\lambda_1 \neq 1$. Show that the barycenter $\lambda_1 a_1 + \dots + \lambda_n a_n$ can be obtained by first determining the barycenter b of the $n - 1$ points a_2, \dots, a_n assigned some appropriate weights, and then the barycenter of

a_1 and b assigned the weights λ_1 and $\lambda_2 + \cdots + \lambda_n$. From this, show that the barycenter of any $n \geq 3$ points can be determined by repeated computations of barycenters of two points. Deduce from the above that a nonempty subset V of E is an affine subspace iff whenever V contains any two points $x, y \in V$, then V contains the entire line $(1 - \lambda)x + \lambda y$, $\lambda \in \mathbb{R}$.

Problem 19.15. Assume that K is a field such that $2 = 1 + 1 \neq 0$, and let E be an affine space over K . In the case where $\lambda_1 + \cdots + \lambda_n = 1$ and $\lambda_i = 1$, for $1 \leq i \leq n$ and $n \geq 3$, show that the barycenter $a_1 + a_2 + \cdots + a_n$ can still be computed by repeated computations of barycenters of two points.

Finally, assume that the field K contains at least three elements (thus, there is some $\mu \in K$ such that $\mu \neq 0$ and $\mu \neq 1$, but $2 = 1 + 1 = 0$ is possible). Prove that the barycenter of any $n \geq 3$ points can be determined by repeated computations of barycenters of two points. Prove that a nonempty subset V of E is an affine subspace iff whenever V contains any two points $x, y \in V$, then V contains the entire line $(1 - \lambda)x + \lambda y$, $\lambda \in K$.

Hint. When $2 = 0$, $\lambda_1 + \cdots + \lambda_n = 1$ and $\lambda_i = 1$, for $1 \leq i \leq n$, show that n must be odd, and that the problem reduces to computing the barycenter of three points in two steps involving two barycenters. Since there is some $\mu \in K$ such that $\mu \neq 0$ and $\mu \neq 1$, note that μ^{-1} and $(1 - \mu)^{-1}$ both exist, and use the fact that

$$\frac{-\mu}{1 - \mu} + \frac{1}{1 - \mu} = 1.$$

Problem 19.16. (i) Let (a, b, c) be three points in \mathbb{A}^2 , and assume that (a, b, c) are not collinear. For any point $x \in \mathbb{A}^2$, if $x = \lambda_0 a + \lambda_1 b + \lambda_2 c$, where $(\lambda_0, \lambda_1, \lambda_2)$ are the barycentric coordinates of x with respect to (a, b, c) , show that

$$\lambda_0 = \frac{\det(\vec{xb}, \vec{bc})}{\det(\vec{ab}, \vec{ac})}, \quad \lambda_1 = \frac{\det(\vec{ax}, \vec{ac})}{\det(\vec{ab}, \vec{ac})}, \quad \lambda_2 = \frac{\det(\vec{ab}, \vec{ax})}{\det(\vec{ab}, \vec{ac})}.$$

Conclude that $\lambda_0, \lambda_1, \lambda_2$ are certain signed ratios of the areas of the triangles (a, b, c) , (x, a, b) , (x, a, c) , and (x, b, c) .

(ii) Let (a, b, c) be three points in \mathbb{A}^3 , and assume that (a, b, c) are not collinear. For any point x in the plane determined by (a, b, c) , if $x = \lambda_0 a + \lambda_1 b + \lambda_2 c$, where $(\lambda_0, \lambda_1, \lambda_2)$ are the barycentric coordinates of x with respect to (a, b, c) , show that

$$\lambda_0 = \frac{\vec{xb} \times \vec{bc}}{\vec{ab} \times \vec{ac}}, \quad \lambda_1 = \frac{\vec{ax} \times \vec{ac}}{\vec{ab} \times \vec{ac}}, \quad \lambda_2 = \frac{\vec{ab} \times \vec{ax}}{\vec{ab} \times \vec{ac}}.$$

Given any point O not in the plane of the triangle (a, b, c) , prove that

$$\lambda_1 = \frac{\det(\vec{Oa}, \vec{Ox}, \vec{Oc})}{\det(\vec{Oa}, \vec{Ob}, \vec{Oc})}, \quad \lambda_2 = \frac{\det(\vec{Oa}, \vec{Ob}, \vec{Ox})}{\det(\vec{Oa}, \vec{Ob}, \vec{Oc})},$$

and

$$\lambda_0 = \frac{\det(\overrightarrow{Ox}, \overrightarrow{Ob}, \overrightarrow{Oc})}{\det(\overrightarrow{Oa}, \overrightarrow{Ob}, \overrightarrow{Oc})}.$$

(iii) Let (a, b, c, d) be four points in \mathbb{A}^3 , and assume that (a, b, c, d) are not coplanar. For any point $x \in \mathbb{A}^3$, if $x = \lambda_0 a + \lambda_1 b + \lambda_2 c + \lambda_3 d$, where $(\lambda_0, \lambda_1, \lambda_2, \lambda_3)$ are the barycentric coordinates of x with respect to (a, b, c, d) , show that

$$\lambda_1 = \frac{\det(\overrightarrow{ax}, \overrightarrow{ac}, \overrightarrow{ad})}{\det(\overrightarrow{ab}, \overrightarrow{ac}, \overrightarrow{ad})}, \quad \lambda_2 = \frac{\det(\overrightarrow{ab}, \overrightarrow{ax}, \overrightarrow{ad})}{\det(\overrightarrow{ab}, \overrightarrow{ac}, \overrightarrow{ad})}, \quad \lambda_3 = \frac{\det(\overrightarrow{ab}, \overrightarrow{ac}, \overrightarrow{ax})}{\det(\overrightarrow{ab}, \overrightarrow{ac}, \overrightarrow{ad})},$$

and

$$\lambda_0 = \frac{\det(\overrightarrow{xb}, \overrightarrow{bc}, \overrightarrow{bd})}{\det(\overrightarrow{ab}, \overrightarrow{ac}, \overrightarrow{ad})}.$$

Conclude that $\lambda_0, \lambda_1, \lambda_2, \lambda_3$ are certain signed ratios of the volumes of the five tetrahedra (a, b, c, d) , (x, a, b, c) , (x, a, b, d) , (x, a, c, d) , and (x, b, c, d) .

(iv) Let (a_0, \dots, a_m) be $m+1$ points in \mathbb{A}^m , and assume that they are affinely independent. For any point $x \in \mathbb{A}^m$, if $x = \lambda_0 a_0 + \dots + \lambda_m a_m$, where $(\lambda_0, \dots, \lambda_m)$ are the barycentric coordinates of x with respect to (a_0, \dots, a_m) , show that

$$\lambda_i = \frac{\det(\overrightarrow{a_0 a_1}, \dots, \overrightarrow{a_0 a_{i-1}}, \overrightarrow{a_0 x}, \overrightarrow{a_0 a_{i+1}}, \dots, \overrightarrow{a_0 a_m})}{\det(\overrightarrow{a_0 a_1}, \dots, \overrightarrow{a_0 a_{i-1}}, \overrightarrow{a_0 a_i}, \overrightarrow{a_0 a_{i+1}}, \dots, \overrightarrow{a_0 a_m})}$$

for every i , $1 \leq i \leq m$, and

$$\lambda_0 = \frac{\det(\overrightarrow{xa_1}, \overrightarrow{a_1 a_2}, \dots, \overrightarrow{a_1 a_m})}{\det(\overrightarrow{a_0 a_1}, \dots, \overrightarrow{a_0 a_i}, \dots, \overrightarrow{a_0 a_m})}.$$

Conclude that λ_i is the signed ratio of the volumes of the simplexes $(a_0, \dots, x, \dots, a_m)$ and $(a_0, \dots, a_i, \dots, a_m)$, where $0 \leq i \leq m$.

Problem 19.17. With respect to the standard affine frame for the plane \mathbb{A}^2 , consider the three geometric transformations f_1, f_2, f_3 defined by

$$\begin{aligned} x' &= -\frac{1}{4}x - \frac{\sqrt{3}}{4}y + \frac{3}{4}, & y' &= \frac{\sqrt{3}}{4}x - \frac{1}{4}y + \frac{\sqrt{3}}{4}, \\ x' &= -\frac{1}{4}x + \frac{\sqrt{3}}{4}y - \frac{3}{4}, & y' &= -\frac{\sqrt{3}}{4}x - \frac{1}{4}y + \frac{\sqrt{3}}{4}, \\ x' &= \frac{1}{2}x, & y' &= \frac{1}{2}y + \frac{\sqrt{3}}{2}. \end{aligned}$$

(a) Prove that these maps are affine. Can you describe geometrically what their action is (rotation, translation, scaling)?

(b) Given any polygonal line L , define the following sequence of polygonal lines:

$$\begin{aligned} S_0 &= L, \\ S_{n+1} &= f_1(S_n) \cup f_2(S_n) \cup f_3(S_n). \end{aligned}$$

Construct S_1 starting from the line segment $L = ((-1, 0), (1, 0))$.

Can you figure out what S_n looks like in general? (You may want to write a computer program.) Do you think that S_n has a limit?

Problem 19.18. In the plane \mathbb{A}^2 , with respect to the standard affine frame, a point of coordinates (x, y) can be represented as the complex number $z = x + iy$. Consider the set of geometric transformations of the form

$$z \mapsto az + b,$$

where a, b are complex numbers such that $a \neq 0$.

- (a) Prove that these maps are affine. Describe what these maps do geometrically.
- (b) Prove that the above set of maps is a group under composition.
- (c) Consider the set of geometric transformations of the form

$$z \mapsto az + b \quad \text{or} \quad z \mapsto a\bar{z} + b,$$

where a, b are complex numbers such that $a \neq 0$, and where $\bar{z} = x - iy$ if $z = x + iy$. Describe what these maps do geometrically. Prove that these maps are affine and that this set of maps is a group under composition.

Problem 19.19. Given a group G , a subgroup H of G is called a *normal subgroup* of G iff $xHx^{-1} = H$ for all $x \in G$ (where $xHx^{-1} = \{xhx^{-1} \mid h \in H\}$).

- (i) Given any two subgroups H and K of a group G , let

$$HK = \{hk \mid h \in H, k \in K\}.$$

Prove that every $x \in HK$ can be written in a unique way as $x = hk$ for $h \in H$ and $k \in K$ iff $H \cap K = \{1\}$, where 1 is the identity element of G .

(ii) If H and K are subgroups of G , and H is a normal subgroup of G , prove that HK is a subgroup of G . Furthermore, if $G = HK$ and $H \cap K = \{1\}$, prove that G is isomorphic to $H \times K$ under the multiplication operation

$$(h_1, k_1) \cdot (h_2, k_2) = (h_1k_1h_2k_1^{-1}, k_1k_2).$$

When $G = HK$, where H, K are subgroups of G , H is a normal subgroup of G , and $H \cap K = \{1\}$, we say that G is the *semidirect product* of H and K .

(iii) Let (E, \vec{E}) be an affine space. Recall that the *affine group* of E , denoted by $\mathbf{GA}(E)$, is the set of affine bijections of E , and that the *linear group* of \vec{E} , denoted by $\mathbf{GL}(\vec{E})$, is the group of bijective linear maps of \vec{E} . The map $f \mapsto \vec{f}$ defines a group homomorphism

$L: \mathbf{GA}(E) \rightarrow \mathbf{GL}(\vec{E})$, and the kernel of this map is the set of translations on E , denoted as $T(E)$. Prove that $T(E)$ is a normal subgroup of $\mathbf{GA}(E)$.

(iv) For any $a \in E$, let

$$\mathbf{GA}_a(E) = \{f \in \mathbf{GA}(E) \mid f(a) = a\},$$

the set of affine bijections leaving a fixed. Prove that $\mathbf{GA}_a(E)$ is a subgroup of $\mathbf{GA}(E)$, and that $\mathbf{GA}_a(E)$ is isomorphic to $\mathbf{GL}(\vec{E})$. Prove that $\mathbf{GA}(E)$ is isomorphic to the direct product of $T(E)$ and $\mathbf{GA}_a(E)$.

Hint. Note that if $u = \overrightarrow{f(a)a}$ and t_u is the translation associated with the vector u , then $t_u \circ f \in \mathbf{GA}_a(E)$ (where the translation t_u is defined such that $t_u(a) = a + u$ for every $a \in E$).

(v) Given a group G , let $\mathbf{Aut}(G)$ denote the set of homomorphisms $f: G \rightarrow G$. Prove that the set $\mathbf{Aut}(G)$ is a group under composition (called the *group of automorphisms of G*). Given any two groups H and K and a homomorphism $\theta: K \rightarrow \mathbf{Aut}(H)$, we define $H \times_\theta K$ as the set $H \times K$ under the multiplication operation

$$(h_1, k_1) \cdot (h_2, k_2) = (h_1 \theta(k_1)(h_2), k_1 k_2).$$

Prove that $H \times_\theta K$ is a group.

Hint. The inverse of (h, k) is $(\theta(k^{-1})(h^{-1}), k^{-1})$.

Prove that the group $H \times_\theta K$ is the semidirect product of the subgroups $\{(h, 1) \mid h \in H\}$ and $\{(1, k) \mid k \in K\}$. The group $H \times_\theta K$ is also called the *semidirect product of H and K relative to θ* .

Note. It is natural to identify $\{(h, 1) \mid h \in H\}$ with H and $\{(1, k) \mid k \in K\}$ with K .

If G is the semidirect product of two subgroups H and K as defined in (ii), prove that the map $\gamma: K \rightarrow \mathbf{Aut}(H)$ defined by conjugation such that

$$\gamma(k)(h) = khk^{-1}$$

is a homomorphism, and that G is isomorphic to $H \times_\gamma K$.

(vi) Define the map $\theta: \mathbf{GL}(\vec{E}) \rightarrow \mathbf{Aut}(\vec{E})$ as follows: $\theta(f) = f$, where $f \in \mathbf{GL}(\vec{E})$ (note that θ can be viewed as an inclusion map). Prove that $\mathbf{GA}(E)$ is isomorphic to the semidirect product $\vec{E} \times_\theta \mathbf{GL}(\vec{E})$.

(vii) Let $\mathbf{SL}(\vec{E})$ be the subgroup of $\mathbf{GL}(\vec{E})$ consisting of the linear maps such that $\det(f) = 1$ (the *special linear group of \vec{E}*), and let $\mathbf{SA}(E)$ be the subgroup of $\mathbf{GA}(E)$ (the *special affine group of E*) consisting of the affine maps f such that $\vec{f} \in \mathbf{SL}(\vec{E})$. Prove that $\mathbf{SA}(E)$ is isomorphic to the semidirect product $\vec{E} \times_\theta \mathbf{SL}(\vec{E})$, where $\theta: \mathbf{SL}(\vec{E}) \rightarrow \mathbf{Aut}(\vec{E})$ is defined as in (vi).

(viii) Assume that (E, \vec{E}) is a Euclidean affine space. Let $\mathbf{SO}(\vec{E})$ be the *special orthogonal group of \vec{E}* (the isometries with determinant +1), and let $\mathbf{SE}(E)$ be the subgroup of $\mathbf{SA}(E)$ (the *special Euclidean group of E*) consisting of the affine isometries f such that $\vec{f} \in \mathbf{SO}(\vec{E})$. Prove that $\mathbf{SE}(E)$ is isomorphic to the semidirect product $\vec{E} \times_\theta \mathbf{SO}(\vec{E})$, where $\theta: \mathbf{SO}(\vec{E}) \rightarrow \mathbf{Aut}(\vec{E})$ is defined as in (vi).

Problem 19.20. The purpose of this problem is to study certain affine maps of \mathbb{A}^2 .

(1) Consider affine maps of the form

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}.$$

Prove that such maps have a unique fixed point c if $\theta \neq 2k\pi$, for all integers k . Show that these are rotations of center c , which means that with respect to a frame with origin c (the unique fixed point), these affine maps are represented by rotation matrices.

(2) Consider affine maps of the form

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} \lambda \cos \theta & -\lambda \sin \theta \\ \mu \sin \theta & \mu \cos \theta \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}.$$

Prove that such maps have a unique fixed point iff $(\lambda + \mu) \cos \theta \neq 1 + \lambda\mu$. Prove that if $\lambda\mu = 1$ and $\lambda > 0$, there is some angle θ for which either there is no fixed point, or there are infinitely many fixed points.

(3) Prove that the affine map

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} 8/5 & -6/5 \\ 3/10 & 2/5 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

has no fixed point.

(4) Prove that an arbitrary affine map

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} a_1 & a_2 \\ a_3 & a_4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

has a unique fixed point iff the matrix

$$\begin{pmatrix} a_1 - 1 & a_2 \\ a_3 & a_4 - 1 \end{pmatrix}$$

is invertible.

Problem 19.21. Let (E, \overrightarrow{E}) be any affine space of finite dimension. For every affine map $f: E \rightarrow E$, let $\text{Fix}(f) = \{a \in E \mid f(a) = a\}$ be the set of fixed points of f .

(i) Prove that if $\text{Fix}(f) \neq \emptyset$, then $\text{Fix}(f)$ is an affine subspace of E such that for every $b \in \text{Fix}(f)$,

$$\text{Fix}(f) = b + \text{Ker}(\overrightarrow{f} - \text{id}).$$

(ii) Prove that $\text{Fix}(f)$ contains a unique fixed point iff $\text{Ker}(\overrightarrow{f} - \text{id}) = \{0\}$, i.e., $\overrightarrow{f}(u) = u$ iff $u = 0$.

Hint. Show that

$$\overrightarrow{\Omega f(a)} - \overrightarrow{\Omega a} = \overrightarrow{\Omega f(\Omega)} + \overrightarrow{f}(\overrightarrow{\Omega a}) - \overrightarrow{\Omega a},$$

for any two points $\Omega, a \in E$.

Problem 19.22. Given two affine spaces (E, \vec{E}) and (F, \vec{F}) , let $\mathcal{A}(E, F)$ be the set of all affine maps $f: E \rightarrow F$.

(i) Prove that the set $\mathcal{A}(E, \vec{F})$ (viewing \vec{F} as an affine space) is a vector space under the operations $f + g$ and λf defined such that

$$\begin{aligned}(f + g)(a) &= f(a) + g(a), \\ (\lambda f)(a) &= \lambda f(a),\end{aligned}$$

for all $a \in E$.

(ii) Define an action

$$+: \mathcal{A}(E, F) \times \mathcal{A}(E, \vec{F}) \rightarrow \mathcal{A}(E, F)$$

of $\mathcal{A}(E, \vec{F})$ on $\mathcal{A}(E, F)$ as follows: For every $a \in E$, every $f \in \mathcal{A}(E, F)$, and every $h \in \mathcal{A}(E, \vec{F})$,

$$(f + h)(a) = f(a) + h(a).$$

Prove that $(\mathcal{A}(E, F), \mathcal{A}(E, \vec{F}), +)$ is an affine space.

Hint. Show that for any two affine maps $f, g \in \mathcal{A}(E, F)$, the map $\vec{f}g$ defined such that

$$\vec{f}g(a) = \overrightarrow{f(a)g(a)}$$

(for every $a \in E$) is affine, and thus $\vec{f}g \in \mathcal{A}(E, \vec{F})$. Furthermore, $\vec{f}g$ is the unique map in $\mathcal{A}(E, \vec{F})$ such that

$$f + \vec{f}g = g.$$

(iii) If \vec{E} has dimension m and \vec{F} has dimension n , prove that $\mathcal{A}(E, \vec{F})$ has dimension $n + mn = n(m + 1)$.

Problem 19.23. Let (c_1, \dots, c_n) be $n \geq 3$ points in \mathbb{A}^m (where $m \geq 2$). Investigate whether there is a closed polygon with n vertices (a_1, \dots, a_n) such that c_i is the middle of the edge (a_i, a_{i+1}) for every i with $1 \leq i \leq n - 1$, and c_n is the middle of the edge (a_n, a_0) .

Hint. The parity (odd or even) of n plays an important role. When n is odd, there is a unique solution, and when n is even, there are no solutions or infinitely many solutions. Clarify under which conditions there are infinitely many solutions.

Problem 19.24. Given an affine space E of dimension n and an affine frame (a_0, \dots, a_n) for E , let $f: E \rightarrow E$ and $g: E \rightarrow E$ be two affine maps represented by the two $(n + 1) \times (n + 1)$ matrices

$$\begin{pmatrix} A & b \\ 0 & 1 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} B & c \\ 0 & 1 \end{pmatrix}$$

w.r.t. the frame (a_0, \dots, a_n) . We also say that f and g are represented by (A, b) and (B, c) .

(1) Prove that the composition $f \circ g$ is represented by the matrix

$$\begin{pmatrix} AB & Ac + b \\ 0 & 1 \end{pmatrix}.$$

We also say that $f \circ g$ is represented by $(A, b)(B, c) = (AB, Ac + b)$.

(2) Prove that f is invertible iff A is invertible and that the matrix representing f^{-1} is

$$\begin{pmatrix} A^{-1} & -A^{-1}b \\ 0 & 1 \end{pmatrix}.$$

We also say that f^{-1} is represented by $(A, b)^{-1} = (A^{-1}, -A^{-1}b)$. Prove that if A is an orthogonal matrix, the matrix associated with f^{-1} is

$$\begin{pmatrix} A^\top & -A^\top b \\ 0 & 1 \end{pmatrix}.$$

Furthermore, denoting the columns of A by A_1, \dots, A_n , prove that the vector $A^\top b$ is the column vector of components

$$(A_1 \cdot b, \dots, A_n \cdot b)$$

(where \cdot denotes the standard inner product of vectors).

(3) Given two affine frames (a_0, \dots, a_n) and (a'_0, \dots, a'_n) for E , any affine map $f: E \rightarrow E$ has a matrix representation (A, b) w.r.t. (a_0, \dots, a_n) and (a'_0, \dots, a'_n) defined such that $b = \overrightarrow{a'_0 f(a_0)}$ is expressed over the basis $(\overrightarrow{a'_0 a'_1}, \dots, \overrightarrow{a'_0 a'_n})$, and a_{ij} is the i th coefficient of $f(\overrightarrow{a_0 a_j})$ over the basis $(\overrightarrow{a'_0 a'_1}, \dots, \overrightarrow{a'_0 a'_n})$. Given any three frames (a_0, \dots, a_n) , (a'_0, \dots, a'_n) , and (a''_0, \dots, a''_n) , for any two affine maps $f: E \rightarrow E$ and $g: E \rightarrow E$, if f has the matrix representation (A, b) w.r.t. (a_0, \dots, a_n) and (a'_0, \dots, a'_n) and g has the matrix representation (B, c) w.r.t. (a'_0, \dots, a'_n) and (a''_0, \dots, a''_n) , prove that $g \circ f$ has the matrix representation $(B, c)(A, b)$ w.r.t. (a_0, \dots, a_n) and (a''_0, \dots, a''_n) .

(4) Given two affine frames (a_0, \dots, a_n) and (a'_0, \dots, a'_n) for E , there is a unique affine map $h: E \rightarrow E$ such that $h(a_i) = a'_i$ for $i = 0, \dots, n$, and we let (P, ω) be its associated matrix representation with respect to the frame (a_0, \dots, a_n) . Note that $\omega = \overrightarrow{a_0 a'_0}$, and that p_{ij} is the i th coefficient of $\overrightarrow{a'_0 a'_j}$ over the basis $(\overrightarrow{a_0 a_1}, \dots, \overrightarrow{a_0 a_n})$. Observe that (P, ω) is also the matrix representation of id_E w.r.t. the frames (a'_0, \dots, a'_n) and (a_0, \dots, a_n) , **in that order**. For any affine map $f: E \rightarrow E$, if f has the matrix representation (A, b) over the frame (a_0, \dots, a_n) and the matrix representation (A', b') over the frame (a'_0, \dots, a'_n) , prove that

$$(A', b') = (P, \omega)^{-1}(A, b)(P, \omega).$$

Given any two affine maps $f: E \rightarrow E$ and $g: E \rightarrow E$, where f is invertible, for any affine frame (a_0, \dots, a_n) for E , if (a'_0, \dots, a'_n) is the affine frame image of (a_0, \dots, a_n) under f (i.e., $f(a_i) = a'_i$ for $i = 0, \dots, n$), letting (A, b) be the matrix representation of f w.r.t. the frame (a_0, \dots, a_n) and (B, c) be the matrix representation of g w.r.t. the frame (a'_0, \dots, a'_n) (**not**

the frame (a_0, \dots, a_n) , prove that $g \circ f$ is represented by the matrix $(A, b)(B, c)$ w.r.t. the frame (a_0, \dots, a_n) .

Remark: Note that this is the **opposite** of what happens if f and g are both represented by matrices w.r.t. the “fixed” frame (a_0, \dots, a_n) , where $g \circ f$ is represented by the matrix $(B, c)(A, b)$. The frame (a'_0, \dots, a'_n) can be viewed as a “moving” frame. The above has applications in robotics, for example to rotation matrices expressed in terms of Euler angles, or “roll, pitch, and yaw.”

Problem 19.25. (a) Let E be a vector space, and let U and V be two subspaces of E so that they form a direct sum $E = U \oplus V$. Recall that this means that every vector $x \in E$ can be written as $x = u + v$, for some unique $u \in U$ and some unique $v \in V$. Define the function $p_U: E \rightarrow U$ (resp. $p_V: E \rightarrow V$) so that $p_U(x) = u$ (resp. $p_V(x) = v$), where $x = u + v$, as explained above. Check that p_U and p_V are linear.

(b) Now assume that E is an affine space (nontrivial), and let U and V be affine subspaces such that $\vec{E} = \vec{U} \oplus \vec{V}$. Pick any $\Omega \in V$, and define $q_U: E \rightarrow \vec{U}$ (resp. $q_V: E \rightarrow \vec{V}$, with $\Omega \in U$) so that

$$q_U(a) = p_{\vec{U}}(\vec{\Omega a}) \quad (\text{resp.} \quad q_V(a) = p_{\vec{V}}(\vec{\Omega a})), \quad \text{for every } a \in E.$$

Prove that q_U does not depend on the choice of $\Omega \in V$ (resp. q_V does not depend on the choice of $\Omega \in U$). Define the map $p_U: E \rightarrow U$ (resp. $p_V: E \rightarrow V$) so that

$$p_U(a) = a - q_V(a) \quad (\text{resp.} \quad p_V(a) = a - q_U(a)), \quad \text{for every } a \in E.$$

Prove that p_U (resp. p_V) is affine.

The map p_U (resp. p_V) is called the *projection onto U parallel to V* (resp. *projection onto V parallel to U*).

(c) Let (a_0, \dots, a_n) be $n+1$ affinely independent points in \mathbb{A}^n and let $\Delta(a_0, \dots, a_n)$ denote the convex hull of (a_0, \dots, a_n) (an n -simplex). Prove that if $f: \mathbb{A}^n \rightarrow \mathbb{A}^n$ is an affine map sending $\Delta(a_0, \dots, a_n)$ inside itself, i.e.,

$$f(\Delta(a_0, \dots, a_n)) \subseteq \Delta(a_0, \dots, a_n),$$

then, f has some fixed point $b \in \Delta(a_0, \dots, a_n)$, i.e., $f(b) = b$.

Hint: Proceed by induction on n . First, treat the case $n = 1$. The affine map is determined by $f(a_0)$ and $f(a_1)$, which are affine combinations of a_0 and a_1 . There is an explicit formula for some fixed point of f . For the induction step, compose f with some suitable projections.

Chapter 20

Polynomials, Ideals and PID's

20.1 Multisets

This chapter contains a review of polynomials and their basic properties. First, multisets are defined. Polynomials in one variable are defined next. The notion of a polynomial function in one argument is defined. Polynomials in several variables are defined, and so is the notion of a polynomial function in several arguments. The Euclidean division algorithm is presented, and the main consequences of its existence are derived. Ideals are defined, and the characterization of greatest common divisors of polynomials in one variable (gcd's) in terms of ideals is shown. We also prove the Bezout identity. Next, we consider the factorization of polynomials in one variable into irreducible factors. The unique factorization of polynomials in one variable into irreducible factors is shown. Roots of polynomials and their multiplicity are defined. It is shown that a nonnull polynomial in one variable and of degree m over an integral domain has at most m roots. The chapter ends with a brief treatment of polynomial interpolation: Lagrange, Newton, and Hermite interpolants are introduced.

In this chapter, it is assumed that all rings considered are commutative. Recall that a (commutative) ring A is an *integral domain* (or an *entire ring*) if $1 \neq 0$, and if $ab = 0$, then either $a = 0$ or $b = 0$, for all $a, b \in A$. This second condition is equivalent to saying that if $a \neq 0$ and $b \neq 0$, then $ab \neq 0$. Also, recall that $a \neq 0$ is *not* a zero divisor if $ab \neq 0$ whenever $b \neq 0$. Observe that a field is an integral domain.

Our goal is to define polynomials in one or more indeterminates (or variables) X_1, \dots, X_n , with coefficients in a ring A . This can be done in several ways, and we choose a definition that has the advantage of extending immediately from one to several variables. First, we need to review the notion of a (finite) multiset.

Definition 20.1. Given a set I , a (*finite*) *multiset over I* is any function $M: I \rightarrow \mathbb{N}$ such that $M(i) \neq 0$ for finitely many $i \in I$. The multiset M such that $M(i) = 0$ for all $i \in I$ is the *empty multiset*, and it is denoted by 0 . If $M(i) = k \neq 0$, we say that i is a *member of M of multiplicity k* . The *union* $M_1 + M_2$ of two multisets M_1 and M_2 is defined such that $(M_1 + M_2)(i) = M_1(i) + M_2(i)$, for every $i \in I$. If I is finite, say $I = \{1, \dots, n\}$, the multiset

M such that $M(i) = k_i$ for every i , $1 \leq i \leq n$, is denoted by $k_1 \cdot 1 + \cdots + k_n \cdot n$, or more simply, by (k_1, \dots, k_n) , and $\deg(k_1 \cdot 1 + \cdots + k_n \cdot n) = k_1 + \cdots + k_n$ is the *size* or *degree* of M . The set of all multisets over I is denoted by $\mathbb{N}^{(I)}$, and when $I = \{1, \dots, n\}$, by $\mathbb{N}^{(n)}$.

Intuitively, the order of the elements of a multiset is irrelevant, but the multiplicity of each element is relevant, contrary to sets. Every $i \in I$ is identified with the multiset M_i such that $M_i(i) = 1$ and $M_i(j) = 0$ for $j \neq i$. When $I = \{1\}$, the set $\mathbb{N}^{(1)}$ of multisets $k \cdot 1$ can be identified with \mathbb{N} and $\{1\}^*$. We will denote $k \cdot 1$ simply by k .



However, beware that when $n \geq 2$, the set $\mathbb{N}^{(n)}$ of multisets cannot be identified with the set of strings in $\{1, \dots, n\}^*$, because multiset union is commutative, but concatenation of strings in $\{1, \dots, n\}^*$ is not commutative when $n \geq 2$. This is because in a multiset $k_1 \cdot 1 + \cdots + k_n \cdot n$, the order is irrelevant, whereas in a string, the order is relevant. For example, $2 \cdot 1 + 3 \cdot 2 = 3 \cdot 2 + 2 \cdot 1$, but $11222 \neq 22211$, as strings over $\{1, 2\}$.

Nevertheless, $\mathbb{N}^{(n)}$ and the set \mathbb{N}^n of ordered n -tuples under component-wise addition are isomorphic under the map

$$k_1 \cdot 1 + \cdots + k_n \cdot n \mapsto (k_1, \dots, k_n).$$

Thus, since the notation (k_1, \dots, k_n) is less cumbersome than $k_1 \cdot 1 + \cdots + k_n \cdot n$, it will be preferred. We just have to remember that the order of the k_i is really irrelevant.



But when I is infinite, beware that $\mathbb{N}^{(I)}$ and the set \mathbb{N}^I of ordered I -tuples are not isomorphic.

We are now ready to define polynomials.

20.2 Polynomials

We begin with polynomials in one variable.

Definition 20.2. Given a ring A , we define the set $\mathcal{P}_A(1)$ of *polynomials over A in one variable* as the set of functions $P: \mathbb{N} \rightarrow A$ such that $P(k) \neq 0$ for finitely many $k \in \mathbb{N}$. The polynomial such that $P(k) = 0$ for all $k \in \mathbb{N}$ is the *null (or zero) polynomial* and it is denoted by 0. We define addition of polynomials, multiplication by a scalar, and multiplication of polynomials, as follows: Given any three polynomials $P, Q, R \in \mathcal{P}_A(1)$, letting $a_k = P(k)$, $b_k = Q(k)$, and $c_k = R(k)$, for every $k \in \mathbb{N}$, we define $R = P + Q$ such that

$$c_k = a_k + b_k,$$

$R = \lambda P$ such that

$$c_k = \lambda a_k,$$

where $\lambda \in A$,

and $R = PQ$ such that

$$c_k = \sum_{i+j=k} a_i b_j.$$

We define the polynomial e_k such that $e_k(k) = 1$ and $e_k(i) = 0$ for $i \neq k$. We also denote e_0 by 1 when $k = 0$. Given a polynomial P , the $a_k = P(k) \in A$ are called the *coefficients of P* . If P is not the null polynomial, there is a greatest $n \geq 0$ such that $a_n \neq 0$ (and thus, $a_k = 0$ for all $k > n$) called the *degree of P* and denoted by $\deg(P)$. Then, P is written uniquely as

$$P = a_0 e_0 + a_1 e_1 + \cdots + a_n e_n.$$

When P is the null polynomial, we let $\deg(P) = -\infty$.

There is an injection of A into $\mathcal{P}_A(1)$ given by the map $a \mapsto a1$ (recall that 1 denotes e_0). There is also an injection of \mathbb{N} into $\mathcal{P}_A(1)$ given by the map $k \mapsto e_k$. Observe that $e_k = e_1^k$ (with $e_1^0 = e_0 = 1$). In order to alleviate the notation, we often denote e_1 by X , and we call X a *variable (or indeterminate)*. Then, $e_k = e_1^k$ is denoted by X^k . Adopting this notation, given a nonnull polynomial P of degree n , if $P(k) = a_k$, P is denoted by

$$P = a_0 + a_1 X + \cdots + a_n X^n,$$

or by

$$P = a_n X^n + a_{n-1} X^{n-1} + \cdots + a_0,$$

if this is more convenient (the order of the terms does not matter anyway). Sometimes, it will also be convenient to write a polynomial as

$$P = a_0 X^n + a_1 X^{n-1} + \cdots + a_n.$$

The set $\mathcal{P}_A(1)$ is also denoted by $A[X]$ and a polynomial P may be denoted by $P(X)$. In denoting polynomials, we will use both upper-case and lower-case letters, usually, P, Q, R, S, p, q, r, s , but also f, g, h , etc., if needed (as long as no ambiguities arise).

Given a nonnull polynomial P of degree n , the nonnull coefficient a_n is called the *leading coefficient of P* . The coefficient a_0 is called the *constant term of P* . A polynomial of the form $a_k X^k$ is called a *monomial*. We say that $a_k X^k$ *occurs in P* if $a_k \neq 0$. A nonzero polynomial P of degree n is called a *monic polynomial (or unitary polynomial, or monic)* if $a_n = 1$, where a_n is its leading coefficient, and such a polynomial can be written as

$$P = X^n + a_{n-1} X^{n-1} + \cdots + a_0 \quad \text{or} \quad P = X^n + a_1 X^{n-1} + \cdots + a_n.$$



The choice of the variable X to denote e_1 is standard practice, but there is nothing special about X . We could have chosen Y, Z , or any other symbol, as long as no ambiguities arise.

Formally, the definition of $\mathcal{P}_A(1)$ has nothing to do with X . The reason for using X is simply convenience. Indeed, it is more convenient to write a polynomial as $P = a_0 + a_1X + \cdots + a_nX^n$ rather than as $P = a_0e_0 + a_1e_1 + \cdots + a_ne_n$.

We have the following simple but crucial proposition.

Proposition 20.1. *Given two nonnull polynomials $P(X) = a_0 + a_1X + \cdots + a_mX^m$ of degree m and $Q(X) = b_0 + b_1X + \cdots + b_nX^n$ of degree n , if either a_m or b_n is not a zero divisor, then $a_mb_n \neq 0$, and thus, $PQ \neq 0$ and*

$$\deg(PQ) = \deg(P) + \deg(Q).$$

In particular, if A is an integral domain, then $A[X]$ is an integral domain.

Proof. Since the coefficient of X^{m+n} in PQ is a_mb_n , and since we assumed that either a_m or b_n is not a zero divisor, we have $a_mb_n \neq 0$, and thus, $PQ \neq 0$ and

$$\deg(PQ) = \deg(P) + \deg(Q).$$

Then, it is obvious that $A[X]$ is an integral domain. □

It is easily verified that $A[X]$ is a commutative ring, with multiplicative identity $1X^0 = 1$. It is also easily verified that $A[X]$ satisfies all the conditions of Definition 2.9, but $A[X]$ is not a vector space, since A is not necessarily a field.

A structure satisfying the axioms of Definition 2.9 when K is a ring (and not necessarily a field) is called a *module*. As we mentioned in Section 4.2, we will not study modules because they fail to have some of the nice properties that vector spaces have, and thus, they are harder to study. For example, there are modules that do not have a basis.

However, when the ring A is a field, $A[X]$ is a vector space. But even when A is just a ring, the family of polynomials $(X^k)_{k \in \mathbb{N}}$ is a basis of $A[X]$, since every polynomial $P(X)$ can be written in a unique way as $P(X) = a_0 + a_1X + \cdots + a_nX^n$ (with $P(X) = 0$ when $P(X)$ is the null polynomial). Thus, $A[X]$ is a free module.

Next, we want to define the notion of evaluating a polynomial $P(X)$ at some $\alpha \in A$. For this, we need a proposition.

Proposition 20.2. *Let A, B be two rings and let $h: A \rightarrow B$ be a ring homomorphism. For any $\beta \in B$, there is a unique ring homomorphism $\varphi: A[X] \rightarrow B$ extending h such that $\varphi(X) = \beta$, as in the following diagram (where we denote by $h+\beta$ the map $h+\beta: A \cup \{X\} \rightarrow B$ such that $(h+\beta)(a) = h(a)$ for all $a \in A$ and $(h+\beta)(X) = \beta$):*

$$\begin{array}{ccc} A \cup \{X\} & \xrightarrow{\iota} & A[X] \\ & \searrow h+\beta & \downarrow \varphi \\ & & B \end{array}$$

Proof. Let $\varphi(0) = 0$, and for every nonnull polynomial $P(X) = a_0 + a_1X + \cdots + a_nX^n$, let

$$\varphi(P(X)) = h(a_0) + h(a_1)\beta + \cdots + h(a_n)\beta^n.$$

It is easily verified that φ is the unique homomorphism $\varphi: A[X] \rightarrow B$ extending h such that $\varphi(X) = \beta$. \square

Taking $A = B$ in Proposition 20.2 and $h: A \rightarrow A$ the identity, for every $\beta \in A$, there is a unique homomorphism $\varphi_\beta: A[X] \rightarrow A$ such that $\varphi_\beta(X) = \beta$, and for every polynomial $P(X)$, we write $\varphi_\beta(P(X))$ as $P(\beta)$ and we call $P(\beta)$ the *value of $P(X)$ at $X = \beta$* . Thus, we can define a function $P_A: A \rightarrow A$ such that $P_A(\beta) = P(\beta)$, for all $\beta \in A$. This function is called the *polynomial function induced by P* .

More generally, P_B can be defined for any (commutative) ring B such that $A \subseteq B$. In general, it is possible that $P_A = Q_A$ for distinct polynomials P, Q . We will see shortly conditions for which the map $P \mapsto P_A$ is injective. In particular, this is true for $A = \mathbb{R}$ (in general, any infinite integral domain). We now define polynomials in n variables.

Definition 20.3. Given $n \geq 1$ and a ring A , the set $\mathcal{P}_A(n)$ of *polynomials over A in n variables* is the set of functions $P: \mathbb{N}^{(n)} \rightarrow A$ such that $P(k_1, \dots, k_n) \neq 0$ for finitely many $(k_1, \dots, k_n) \in \mathbb{N}^{(n)}$. The polynomial such that $P(k_1, \dots, k_n) = 0$ for all (k_1, \dots, k_n) is the *null (or zero) polynomial* and it is denoted by 0. We define addition of polynomials, multiplication by a scalar, and multiplication of polynomials, as follows: Given any three polynomials $P, Q, R \in \mathcal{P}_A(n)$, letting $a_{(k_1, \dots, k_n)} = P(k_1, \dots, k_n)$, $b_{(k_1, \dots, k_n)} = Q(k_1, \dots, k_n)$, $c_{(k_1, \dots, k_n)} = R(k_1, \dots, k_n)$, for every $(k_1, \dots, k_n) \in \mathbb{N}^{(n)}$, we define $R = P + Q$ such that

$$c_{(k_1, \dots, k_n)} = a_{(k_1, \dots, k_n)} + b_{(k_1, \dots, k_n)},$$

$R = \lambda P$, where $\lambda \in A$, such that

$$c_{(k_1, \dots, k_n)} = \lambda a_{(k_1, \dots, k_n)},$$

and $R = PQ$, such that

$$c_{(k_1, \dots, k_n)} = \sum_{(i_1, \dots, i_n) + (j_1, \dots, j_n) = (k_1, \dots, k_n)} a_{(i_1, \dots, i_n)} b_{(j_1, \dots, j_n)}.$$

For every $(k_1, \dots, k_n) \in \mathbb{N}^{(n)}$, we let $e_{(k_1, \dots, k_n)}$ be the polynomial such that

$$e_{(k_1, \dots, k_n)}(k_1, \dots, k_n) = 1 \quad \text{and} \quad e_{(k_1, \dots, k_n)}(h_1, \dots, h_n) = 0,$$

for $(h_1, \dots, h_n) \neq (k_1, \dots, k_n)$. We also denote $e_{(0, \dots, 0)}$ by 1. Given a polynomial P , the $a_{(k_1, \dots, k_n)} = P(k_1, \dots, k_n) \in A$, are called the *coefficients of P* . If P is not the null polynomial, there is a greatest $d \geq 0$ such that $a_{(k_1, \dots, k_n)} \neq 0$ for some $(k_1, \dots, k_n) \in \mathbb{N}^{(n)}$, with $d = k_1 + \cdots + k_n$, called the *total degree of P* and denoted by $\deg(P)$. Then, P is written uniquely as

$$P = \sum_{(k_1, \dots, k_n) \in \mathbb{N}^{(n)}} a_{(k_1, \dots, k_n)} e_{(k_1, \dots, k_n)}.$$

When P is the null polynomial, we let $\deg(P) = -\infty$.

There is an injection of A into $\mathcal{P}_A(n)$ given by the map $a \mapsto a1$ (where 1 denotes $e_{(0,\dots,0)}$). There is also an injection of $\mathbb{N}^{(n)}$ into $\mathcal{P}_A(n)$ given by the map $(h_1, \dots, h_n) \mapsto e_{(h_1,\dots,h_n)}$. Note that $e_{(h_1,\dots,h_n)}e_{(k_1,\dots,k_n)} = e_{(h_1+k_1,\dots,h_n+k_n)}$. In order to alleviate the notation, let X_1, \dots, X_n be n distinct variables and denote $e_{(0,\dots,0,1,0,\dots,0)}$, where 1 occurs in the position i , by X_i (where $1 \leq i \leq n$). With this convention, in view of $e_{(h_1,\dots,h_n)}e_{(k_1,\dots,k_n)} = e_{(h_1+k_1,\dots,h_n+k_n)}$, the polynomial $e_{(k_1,\dots,k_n)}$ is denoted by $X_1^{k_1} \cdots X_n^{k_n}$ (with $e_{(0,\dots,0)} = X_1^0 \cdots X_n^0 = 1$) and it is called a *primitive monomial*. Then, P is also written as

$$P = \sum_{(k_1,\dots,k_n) \in \mathbb{N}^{(n)}} a_{(k_1,\dots,k_n)} X_1^{k_1} \cdots X_n^{k_n}.$$

We also denote $\mathcal{P}_A(n)$ by $A[X_1, \dots, X_n]$. A polynomial $P \in A[X_1, \dots, X_n]$ is also denoted by $P(X_1, \dots, X_n)$.

As in the case $n = 1$, there is nothing special about the choice of X_1, \dots, X_n as variables (or indeterminates). It is just a convenience. After all, the construction of $\mathcal{P}_A(n)$ has nothing to do with X_1, \dots, X_n .

Given a nonnull polynomial P of degree d , the nonnull coefficients $a_{(k_1,\dots,k_n)} \neq 0$ such that $d = k_1 + \cdots + k_n$ are called the *leading coefficients of P* . A polynomial of the form $a_{(k_1,\dots,k_n)} X_1^{k_1} \cdots X_n^{k_n}$ is called a *monomial*. Note that $\deg(a_{(k_1,\dots,k_n)} X_1^{k_1} \cdots X_n^{k_n}) = k_1 + \cdots + k_n$. Given a polynomial

$$P = \sum_{(k_1,\dots,k_n) \in \mathbb{N}^{(n)}} a_{(k_1,\dots,k_n)} X_1^{k_1} \cdots X_n^{k_n},$$

a monomial $a_{(k_1,\dots,k_n)} X_1^{k_1} \cdots X_n^{k_n}$ occurs in the polynomial P if $a_{(k_1,\dots,k_n)} \neq 0$.

A polynomial

$$P = \sum_{(k_1,\dots,k_n) \in \mathbb{N}^{(n)}} a_{(k_1,\dots,k_n)} X_1^{k_1} \cdots X_n^{k_n}$$

is *homogeneous of degree d* if

$$\deg(X_1^{k_1} \cdots X_n^{k_n}) = d,$$

for every monomial $a_{(k_1,\dots,k_n)} X_1^{k_1} \cdots X_n^{k_n}$ occurring in P . If P is a polynomial of total degree d , it is clear that P can be written uniquely as

$$P = P^{(0)} + P^{(1)} + \cdots + P^{(d)},$$

where $P^{(i)}$ is the sum of all monomials of degree i occurring in P , where $0 \leq i \leq d$.

It is easily verified that $A[X_1, \dots, X_n]$ is a commutative ring, with multiplicative identity $1X_1^0 \cdots X_n^0 = 1$. It is also easily verified that $A[X]$ is a module. When A is a field, $A[X]$ is a vector space.

Even when A is just a ring, the family of polynomials

$$(X_1^{k_1} \cdots X_n^{k_n})_{(k_1,\dots,k_n) \in \mathbb{N}^{(n)}}$$

is a basis of $A[X_1, \dots, X_n]$, since every polynomial $P(X_1, \dots, X_n)$ can be written in a unique way as

$$P(X_1, \dots, X_n) = \sum_{(k_1, \dots, k_n) \in \mathbb{N}^{(n)}} a_{(k_1, \dots, k_n)} X_1^{k_1} \cdots X_n^{k_n}.$$

Thus, $A[X_1, \dots, X_n]$ is a free module.

Remark: The construction of Definition 20.3 can be immediately extended to an arbitrary set I , and not just $I = \{1, \dots, n\}$. It can also be applied to monoids more general than $\mathbb{N}^{(I)}$.

Proposition 20.2 is generalized as follows.

Proposition 20.3. *Let A, B be two rings and let $h: A \rightarrow B$ be a ring homomorphism. For any $\beta = (\beta_1, \dots, \beta_n) \in B^n$, there is a unique ring homomorphism $\varphi: A[X_1, \dots, X_n] \rightarrow B$ extending h such that $\varphi(X_i) = \beta_i$, $1 \leq i \leq n$, as in the following diagram (where we denote by $h + \beta$ the map $h + \beta: A \cup \{X_1, \dots, X_n\} \rightarrow B$ such that $(h + \beta)(a) = h(a)$ for all $a \in A$ and $(h + \beta)(X_i) = \beta_i$, $1 \leq i \leq n$):*

$$\begin{array}{ccc} A \cup \{X_1, \dots, X_n\} & \xrightarrow{\iota} & A[X_1, \dots, X_n] \\ & \searrow h + \beta & \downarrow \varphi \\ & & B \end{array}$$

Proof. Let $\varphi(0) = 0$, and for every nonnull polynomial

$$P(X_1, \dots, X_n) = \sum_{(k_1, \dots, k_n) \in \mathbb{N}^{(n)}} a_{(k_1, \dots, k_n)} X_1^{k_1} \cdots X_n^{k_n},$$

let

$$\varphi(P(X_1, \dots, X_n)) = \sum h(a_{(k_1, \dots, k_n)}) \beta_1^{k_1} \cdots \beta_n^{k_n}.$$

It is easily verified that φ is the unique homomorphism $\varphi: A[X_1, \dots, X_n] \rightarrow B$ extending h such that $\varphi(X_i) = \beta_i$. \square

Taking $A = B$ in Proposition 20.3 and $h: A \rightarrow A$ the identity, for every $\beta_1, \dots, \beta_n \in A$, there is a unique homomorphism $\varphi: A[X_1, \dots, X_n] \rightarrow A$ such that $\varphi(X_i) = \beta_i$, and for every polynomial $P(X_1, \dots, X_n)$, we write $\varphi(P(X_1, \dots, X_n))$ as $P(\beta_1, \dots, \beta_n)$ and we call $P(\beta_1, \dots, \beta_n)$ the *value of $P(X_1, \dots, X_n)$ at $X_1 = \beta_1, \dots, X_n = \beta_n$* . Thus, we can define a function $P_A: A^n \rightarrow A$ such that $P_A(\beta_1, \dots, \beta_n) = P(\beta_1, \dots, \beta_n)$, for all $\beta_1, \dots, \beta_n \in A$. This function is called the *polynomial function induced by P* .

More generally, P_B can be defined for any (commutative) ring B such that $A \subseteq B$. As in the case of a single variable, it is possible that $P_A = Q_A$ for distinct polynomials P, Q . We will see shortly that the map $P \mapsto P_A$ is injective when $A = \mathbb{R}$ (in general, any infinite integral domain).

Given any nonnull polynomial $P(X_1, \dots, X_n) = \sum_{(k_1, \dots, k_n) \in \mathbb{N}^{(n)}} a_{(k_1, \dots, k_n)} X_1^{k_1} \cdots X_n^{k_n}$ in $A[X_1, \dots, X_n]$, where $n \geq 2$, $P(X_1, \dots, X_n)$ can be uniquely written as

$$P(X_1, \dots, X_n) = \sum Q_{k_n}(X_1, \dots, X_{n-1}) X_n^{k_n},$$

where each polynomial $Q_{k_n}(X_1, \dots, X_{n-1})$ is in $A[X_1, \dots, X_{n-1}]$. Thus, even if A is a field, $A[X_1, \dots, X_{n-1}]$ is not a field, which confirms that it is useful (and necessary!) to consider polynomials over rings that are not necessarily fields.

It is not difficult to show that $A[X_1, \dots, X_n]$ and $A[X_1, \dots, X_{n-1}][X_n]$ are isomorphic rings. This way, it is often possible to prove properties of polynomials in several variables X_1, \dots, X_n , by induction on the number n of variables. For example, given two nonnull polynomials $P(X_1, \dots, X_n)$ of total degree p and $Q(X_1, \dots, X_n)$ of total degree q , since we assumed that A is an integral domain, we can prove that

$$\deg(PQ) = \deg(P) + \deg(Q),$$

and that $A[X_1, \dots, X_n]$ is an integral domain.

Next, we will consider the division of polynomials (in one variable).

20.3 Euclidean Division of Polynomials

We know that every natural number $n \geq 2$ can be written uniquely as a product of powers of prime numbers and that prime numbers play a very important role in arithmetic. It would be nice if every polynomial could be expressed (uniquely) as a product of “irreducible” factors. This is indeed the case for polynomials over a field. The fact that there is a division algorithm for the natural numbers is essential for obtaining many of the arithmetical properties of the natural numbers. As we shall see next, there is also a division algorithm for polynomials in $A[X]$, when A is a field.

Proposition 20.4. *Let A be a ring, let $f(X), g(X) \in A[X]$ be two polynomials of degree $m = \deg(f)$ and $n = \deg(g)$ with $f(X) \neq 0$, and assume that the leading coefficient a_m of $f(X)$ is invertible. Then, there exist unique polynomials $q(X)$ and $r(X)$ in $A[X]$ such that*

$$g = fq + r \quad \text{and} \quad \deg(r) < \deg(f) = m.$$

Proof. We first prove the existence of q and r . Let

$$f = a_m X^m + a_{m-1} X^{m-1} + \cdots + a_0,$$

and

$$g = b_n X^n + b_{n-1} X^{n-1} + \cdots + b_0.$$

If $n < m$, then let $q = 0$ and $r = g$. Since $\deg(g) < \deg(f)$ and $r = g$, we have $\deg(r) < \deg(f)$.

If $n \geq m$, we proceed by induction on n . If $n = 0$, then $g = b_0$, $m = 0$, $f = a_0 \neq 0$, and we let $q = a_0^{-1}b_0$ and $r = 0$. Since $\deg(r) = \deg(0) = -\infty$ and $\deg(f) = \deg(a_0) = 0$ because $a_0 \neq 0$, we have $\deg(r) < \deg(f)$.

If $n \geq 1$, since $n \geq m$, note that

$$\begin{aligned} g_1(X) &= g(X) - b_n a_m^{-1} X^{n-m} f(X) \\ &= b_n X^n + b_{n-1} X^{n-1} + \cdots + b_0 - b_n a_m^{-1} X^{n-m} (a_m X^m + a_{m-1} X^{m-1} + \cdots + a_0) \end{aligned}$$

is a polynomial of degree $\deg(g_1) < n$, since the terms $b_n X^n$ and $b_n a_m^{-1} X^{n-m} a_m X^m$ of degree n cancel out. Now, since $\deg(g_1) < n$, by the induction hypothesis, we can find q_1 and r such that

$$g_1 = f q_1 + r \quad \text{and} \quad \deg(r) < \deg(f) = m,$$

and thus,

$$g_1(X) = g(X) - b_n a_m^{-1} X^{n-m} f(X) = f(X) q_1(X) + r(X),$$

from which, letting $q(X) = b_n a_m^{-1} X^{n-m} + q_1(X)$, we get

$$g = f q + r \quad \text{and} \quad \deg(r) < m = \deg(f).$$

We now prove uniqueness. If

$$g = f q_1 + r_1 = f q_2 + r_2,$$

with $\deg(r_1) < \deg(f)$ and $\deg(r_2) < \deg(f)$, we get

$$f(q_1 - q_2) = r_2 - r_1.$$

If $q_2 - q_1 \neq 0$, since the leading coefficient a_m of f is invertible, by Proposition 20.1, we have

$$\deg(r_2 - r_1) = \deg(f(q_1 - q_2)) = \deg(f) + \deg(q_2 - q_1),$$

and so, $\deg(r_2 - r_1) \geq \deg(f)$, which contradicts the fact that $\deg(r_1) < \deg(f)$ and $\deg(r_2) < \deg(f)$. Thus, $q_1 = q_2$, and then also $r_1 = r_2$. \square

It should be noted that the proof of Proposition 20.4 actually provides an algorithm for finding the *quotient* q and the *remainder* r of the division of g by f . This algorithm is called the *Euclidean algorithm*, or *division algorithm*. Note that the division of g by f is always possible when f is a monic polynomial, since 1 is invertible. Also, when A is a field, $a_m \neq 0$ is always invertible, and thus, the division can always be performed. We say that f *divides* g when $r = 0$ in the result of the division $g = f q + r$. We now draw some important consequences of the existence of the Euclidean algorithm.

20.4 Ideals, PID's, and Greatest Common Divisors

First, we introduce the fundamental concept of an ideal.

Definition 20.4. Given a ring A , an *ideal* of A is any nonempty subset \mathfrak{I} of A satisfying the following two properties:

(ID1) If $a, b \in \mathfrak{I}$, then $b - a \in \mathfrak{I}$.

(ID2) If $a \in \mathfrak{I}$, then $ax \in \mathfrak{I}$ for every $x \in A$.

An ideal \mathfrak{I} is a *principal ideal* if there is some $a \in \mathfrak{I}$, called a *generator*, such that

$$\mathfrak{I} = \{ax \mid x \in A\}.$$

The equality $\mathfrak{I} = \{ax \mid x \in A\}$ is also written as $\mathfrak{I} = aA$ or as $\mathfrak{I} = (a)$. The ideal $\mathfrak{I} = (0) = \{0\}$ is called the *null ideal* (or *zero ideal*).

An ideal \mathfrak{I} is a *maximal ideal* if $\mathfrak{I} \neq A$ and for every ideal $\mathfrak{J} \neq A$, if $\mathfrak{I} \subseteq \mathfrak{J}$, then $\mathfrak{J} = \mathfrak{I}$. An ideal \mathfrak{I} is a *prime ideal* if $\mathfrak{I} \neq A$ and if $ab \in \mathfrak{I}$, then $a \in \mathfrak{I}$ or $b \in \mathfrak{I}$, for all $a, b \in A$. Equivalently, \mathfrak{I} is a prime ideal if $\mathfrak{I} \neq A$ and if $a, b \in A - \mathfrak{I}$, then $ab \in A - \mathfrak{I}$, for all $a, b \in A$. In other words, $A - \mathfrak{I}$ is closed under multiplication and $1 \in A - \mathfrak{I}$.

Note that if \mathfrak{I} is an ideal, then $\mathfrak{I} = A$ iff $1 \in \mathfrak{I}$. Since by definition, an ideal \mathfrak{I} is nonempty, there is some $a \in \mathfrak{I}$, and by (ID1) we get $0 = a - a \in \mathfrak{I}$. Then, for every $a \in \mathfrak{I}$, since $0 \in \mathfrak{I}$, by (ID1) we get $-a \in \mathfrak{I}$. Thus, an ideal is an additive subgroup of A . Because of (ID2), an ideal is also a subring.

Observe that if A is a field, then A only has two ideals, namely, the trivial ideal (0) and A itself. Indeed, if $\mathfrak{I} \neq (0)$, because every nonnull element has an inverse, then $1 \in \mathfrak{I}$, and thus, $\mathfrak{I} = A$.

Given $a, b \in A$, we say that b is a *multiple of a* and that a *divides b* if $b = ac$ for some $c \in A$; this is usually denoted by $a \mid b$. Note that the principal ideal (a) is the set of all multiples of a , and that a divides b iff b is a multiple of a iff $b \in (a)$ iff $(b) \subseteq (a)$.

Note that every $a \in A$ divides 0. However, it is customary to say that a is a *zero divisor* iff $ac = 0$ for some $c \neq 0$. With this convention, 0 is a zero divisor unless $A = \{0\}$ (the trivial ring), and A is an integral domain iff 0 is the only zero divisor in A .

Given $a, b \in A$ with $a, b \neq 0$, if $(a) = (b)$ then there exist $c, d \in A$ such that $a = bc$ and $b = ad$. From this, we get $a = adc$ and $b = bcd$, that is, $a(1 - dc) = 0$ and $b(1 - cd) = 0$. If A is an integral domain, we get $dc = 1$ and $cd = 1$, that is, c is invertible with inverse d . Thus, when A is an integral domain, we have $b = ad$, with d invertible. The converse is obvious, if $b = ad$ with d invertible, then $(a) = (b)$.

As a summary, if A is an integral domain, for any $a, b \in A$ with $a, b \neq 0$, we have $(a) = (b)$ iff there exists some invertible $d \in A$ such that $b = ad$. An invertible element $u \in A$ is also called a *unit*.

Given two ideals \mathfrak{I} and \mathfrak{J} , their sum

$$\mathfrak{I} + \mathfrak{J} = \{a + b \mid a \in \mathfrak{I}, b \in \mathfrak{J}\}$$

is clearly an ideal. Given any nonempty subset J of A , the set

$$\{a_1x_1 + \cdots + a_nx_n \mid x_1, \dots, x_n \in A, a_1, \dots, a_n \in J, n \geq 1\}$$

is easily seen to be an ideal, and in fact, it is the smallest ideal containing J . It is usually denoted by (J) .

Ideals play a very important role in the study of rings. They tend to show up everywhere. For example, they arise naturally from homomorphisms.

Proposition 20.5. *Given any ring homomorphism $h: A \rightarrow B$, the kernel $\text{Ker } h = \{a \in A \mid h(a) = 0\}$ of h is an ideal.*

Proof. Given $a, b \in A$, we have $a, b \in \text{Ker } h$ iff $h(a) = h(b) = 0$, and since h is a homomorphism, we get

$$h(b - a) = h(b) - h(a) = 0,$$

and

$$h(ax) = h(a)h(x) = 0$$

for all $x \in A$, which shows that $\text{Ker } h$ is an ideal. □

There is a sort of converse property. Given a ring A and an ideal $\mathfrak{I} \subseteq A$, we can define the quotient ring A/\mathfrak{I} , and there is a surjective homomorphism $\pi: A \rightarrow A/\mathfrak{I}$ whose kernel is precisely \mathfrak{I} .

Proposition 20.6. *Given any ring A and any ideal $\mathfrak{I} \subseteq A$, the equivalence relation $\equiv_{\mathfrak{I}}$ defined by $a \equiv_{\mathfrak{I}} b$ iff $b - a \in \mathfrak{I}$ is a congruence, which means that if $a_1 \equiv_{\mathfrak{I}} b_1$ and $a_2 \equiv_{\mathfrak{I}} b_2$, then*

1. $a_1 + a_2 \equiv_{\mathfrak{I}} b_1 + b_2$, and
2. $a_1a_2 \equiv_{\mathfrak{I}} b_1b_2$.

Then, the set A/\mathfrak{I} of equivalence classes modulo \mathfrak{I} is a ring under the operations

$$\begin{aligned} [a] + [b] &= [a + b] \\ [a][b] &= [ab]. \end{aligned}$$

The map $\pi: A \rightarrow A/\mathfrak{I}$ such that $\pi(a) = [a]$ is a surjective homomorphism whose kernel is precisely \mathfrak{I} .

Proof. Everything is straightforward. For example, if $a_1 \equiv_{\mathfrak{J}} b_1$ and $a_2 \equiv_{\mathfrak{J}} b_2$, then $b_1 - a_1 \in \mathfrak{J}$ and $b_2 - a_2 \in \mathfrak{J}$. Since \mathfrak{J} is an ideal, we get

$$(b_1 - a_1)b_2 = b_1b_2 - a_1b_2 \in \mathfrak{J}$$

and

$$(b_2 - a_2)a_1 = a_1b_2 - a_1a_2 \in \mathfrak{J}.$$

Since \mathfrak{J} is an ideal, and thus, an additive group, we get

$$b_1b_2 - a_1a_2 \in \mathfrak{J},$$

i.e., $a_1a_2 \equiv_{\mathfrak{J}} b_1b_2$. The equality $\text{Ker } \pi = \mathfrak{J}$ holds because \mathfrak{J} is an ideal. \square

Example 20.1.

1. In the ring \mathbb{Z} , for every $p \in \mathbb{Z}$, the subgroup $p\mathbb{Z}$ is an ideal, and $\mathbb{Z}/p\mathbb{Z}$ is a ring, the ring of residues modulo p . This ring is a field iff p is a prime number.
2. The quotient of the polynomial ring $\mathbb{R}[X]$ by a prime ideal \mathfrak{J} is an integral domain.
3. The quotient of the polynomial ring $\mathbb{R}[X]$ by a maximal ideal \mathfrak{J} is a field. For example, if $\mathfrak{J} = (X^2 + 1)$, the principal ideal generated by $X^2 + 1$ (which is indeed a maximal ideal since $X^2 + 1$ has no real roots), then $\mathbb{R}[X]/(X^2 + 1) \cong \mathbb{C}$.

The following proposition yields a characterization of prime ideals and maximal ideals in terms of quotients.

Proposition 20.7. *Given a ring A , for any ideal $\mathfrak{J} \subseteq A$, the following properties hold.*

- (1) *The ideal \mathfrak{J} is a prime ideal iff A/\mathfrak{J} is an integral domain.*
- (2) *The ideal \mathfrak{J} is a maximal ideal iff A/\mathfrak{J} is a field.*

Proof. (1) Assume that \mathfrak{J} is a prime ideal. Since \mathfrak{J} is prime, $\mathfrak{J} \neq A$, and thus, A/\mathfrak{J} is not the trivial ring (0). If $[a][b] = 0$, since $[a][b] = [ab]$, we have $ab \in \mathfrak{J}$, and since \mathfrak{J} is prime, then either $a \in \mathfrak{J}$ or $b \in \mathfrak{J}$, so that either $[a] = 0$ or $[b] = 0$. Thus, A/\mathfrak{J} is an integral domain.

Conversely, assume that A/\mathfrak{J} is an integral domain. Since A/\mathfrak{J} is not the trivial ring, $\mathfrak{J} \neq A$. Assume that $ab \in \mathfrak{J}$. Then, we have

$$\pi(ab) = \pi(a)\pi(b) = 0,$$

which implies that either $\pi(a) = 0$ or $\pi(b) = 0$, since A/\mathfrak{J} is an integral domain (where $\pi: A \rightarrow A/\mathfrak{J}$ is the quotient map). Thus, either $a \in \mathfrak{J}$ or $b \in \mathfrak{J}$, and \mathfrak{J} is a prime ideal.

(2) Assume that \mathfrak{I} is a maximal ideal. As in (1), A/\mathfrak{I} is not the trivial ring (0). Let $[a] \neq 0$ in A/\mathfrak{I} . We need to prove that $[a]$ has a multiplicative inverse. Since $[a] \neq 0$, we have $a \notin \mathfrak{I}$. Let \mathfrak{I}_a be the ideal generated by \mathfrak{I} and a . We have

$$\mathfrak{I} \subseteq \mathfrak{I}_a \quad \text{and} \quad \mathfrak{I} \neq \mathfrak{I}_a,$$

since $a \notin \mathfrak{I}$, and since \mathfrak{I} is maximal, this implies that

$$\mathfrak{I}_a = A.$$

However, we know that

$$\mathfrak{I}_a = \{ax + h \mid x \in A, h \in \mathfrak{I}\},$$

and thus, there is some $x \in A$ so that

$$ax + h = 1,$$

which proves that $[a][x] = 1$, as desired.

Conversely, assume that A/\mathfrak{I} is a field. Again, since A/\mathfrak{I} is not the trivial ring, $\mathfrak{I} \neq A$. Let \mathfrak{J} be any proper ideal such that $\mathfrak{I} \subseteq \mathfrak{J}$, and assume that $\mathfrak{I} \neq \mathfrak{J}$. Thus, there is some $j \in \mathfrak{J} - \mathfrak{I}$, and since $\text{Ker } \pi = \mathfrak{I}$, we have $\pi(j) \neq 0$. Since A/\mathfrak{I} is a field and π is surjective, there is some $k \in A$ so that $\pi(j)\pi(k) = 1$, which implies that

$$jk - 1 = i$$

for some $i \in \mathfrak{I}$, and since $\mathfrak{I} \subset \mathfrak{J}$ and \mathfrak{J} is an ideal, it follows that $1 = jk - i \in \mathfrak{J}$, showing that $\mathfrak{J} = A$, a contradiction. Therefore, $\mathfrak{I} = \mathfrak{J}$, and \mathfrak{I} is a maximal ideal. \square

As a corollary, we obtain the following useful result. It emphasizes the importance of maximal ideals.

Corollary 20.8. *Given any ring A , every maximal ideal \mathfrak{I} in A is a prime ideal.*

Proof. If \mathfrak{I} is a maximal ideal, then, by Proposition 20.7, the quotient ring A/\mathfrak{I} is a field. However, a field is an integral domain, and by Proposition 20.7 (again), \mathfrak{I} is a prime ideal. \square

Observe that a ring A is an integral domain iff (0) is a prime ideal. This is an example of a prime ideal which is not a maximal ideal, as immediately seen in $A = \mathbb{Z}$, where (p) is a maximal ideal for every prime number p .



A less obvious example of a prime ideal which is not a maximal ideal, is the ideal (X) in the ring of polynomials $\mathbb{Z}[X]$. Indeed, $(X, 2)$ is also a prime ideal, but (X) is properly contained in $(X, 2)$.

Definition 20.5. An integral domain in which every ideal is a principal ideal is called a *principal ring or principal ideal domain*, for short, a *PID*.

The ring \mathbb{Z} is a PID. This is a consequence of the existence of a (Euclidean) division algorithm. As we shall see next, when K is a field, the ring $K[X]$ is also a principal ring.



However, when $n \geq 2$, the ring $K[X_1, \dots, X_n]$ is not principal. For example, in the ring $K[X, Y]$, the ideal (X, Y) generated by X and Y is not principal. First, since (X, Y) is the set of all polynomials of the form $Xq_1 + Yq_2$, where $q_1, q_2 \in K[X, Y]$, except when $Xq_1 + Yq_2 = 0$, we have $\deg(Xq_1 + Yq_2) \geq 1$. Thus, $1 \notin (X, Y)$. Now if there was some $p \in K[X, Y]$ such that $(X, Y) = (p)$, since $1 \notin (X, Y)$, we must have $\deg(p) \geq 1$. But we would also have $X = pq_1$ and $Y = pq_2$, for some $q_1, q_2 \in K[X, Y]$. Since $\deg(X) = \deg(Y) = 1$, this is impossible.

Even though $K[X, Y]$ is not a principal ring, a suitable version of unique factorization in terms of irreducible factors holds. The ring $K[X, Y]$ (and more generally $K[X_1, \dots, X_n]$) is what is called a *unique factorization domain*, for short, UFD, or a *factorial ring*.

From this point until Definition 20.10, we consider polynomials in one variable over a field K .

Remark: Although we already proved part (1) of Proposition 20.9 in a more general situation above, we reprove it in the special case of polynomials. This may offend the purists, but most readers will probably not mind.

Proposition 20.9. *Let K be a field. The following properties hold:*

- (1) *For any two nonzero polynomials $f, g \in K[X]$, $(f) = (g)$ iff there is some $\lambda \neq 0$ in K such that $g = \lambda f$.*
- (2) *For every nonnull ideal \mathfrak{J} in $K[X]$, there is a unique monic polynomial $f \in K[X]$ such that $\mathfrak{J} = (f)$.*

Proof. (1) If $(f) = (g)$, there are some nonzero polynomials $q_1, q_2 \in K[X]$ such that $g = fq_1$ and $f = gq_2$. Thus, we have $f = fq_1q_2$, which implies $f(1 - q_1q_2) = 0$. Since K is a field, by Proposition 20.1, $K[X]$ has no zero divisor, and since we assumed $f \neq 0$, we must have $q_1q_2 = 1$. However, if either q_1 or q_2 is not a constant, by Proposition 20.1 again, $\deg(q_1q_2) = \deg(q_1) + \deg(q_2) \geq 1$, contradicting $q_1q_2 = 1$, since $\deg(1) = 0$. Thus, both $q_1, q_2 \in K - \{0\}$, and (1) holds with $\lambda = q_1$. In the other direction, it is obvious that $g = \lambda f$ implies that $(f) = (g)$.

(2) Since we are assuming that \mathfrak{J} is not the null ideal, there is some polynomial of smallest degree in \mathfrak{J} , and since K is a field, by suitable multiplication by a scalar, we can make sure that this polynomial is monic. Thus, let f be a monic polynomial of smallest degree in \mathfrak{J} . By (ID2), it is clear that $(f) \subseteq \mathfrak{J}$. Now, let $g \in \mathfrak{J}$. Using the Euclidean algorithm, there exist unique $q, r \in K[X]$ such that

$$g = qf + r \quad \text{and} \quad \deg(r) < \deg(f).$$

If $r \neq 0$, there is some $\lambda \neq 0$ in K such that λr is a monic polynomial, and since $\lambda r = \lambda g - \lambda q f$, with $f, g \in \mathfrak{J}$, by (ID1) and (ID2), we have $\lambda r \in \mathfrak{J}$, where $\deg(\lambda r) < \deg(f)$ and λr is a monic polynomial, contradicting the minimality of the degree of f . Thus, $r = 0$, and $g \in (f)$. The uniqueness of the monic polynomial f follows from (1). \square

Proposition 20.9 shows that $K[X]$ is a principal ring when K is a field.

We now investigate the existence of a greatest common divisor (gcd) for two nonzero polynomials. Given any two nonzero polynomials $f, g \in K[X]$, recall that f divides g if $g = fq$ for some $q \in K[X]$.

Definition 20.6. Given any two nonzero polynomials $f, g \in K[X]$, a polynomial $d \in K[X]$ is a *greatest common divisor of f and g* (for short, a *gcd of f and g*) if d divides f and g and whenever $h \in K[X]$ divides f and g , then h divides d . We say that f and g are *relatively prime* if 1 is a gcd of f and g .

Note that f and g are relatively prime iff all of their gcd's are constants (scalars in K), or equivalently, if f, g have no divisor q of degree $\deg(q) \geq 1$.



In particular, note that f and g are relatively prime when f is a nonzero constant polynomial (a scalar $\lambda \neq 0$ in K) and g is any nonzero polynomial.

We can characterize gcd's of polynomials as follows.

Proposition 20.10. Let K be a field and let $f, g \in K[X]$ be any two nonzero polynomials. For every polynomial $d \in K[X]$, the following properties are equivalent:

- (1) The polynomial d is a gcd of f and g .
- (2) The polynomial d divides f and g and there exist $u, v \in K[X]$ such that

$$d = uf + vg.$$

- (3) The ideals (f) , (g) , and (d) satisfy the equation

$$(d) = (f) + (g).$$

In addition, $d \neq 0$, and d is unique up to multiplication by a nonzero scalar in K .

Proof. Given any two nonzero polynomials $u, v \in K[X]$, observe that u divides v iff $(v) \subseteq (u)$. Now, (2) can be restated as $(f) \subseteq (d)$, $(g) \subseteq (d)$, and $d \in (f) + (g)$, which is equivalent to $(d) = (f) + (g)$, namely (3).

If (2) holds, since $d = uf + vg$, whenever $h \in K[X]$ divides f and g , then h divides d , and d is a gcd of f and g .

Assume that d is a gcd of f and g . Then, since d divides f and d divides g , we have $(f) \subseteq (d)$ and $(g) \subseteq (d)$, and thus $(f) + (g) \subseteq (d)$, and $(f) + (g)$ is nonempty since f and g are nonzero. By Proposition 20.9, there exists a monic polynomial $d_1 \in K[X]$ such that $(d_1) = (f) + (g)$. Then, d_1 divides both f and g , and since d is a gcd of f and g , then d_1 divides d , which shows that $(d) \subseteq (d_1) = (f) + (g)$. Consequently, $(f) + (g) = (d)$, and (3) holds.

Since $(d) = (f) + (g)$ and f and g are nonzero, the last part of the proposition is obvious. \square

As a consequence of Proposition 20.10, two nonzero polynomials $f, g \in K[X]$ are relatively prime iff there exist $u, v \in K[X]$ such that

$$uf + vg = 1.$$

The identity

$$d = uf + vg$$

of part (2) of Proposition 20.10 is often called the *Bezout identity*.

We derive more useful consequences of Proposition 20.10.

Proposition 20.11. *Let K be a field and let $f, g \in K[X]$ be any two nonzero polynomials. For every gcd $d \in K[X]$ of f and g , the following properties hold:*

- (1) *For every nonzero polynomial $q \in K[X]$, the polynomial dq is a gcd of fq and gq .*
- (2) *For every nonzero polynomial $q \in K[X]$, if q divides f and g , then d/q is a gcd of f/q and g/q .*

Proof. (1) By Proposition 20.10 (2), d divides f and g , and there exist $u, v \in K[X]$, such that

$$d = uf + vg.$$

Then, dq divides fq and gq , and

$$dq = ufq + vgg.$$

By Proposition 20.10 (2), dq is a gcd of fq and gq . The proof of (2) is similar. \square

The following proposition is used often.

Proposition 20.12. (*Euclid's proposition*) *Let K be a field and let $f, g, h \in K[X]$ be any nonzero polynomials. If f divides gh and f is relatively prime to g , then f divides h .*

Proof. From Proposition 20.10, f and g are relatively prime iff there exist some polynomials $u, v \in K[X]$ such that

$$uf + vg = 1.$$

Then, we have

$$ufh + vgh = h,$$

and since f divides gh , it divides both ufh and vgh , and so, f divides h . \square

Proposition 20.13. *Let K be a field and let $f, g_1, \dots, g_m \in K[X]$ be some nonzero polynomials. If f and g_i are relatively prime for all i , $1 \leq i \leq m$, then f and $g_1 \cdots g_m$ are relatively prime.*

Proof. We proceed by induction on m . The case $m = 1$ is trivial. Let $h = g_2 \cdots g_m$. By the induction hypothesis, f and h are relatively prime. Let d be a gcd of f and $g_1 h$. We claim that d is relatively prime to g_1 . Otherwise, d and g_1 would have some nonconstant gcd d_1 which would divide both f and g_1 , contradicting the fact that f and g_1 are relatively prime. Now, by Proposition 20.12, since d divides $g_1 h$ and d and g_1 are relatively prime, d divides $h = g_2 \cdots g_m$. But then, d is a divisor of f and h , and since f and h are relatively prime, d must be a constant, and f and $g_1 \cdots g_m$ are relatively prime. \square

Definition 20.6 is generalized to any finite number of polynomials as follows.

Definition 20.7. Given any nonzero polynomials $f_1, \dots, f_n \in K[X]$, where $n \geq 2$, a polynomial $d \in K[X]$ is a *greatest common divisor* of f_1, \dots, f_n (for short, a *gcd* of f_1, \dots, f_n) if d divides each f_i and whenever $h \in K[X]$ divides each f_i , then h divides d . We say that f_1, \dots, f_n are *relatively prime* if 1 is a gcd of f_1, \dots, f_n .

It is easily shown that Proposition 20.10 can be generalized to any finite number of polynomials, and similarly for its relevant corollaries. The details are left as an exercise.

Proposition 20.14. *Let K be a field and let $f_1, \dots, f_n \in K[X]$ be any $n \geq 2$ nonzero polynomials. For every polynomial $d \in K[X]$, the following properties are equivalent:*

- (1) *The polynomial d is a gcd of f_1, \dots, f_n .*
- (2) *The polynomial d divides each f_i and there exist $u_1, \dots, u_n \in K[X]$ such that*

$$d = u_1 f_1 + \cdots + u_n f_n.$$

- (3) *The ideals (f_i) , and (d) satisfy the equation*

$$(d) = (f_1) + \cdots + (f_n).$$

In addition, $d \neq 0$, and d is unique up to multiplication by a nonzero scalar in K .

As a consequence of Proposition 20.14, some polynomials $f_1, \dots, f_n \in K[X]$ are relatively prime iff there exist $u_1, \dots, u_n \in K[X]$ such that

$$u_1 f_1 + \cdots + u_n f_n = 1.$$

The identity

$$u_1 f_1 + \cdots + u_n f_n = 1$$

of part (2) of Proposition 20.14 is also called the *Bezout identity*.

We now consider the factorization of polynomials of a single variable into irreducible factors.

20.5 Factorization and Irreducible Factors in $K[X]$

Definition 20.8. Given a field K , a polynomial $p \in K[X]$ is *irreducible* or *indecomposable* or *prime* if $\deg(p) \geq 1$ and if p is not divisible by any polynomial $q \in K[X]$ such that $1 \leq \deg(q) < \deg(p)$. Equivalently, p is irreducible if $\deg(p) \geq 1$ and if $p = q_1 q_2$, then either $q_1 \in K$ or $q_2 \in K$ (and of course, $q_1 \neq 0$, $q_2 \neq 0$).

Example 20.2. Every polynomial $aX + b$ of degree 1 is irreducible. Over the field \mathbb{R} , the polynomial $X^2 + 1$ is irreducible (why?), but $X^3 + 1$ is not irreducible, since

$$X^3 + 1 = (X + 1)(X^2 - X + 1).$$

The polynomial $X^2 - X + 1$ is irreducible over \mathbb{R} (why?). It would seem that $X^4 + 1$ is irreducible over \mathbb{R} , but in fact,

$$X^4 + 1 = (X^2 - \sqrt{2}X + 1)(X^2 + \sqrt{2}X + 1).$$

However, in view of the above factorization, $X^4 + 1$ is irreducible over \mathbb{Q} .

It can be shown that the irreducible polynomials over \mathbb{R} are the polynomials of degree 1, or the polynomials of degree 2 of the form $aX^2 + bX + c$, for which $b^2 - 4ac < 0$ (i.e., those having no real roots). This is not easy to prove! Over the complex numbers \mathbb{C} , the only irreducible polynomials are those of degree 1. This is a version of a fact often referred to as the “Fundamental theorem of Algebra”, or, as the French sometimes say, as “d’Alembert’s theorem”!

We already observed that for any two nonzero polynomials $f, g \in K[X]$, f divides g iff $(g) \subseteq (f)$. In view of the definition of a maximal ideal given in Definition 20.4, we now prove that a polynomial $p \in K[X]$ is irreducible iff (p) is a maximal ideal in $K[X]$.

Proposition 20.15. *A polynomial $p \in K[X]$ is irreducible iff (p) is a maximal ideal in $K[X]$.*

Proof. Since $K[X]$ is an integral domain, for all nonzero polynomials $p, q \in K[X]$, $\deg(pq) = \deg(p) + \deg(q)$, and thus, $(p) \neq K[X]$ iff $\deg(p) \geq 1$. Assume that $p \in K[X]$ is irreducible. Since every ideal in $K[X]$ is a principal ideal, every ideal in $K[X]$ is of the form (q) , for some $q \in K[X]$. If $(p) \subseteq (q)$, with $\deg(q) \geq 1$, then q divides p , and since $p \in K[X]$ is irreducible, this implies that $p = \lambda q$ for some $\lambda \neq 0$ in K , and so, $(p) = (q)$. Thus, (p) is a maximal ideal. Conversely, assume that (p) is a maximal ideal. Then, as we showed above, $\deg(p) \geq 1$, and if q divides p , with $\deg(q) \geq 1$, then $(p) \subseteq (q)$, and since (p) is a maximal ideal, this implies that $(p) = (q)$, which means that $p = \lambda q$ for some $\lambda \neq 0$ in K , and so, p is irreducible. \square

Let $p \in K[X]$ be irreducible. Then, for every nonzero polynomial $g \in K[X]$, either p and g are relatively prime, or p divides g . Indeed, if d is any gcd of p and g , if d is a constant, then

p and g are relatively prime, and if not, because p is irreducible, we have $d = \lambda p$ for some $\lambda \neq 0$ in K , and thus, p divides g . As a consequence, if $p, q \in K[X]$ are both irreducible, then either p and q are relatively prime, or $p = \lambda q$ for some $\lambda \neq 0$ in K . In particular, if $p, q \in K[X]$ are both irreducible monic polynomials and $p \neq q$, then p and q are relatively prime.

We now prove the (unique) factorization of polynomials into irreducible factors.

Theorem 20.16. *Given any field K , for every nonzero polynomial*

$$f = a_d X^d + a_{d-1} X^{d-1} + \cdots + a_0$$

of degree $d = \deg(f) \geq 1$ in $K[X]$, there exists a unique set $\{\langle p_1, k_1 \rangle, \dots, \langle p_m, k_m \rangle\}$ such that

$$f = a_d p_1^{k_1} \cdots p_m^{k_m},$$

where the $p_i \in K[X]$ are distinct irreducible monic polynomials, the k_i are (not necessarily distinct) integers, and $m \geq 1, k_i \geq 1$.

Proof. First, we prove the existence of such a factorization by induction on $d = \deg(f)$. Clearly, it is enough to prove the result for monic polynomials f of degree $d = \deg(f) \geq 1$. If $d = 1$, then $f = X + a_0$, which is an irreducible monic polynomial.

Assume $d \geq 2$, and assume the induction hypothesis for all monic polynomials of degree $< d$. Consider the set S of all monic polynomials g such that $\deg(g) \geq 1$ and g divides f . Since $f \in S$, the set S is nonempty, and thus, S contains some monic polynomial p_1 of minimal degree. Since $\deg(p_1) \geq 1$, the monic polynomial p_1 must be irreducible. Otherwise we would have $p_1 = g_1 g_2$, for some monic polynomials g_1, g_2 such that $\deg(p_1) > \deg(g_1) \geq 1$ and $\deg(p_1) > \deg(g_2) \geq 1$, and since p_1 divide f , then g_1 would divide f , contradicting the minimality of the degree of p_1 . Thus, we have $f = p_1 q$, for some irreducible monic polynomial p_1 , with q also monic. Since $\deg(p_1) \geq 1$, we have $\deg(q) < \deg(f)$, and we can apply the induction hypothesis to q . Thus, we obtain a factorization of the desired form.

We now prove uniqueness. Assume that

$$f = a_d p_1^{k_1} \cdots p_m^{k_m},$$

and

$$f = a_d q_1^{h_1} \cdots q_n^{h_n}.$$

Thus, we have

$$a_d p_1^{k_1} \cdots p_m^{k_m} = a_d q_1^{h_1} \cdots q_n^{h_n}.$$

We prove that $m = n$, $p_i = q_i$ and $h_i = k_i$, for all i , with $1 \leq i \leq n$.

The proof proceeds by induction on $h_1 + \cdots + h_n$.

If $h_1 + \cdots + h_n = 1$, then $n = 1$ and $h_1 = 1$. Then, since $K[X]$ is an integral domain, we have

$$p_1^{k_1} \cdots p_m^{k_m} = q_1,$$

and since q_1 and the p_i are irreducible monic, we must have $m = 1$ and $p_1 = q_1$.

If $h_1 + \cdots + h_n \geq 2$, since $K[X]$ is an integral domain and since $h_1 \geq 1$, we have

$$p_1^{k_1} \cdots p_m^{k_m} = q_1 q,$$

with

$$q = q_1^{h_1-1} \cdots q_n^{h_n},$$

where $(h_1 - 1) + \cdots + h_n \geq 1$ (and $q_1^{h_1-1} = 1$ if $h_1 = 1$). Now, if q_1 is not equal to any of the p_i , by a previous remark, q_1 and p_i are relatively prime, and by Proposition 20.13, q_1 and $p_1^{k_1} \cdots p_m^{k_m}$ are relatively prime. But this contradicts the fact that q_1 divides $p_1^{k_1} \cdots p_m^{k_m}$. Thus, q_1 is equal to one of the p_i . Without loss of generality, we can assume that $q_1 = p_1$. Then, since $K[X]$ is an integral domain, we have

$$p_1^{k_1-1} \cdots p_m^{k_m} = q_1^{h_1-1} \cdots q_n^{h_n},$$

where $p_1^{k_1-1} = 1$ if $k_1 = 1$, and $q_1^{h_1-1} = 1$ if $h_1 = 1$. Now, $(h_1 - 1) + \cdots + h_n < h_1 + \cdots + h_n$, and we can apply the induction hypothesis to conclude that $m = n$, $p_i = q_i$ and $h_i = k_i$, with $1 \leq i \leq n$. \square

The above considerations about unique factorization into irreducible factors can be extended almost without changes to more general rings known as *Euclidean domains*. In such rings, some abstract version of the division theorem is assumed to hold.

Definition 20.9. A *Euclidean domain* (or *Euclidean ring*) is an integral domain A such that there exists a function $\varphi: A \rightarrow \mathbb{N}$ with the following property: For all $a, b \in A$ with $b \neq 0$, there are some $q, r \in A$ such that

$$a = bq + r \quad \text{and} \quad \varphi(r) < \varphi(b).$$

Note that the pair (q, r) is not necessarily unique.

Actually, unique factorization holds in principal ideal domains (PID's), see Theorem 21.12. As shown below, every Euclidean domain is a PID, and thus, unique factorization holds for Euclidean domains.

Proposition 20.17. *Every Euclidean domain A is a PID.*

Proof. Let \mathfrak{I} be a nonnull ideal in A . Then, the set

$$\{\varphi(a) \mid a \in \mathfrak{I}\}$$

is nonempty, and thus, has a smallest element m . Let b be any (nonnull) element of \mathfrak{I} such that $m = \varphi(b)$. We claim that $\mathfrak{I} = (b)$. Given any $a \in \mathfrak{I}$, we can write

$$a = bq + r$$

for some $q, r \in A$, with $\varphi(r) < \varphi(b)$. Since $b \in \mathfrak{I}$ and \mathfrak{I} is an ideal, we also have $bq \in \mathfrak{I}$, and since $a, bq \in \mathfrak{I}$ and \mathfrak{I} is an ideal, then $r \in \mathfrak{I}$ with $\varphi(r) < \varphi(b) = m$, contradicting the minimality of m . Thus, $r = 0$ and $a \in (b)$. But then,

$$\mathfrak{I} \subseteq (b),$$

and since $b \in \mathfrak{I}$, we get

$$\mathfrak{I} = (b),$$

and A is a PID. □

As a corollary of Proposition 20.17, the ring \mathbb{Z} is a Euclidean domain (using the function $\varphi(a) = |a|$) and thus, a PID. If K is a field, the function φ on $K[X]$ defined such that

$$\varphi(f) = \begin{cases} 0 & \text{if } f = 0, \\ \deg(f) + 1 & \text{if } f \neq 0, \end{cases}$$

shows that $K[X]$ is a Euclidean domain.

Example 20.3. A more interesting example of a Euclidean domain is the ring $\mathbb{Z}[i]$ of *Gaussian integers*, i.e., the subring of \mathbb{C} consisting of all complex numbers of the form $a + ib$, where $a, b \in \mathbb{Z}$. Using the function φ defined such that

$$\varphi(a + ib) = a^2 + b^2,$$

we leave it as an interesting exercise to prove that $\mathbb{Z}[i]$ is a Euclidean domain.



Not every PID is a Euclidean ring.

Remark: Given any integer, $d \in \mathbb{Z}$, such that $d \neq 0, 1$ and d does not have any square factor greater than one, the *quadratic field*, $\mathbb{Q}(\sqrt{d})$, is the field consisting of all complex numbers of the form $a + ib\sqrt{-d}$ if $d < 0$, and of all the real numbers of the form $a + b\sqrt{d}$ if $d > 0$, with $a, b \in \mathbb{Q}$. The subring of $\mathbb{Q}(\sqrt{d})$ consisting of all elements as above for which $a, b \in \mathbb{Z}$ is denoted by $\mathbb{Z}[\sqrt{d}]$. We define the *ring of integers* of the field $\mathbb{Q}(\sqrt{d})$ as the subring of $\mathbb{Q}(\sqrt{d})$ consisting of the following elements:

- (1) If $d \equiv 2 \pmod{4}$ or $d \equiv 3 \pmod{4}$, then all elements of the form $a + ib\sqrt{-d}$ (if $d < 0$) or all elements of the form $a + b\sqrt{d}$ (if $d > 0$), with $a, b \in \mathbb{Z}$;
- (2) If $d \equiv 1 \pmod{4}$, then all elements of the form $(a + ib\sqrt{-d})/2$ (if $d < 0$) or all elements of the form $(a + b\sqrt{d})/2$ (if $d > 0$), with $a, b \in \mathbb{Z}$ and with a, b either both even or both odd.

Observe that when $d \equiv 2 \pmod{4}$ or $d \equiv 3 \pmod{4}$, the ring of integers of $\mathbb{Q}(\sqrt{d})$ is equal to $\mathbb{Z}[\sqrt{d}]$. For more on quadratic fields and their rings of integers, see Stark [100] (Chapter 8) or Niven, Zuckerman and Montgomery [86] (Chapter 9). It can be shown that the rings of integers, $\mathbb{Z}[\sqrt{-d}]$, where $d = 19, 43, 67, 163$, are PID's, but not Euclidean rings.

Actually the rings of integers of $\mathbb{Q}(\sqrt{d})$ that are Euclidean domains are completely determined but the proof is quite difficult. It turns out that there are twenty one such rings corresponding to the integers: $-11, -7, -3, -2, -1, 2, 3, 5, 6, 7, 11, 13, 17, 19, 21, 29, 33, 37, 41, 57$ and 73 , see Stark [100] (Chapter 8).

It is possible to characterize a larger class of rings (in terms of ideals), *factorial rings (or unique factorization domains)*, for which unique factorization holds (see Section 21.1). We now consider zeros (or roots) of polynomials.

20.6 Roots of Polynomials

We go back to the general case of an arbitrary ring for a little while.

Definition 20.10. Given a ring A and any polynomial $f \in A[X]$, we say that some $\alpha \in A$ is a *zero of f , or a root of f* , if $f(\alpha) = 0$. Similarly, given a polynomial $f \in A[X_1, \dots, X_n]$, we say that $(\alpha_1, \dots, \alpha_n) \in A^n$ is a *zero of f , or a root of f* , if $f(\alpha_1, \dots, \alpha_n) = 0$.

When $f \in A[X]$ is the null polynomial, every $\alpha \in A$ is trivially a zero of f . This case being trivial, we usually assume that we are considering zeros of nonnull polynomials.

Example 20.4. Considering the polynomial $f(X) = X^2 - 1$, both $+1$ and -1 are zeros of $f(X)$. Over the field of reals, the polynomial $g(X) = X^2 + 1$ has no zeros. Over the field \mathbb{C} of complex numbers, $g(X) = X^2 + 1$ has two roots i and $-i$, the square roots of -1 , which are “imaginary numbers.”

We have the following basic proposition showing the relationship between polynomial division and roots.

Proposition 20.18. *Let $f \in A[X]$ be any polynomial and $\alpha \in A$ any element of A . If the result of dividing f by $X - \alpha$ is $f = (X - \alpha)q + r$, then $r = 0$ iff $f(\alpha) = 0$, i.e., α is a root of f iff $r = 0$.*

Proof. We have $f = (X - \alpha)q + r$, with $\deg(r) < 1 = \deg(X - \alpha)$. Thus, r is a constant in K , and since $f(\alpha) = (\alpha - \alpha)q(\alpha) + r$, we get $f(\alpha) = r$, and the proposition is trivial. \square

We now consider the issue of multiplicity of a root.

Proposition 20.19. *Let $f \in A[X]$ be any nonnull polynomial and $h \geq 0$ any integer. The following conditions are equivalent.*

- (1) f is divisible by $(X - \alpha)^h$ but not by $(X - \alpha)^{h+1}$.

(2) There is some $g \in A[X]$ such that $f = (X - \alpha)^h g$ and $g(\alpha) \neq 0$.

Proof. Assume (1). Then, we have $f = (X - \alpha)^h g$ for some $g \in A[X]$. If we had $g(\alpha) = 0$, by Proposition 20.18, g would be divisible by $(X - \alpha)$, and then f would be divisible by $(X - \alpha)^{h+1}$, contradicting (1).

Assume (2), that is, $f = (X - \alpha)^h g$ and $g(\alpha) \neq 0$. If f is divisible by $(X - \alpha)^{h+1}$, then we have $f = (X - \alpha)^{h+1} g_1$, for some $g_1 \in A[X]$. Then, we have

$$(X - \alpha)^h g = (X - \alpha)^{h+1} g_1,$$

and thus

$$(X - \alpha)^h (g - (X - \alpha)g_1) = 0,$$

and since the leading coefficient of $(X - \alpha)^h$ is 1 (show this by induction), by Proposition 20.1, $(X - \alpha)^h$ is not a zero divisor, and we get $g - (X - \alpha)g_1 = 0$, i.e., $g = (X - \alpha)g_1$, and so $g(\alpha) = 0$, contrary to the hypothesis. \square

As a consequence of Proposition 20.19, for every nonnull polynomial $f \in A[X]$ and every $\alpha \in A$, there is a unique integer $h \geq 0$ such that f is divisible by $(X - \alpha)^h$ but not by $(X - \alpha)^{h+1}$. Indeed, since f is divisible by $(X - \alpha)^h$, we have $h \leq \deg(f)$. When $h = 0$, α is not a root of f , i.e., $f(\alpha) \neq 0$. The interesting case is when α is a root of f .

Definition 20.11. Given a ring A and any nonnull polynomial $f \in A[X]$, given any $\alpha \in A$, the unique $h \geq 0$ such that f is divisible by $(X - \alpha)^h$ but not by $(X - \alpha)^{h+1}$ is called the *order, or multiplicity, of α* . We have $h = 0$ iff α is not a root of f , and when α is a root of f , if $h = 1$, we call α a *simple root*, if $h = 2$, a *double root*, and generally, a root of multiplicity $h \geq 2$ is called a *multiple root*.

Observe that Proposition 20.19 (2) implies that if $A \subseteq B$, where A and B are rings, for every nonnull polynomial $f \in A[X]$, if $\alpha \in A$ is a root of f , then the multiplicity of α with respect to $f \in A[X]$ and the multiplicity of α with respect to f considered as a polynomial in $B[X]$, is the same.

We now show that if the ring A is an integral domain, the number of roots of a nonzero polynomial is at most its degree.

Proposition 20.20. Let $f, g \in A[X]$ be nonnull polynomials, let $\alpha \in A$, and let $h \geq 0$ and $k \geq 0$ be the multiplicities of α with respect to f and g . The following properties hold.

- (1) If l is the multiplicity of α with respect to $(f + g)$, then $l \geq \min(h, k)$. If $h \neq k$, then $l = \min(h, k)$.
- (2) If m is the multiplicity of α with respect to fg , then $m \geq h + k$. If A is an integral domain, then $m = h + k$.

Proof. (1) We have $f(X) = (X - \alpha)^h f_1(X)$, $g(X) = (X - \alpha)^k g_1(X)$, with $f_1(\alpha) \neq 0$ and $g_1(\alpha) \neq 0$. Clearly, $l \geq \min(h, k)$. If $h \neq k$, assume $h < k$. Then, we have

$$f(X) + g(X) = (X - \alpha)^h f_1(X) + (X - \alpha)^k g_1(X) = (X - \alpha)^h (f_1(X) + (X - \alpha)^{k-h} g_1(X)),$$

and since $(f_1(X) + (X - \alpha)^{k-h} g_1(X))(\alpha) = f_1(\alpha) \neq 0$, we have $l = h = \min(h, k)$.

(2) We have

$$f(X)g(X) = (X - \alpha)^{h+k} f_1(X)g_1(X),$$

with $f_1(\alpha) \neq 0$ and $g_1(\alpha) \neq 0$. Clearly, $m \geq h + k$. If A is an integral domain, then $f_1(\alpha)g_1(\alpha) \neq 0$, and so $m = h + k$. \square

Proposition 20.21. *Let A be an integral domain. Let f be any nonnull polynomial $f \in A[X]$ and let $\alpha_1, \dots, \alpha_m \in A$ be $m \geq 1$ distinct roots of f of respective multiplicities k_1, \dots, k_m . Then, we have*

$$f(X) = (X - \alpha_1)^{k_1} \cdots (X - \alpha_m)^{k_m} g(X),$$

where $g \in A[X]$ and $g(\alpha_i) \neq 0$ for all i , $1 \leq i \leq m$.

Proof. We proceed by induction on m . The case $m = 1$ is obvious in view of Definition 20.11 (which itself, is justified by Proposition 20.19). If $m \geq 2$, by the induction hypothesis, we have

$$f(X) = (X - \alpha_1)^{k_1} \cdots (X - \alpha_{m-1})^{k_{m-1}} g_1(X),$$

where $g_1 \in A[X]$ and $g_1(\alpha_i) \neq 0$, for $1 \leq i \leq m - 1$. Since A is an integral domain and $\alpha_i \neq \alpha_j$ for $i \neq j$, since α_m is a root of f , we have

$$0 = (\alpha_m - \alpha_1)^{k_1} \cdots (\alpha_m - \alpha_{m-1})^{k_{m-1}} g_1(\alpha_m),$$

which implies that $g_1(\alpha_m) = 0$. Now, by Proposition 20.20 (2), since α_m is not a root of the polynomial $(X - \alpha_1)^{k_1} \cdots (X - \alpha_{m-1})^{k_{m-1}}$ and since A is an integral domain, α_m must be a root of multiplicity k_m of g_1 , which means that

$$g_1(X) = (X - \alpha_m)^{k_m} g(X),$$

with $g(\alpha_m) \neq 0$. Since $g_1(\alpha_i) \neq 0$ for $1 \leq i \leq m - 1$ and A is an integral domain, we must also have $g(\alpha_i) \neq 0$, for $1 \leq i \leq m - 1$. Thus, we have

$$f(X) = (X - \alpha_1)^{k_1} \cdots (X - \alpha_m)^{k_m} g(X),$$

where $g \in A[X]$, and $g(\alpha_i) \neq 0$ for $1 \leq i \leq m$. \square

As a consequence of Proposition 20.21, we get the following important result.

Theorem 20.22. *Let A be an integral domain. For every nonnull polynomial $f \in A[X]$, if the degree of f is $n = \deg(f)$ and k_1, \dots, k_m are the multiplicities of all the distinct roots of f (where $m \geq 0$), then $k_1 + \cdots + k_m \leq n$.*

Proof. Immediate from Proposition 20.21. \square

Since fields are integral domains, Theorem 20.22 holds for nonnull polynomials over fields and in particular, for \mathbb{R} and \mathbb{C} . An important consequence of Theorem 20.22 is the following.

Proposition 20.23. *Let A be an integral domain. For any two polynomials $f, g \in A[X]$, if $\deg(f) \leq n$, $\deg(g) \leq n$, and if there are $n + 1$ distinct elements $\alpha_1, \alpha_2, \dots, \alpha_{n+1} \in A$ (with $\alpha_i \neq \alpha_j$ for $i \neq j$) such that $f(\alpha_i) = g(\alpha_i)$ for all i , $1 \leq i \leq n + 1$, then $f = g$.*

Proof. Assume $f \neq g$, then, $(f - g)$ is nonnull, and since $f(\alpha_i) = g(\alpha_i)$ for all i , $1 \leq i \leq n + 1$, the polynomial $(f - g)$ has $n + 1$ distinct roots. Thus, $(f - g)$ has $n + 1$ distinct roots and is of degree at most n , which contradicts Theorem 20.22. \square

Proposition 20.23 is often used to show that polynomials coincide. We will use it to show some interpolation formulae due to Lagrange and Hermite. But first, we characterize the multiplicity of a root of a polynomial. For this, we need the notion of derivative familiar in analysis. Actually, we can simply define this notion algebraically.

First, we need to rule out some undesirable behaviors. Given a field K , as we saw in Example 2.4, we can define a homomorphism $\chi: \mathbb{Z} \rightarrow K$ given by

$$\chi(n) = n \cdot 1,$$

where 1 is the multiplicative identity of K . Recall that we define $n \cdot a$ by

$$n \cdot a = \underbrace{a + \cdots + a}_n$$

if $n \geq 0$ (with $0 \cdot a = 0$) and

$$n \cdot a = -(-n) \cdot a$$

if $n < 0$. We say that the field K is of *characteristic zero* if the homomorphism χ is injective. Then, for any $a \in K$ with $a \neq 0$, we have $n \cdot a \neq 0$ for all $n \neq 0$.

The fields \mathbb{Q} , \mathbb{R} , and \mathbb{C} are of characteristic zero. In fact, it is easy to see that every field of characteristic zero contains a subfield isomorphic to \mathbb{Q} . Thus, finite fields can't be of characteristic zero.

Remark: If a field is not of characteristic zero, it is not hard to show that its characteristic, that is, the smallest $n \geq 2$ such that $n \cdot 1 = 0$, is a prime number p . The characteristic p of K is the generator of the principal ideal $p\mathbb{Z}$, the kernel of the homomorphism $\chi: \mathbb{Z} \rightarrow K$. Thus, every finite field is of characteristic some prime p . Infinite fields of nonzero characteristic also exist.

Definition 20.12. Let A be a ring. The *derivative* f' , or Df , or D^1f , of a polynomial $f \in A[X]$ is defined inductively as follows:

$$\begin{aligned} f' &= 0, & \text{if } f = 0, \text{ the null polynomial,} \\ f' &= 0, & \text{if } f = a, a \neq 0, a \in A, \\ f' &= na_nX^{n-1} + (n-1)a_{n-1}X^{n-2} + \cdots + 2a_2X + a_1, \\ & & \text{if } f = a_nX^n + a_{n-1}X^{n-1} + \cdots + a_0, \text{ with } n = \deg(f) \geq 1. \end{aligned}$$

If $A = K$ is a field of characteristic zero, if $\deg(f) \geq 1$, the leading coefficient na_n of f' is nonzero, and thus, f' is not the null polynomial. Thus, if $A = K$ is a field of characteristic zero, when $n = \deg(f) \geq 1$, we have $\deg(f') = n - 1$.



For rings or for fields of characteristic $p \geq 2$, we could have $f' = 0$, for a polynomial f of degree ≥ 1 .

The following standard properties of derivatives are recalled without proof (prove them as an exercise).

Given any two polynomials, $f, g \in A[X]$, we have

$$\begin{aligned} (f + g)' &= f' + g', \\ (fg)' &= f'g + fg'. \end{aligned}$$

For example, if $f = (X - \alpha)^k g$ and $k \geq 1$, we have

$$f' = k(X - \alpha)^{k-1}g + (X - \alpha)^k g'.$$

We can now give a criterion for the existence of simple roots. The first proposition holds for any ring.

Proposition 20.24. *Let A be any ring. For every nonnull polynomial $f \in A[X]$, $\alpha \in A$ is a simple root of f iff α is a root of f and α is not a root of f' .*

Proof. Since $\alpha \in A$ is a root of f , we have $f = (X - \alpha)g$ for some $g \in A[X]$. Now, α is a simple root of f iff $g(\alpha) \neq 0$. However, we have $f' = g + (X - \alpha)g'$, and so $f'(\alpha) = g(\alpha)$. Thus, α is a simple root of f iff $f'(\alpha) \neq 0$. \square

We can improve the previous proposition as follows.

Proposition 20.25. *Let A be any ring. For every nonnull polynomial $f \in A[X]$, let $\alpha \in A$ be a root of multiplicity $k \geq 1$ of f . Then, α is a root of multiplicity at least $k - 1$ of f' . If A is a field of characteristic zero, then α is a root of multiplicity $k - 1$ of f' .*

Proof. Since $\alpha \in A$ is a root of multiplicity k of f , we have $f = (X - \alpha)^k g$ for some $g \in A[X]$ and $g(\alpha) \neq 0$. Since

$$f' = k(X - \alpha)^{k-1}g + (X - \alpha)^k g' = (X - \alpha)^{k-1}(kg + (X - \alpha)g'),$$

it is clear that the multiplicity of α w.r.t. f' is at least $k-1$. Now, $(kg + (X - \alpha)g')(\alpha) = kg(\alpha)$, and if A is of characteristic zero, since $g(\alpha) \neq 0$, then $kg(\alpha) \neq 0$. Thus, α is a root of multiplicity $k-1$ of f' . \square

As a consequence, we obtain the following test for the existence of a root of multiplicity k for a polynomial f :

Given a field K of characteristic zero, for any nonnull polynomial $f \in K[X]$, any $\alpha \in K$ is a root of multiplicity $k \geq 1$ of f iff α is a root of $f, D^1 f, D^2 f, \dots, D^{k-1} f$, but not a root of $D^k f$.

We can now return to polynomial functions and tie up some loose ends. Given a ring A , recall that every polynomial $f \in A[X_1, \dots, X_n]$ induces a function $f_A: A^n \rightarrow A$ defined such that $f_A(\alpha_1, \dots, \alpha_n) = f(\alpha_1, \dots, \alpha_n)$, for every $(\alpha_1, \dots, \alpha_n) \in A^n$. We now give a sufficient condition for the mapping $f \mapsto f_A$ to be injective.

Proposition 20.26. *Let A be an integral domain. For every polynomial $f \in A[X_1, \dots, X_n]$, if A_1, \dots, A_n are n infinite subsets of A such that $f(\alpha_1, \dots, \alpha_n) = 0$ for all $(\alpha_1, \dots, \alpha_n) \in A_1 \times \dots \times A_n$, then $f = 0$, i.e., f is the null polynomial. As a consequence, if A is an infinite integral domain, then the map $f \mapsto f_A$ is injective.*

Proof. We proceed by induction on n . Assume $n = 1$. If $f \in A[X_1]$ is nonnull, let $m = \deg(f)$ be its degree. Since A_1 is infinite and $f(\alpha_1) = 0$ for all $\alpha_1 \in A_1$, then f has an infinite number of roots. But since f is of degree m , this contradicts Theorem 20.22. Thus, $f = 0$.

If $n \geq 2$, we can view $f \in A[X_1, \dots, X_n]$ as a polynomial

$$f = g_m X_n^m + g_{m-1} X_n^{m-1} + \dots + g_0,$$

where the coefficients g_i are polynomials in $A[X_1, \dots, X_{n-1}]$. Now, for every $(\alpha_1, \dots, \alpha_{n-1}) \in A_1 \times \dots \times A_{n-1}$, $f(\alpha_1, \dots, \alpha_{n-1}, X_n)$ determines a polynomial $h(X_n) \in A[X_n]$, and since A_n is infinite and $h(\alpha_n) = f(\alpha_1, \dots, \alpha_{n-1}, \alpha_n) = 0$ for all $\alpha_n \in A_n$, by the induction hypothesis, we have $g_i(\alpha_1, \dots, \alpha_{n-1}) = 0$. Now, since A_1, \dots, A_{n-1} are infinite, using the induction hypothesis again, we get $g_i = 0$, which shows that f is the null polynomial. The second part of the proposition follows immediately from the first, by letting $A_i = A$. \square

When A is an infinite integral domain, in particular an infinite field, since the map $f \mapsto f_A$ is injective, we identify the polynomial f with the polynomial function f_A , and we write f_A simply as f .

The following proposition can be very useful to show polynomial identities.

Proposition 20.27. *Let A be an infinite integral domain and $f, g_1, \dots, g_m \in A[X_1, \dots, X_n]$ be polynomials. If the g_i are nonnull polynomials and if*

$$f(\alpha_1, \dots, \alpha_n) = 0 \text{ whenever } g_i(\alpha_1, \dots, \alpha_n) \neq 0 \text{ for all } i, 1 \leq i \leq m,$$

for every $(\alpha_1, \dots, \alpha_n) \in A^n$, then

$$f = 0,$$

i.e., f is the null polynomial.

Proof. If f is not the null polynomial, since the g_i are nonnull and A is an integral domain, then $fg_1 \cdots g_m$ is nonnull. By Proposition 20.26, only the null polynomial maps to the zero function, and thus there must be some $(\alpha_1, \dots, \alpha_n) \in A^n$, such that

$$f(\alpha_1, \dots, \alpha_n)g_1(\alpha_1, \dots, \alpha_n) \cdots g_m(\alpha_1, \dots, \alpha_n) \neq 0,$$

but this contradicts the hypothesis. □

Proposition 20.27 is often called the *principle of extension of algebraic identities*. Another perhaps more illuminating way of stating this proposition is as follows: For any polynomial $g \in A[X_1, \dots, X_n]$, let

$$V(g) = \{(\alpha_1, \dots, \alpha_n) \in A^n \mid g(\alpha_1, \dots, \alpha_n) = 0\},$$

the set of zeros of g . Note that $V(g_1) \cup \cdots \cup V(g_m) = V(g_1 \cdots g_m)$. Then, Proposition 20.27 can be stated as:

If $f(\alpha_1, \dots, \alpha_n) = 0$ for every $(\alpha_1, \dots, \alpha_n) \in A^n - V(g_1 \cdots g_m)$, then $f = 0$.

In other words, if the algebraic identity $f(\alpha_1, \dots, \alpha_n) = 0$ holds on the complement of $V(g_1) \cup \cdots \cup V(g_m) = V(g_1 \cdots g_m)$, then $f(\alpha_1, \dots, \alpha_n) = 0$ holds everywhere in A^n . With this second formulation, we understand better the terminology “principle of extension of algebraic identities.”

Remark: Letting $U(g) = A - V(g)$, the identity $V(g_1) \cup \cdots \cup V(g_m) = V(g_1 \cdots g_m)$ translates to $U(g_1) \cap \cdots \cap U(g_m) = U(g_1 \cdots g_m)$. This suggests to define a topology on A whose basis of open sets consists of the sets $U(g)$. In this topology (called the Zariski topology), the sets of the form $V(g)$ are closed sets. Also, when $g_1, \dots, g_m \in A[X_1, \dots, X_n]$ and $n \geq 2$, understanding the structure of the closed sets of the form $V(g_1) \cap \cdots \cap V(g_m)$ is quite difficult, and it is the object of algebraic geometry (at least, its classical part).



When $f \in A[X_1, \dots, X_n]$ and $n \geq 2$, one should not apply Proposition 20.26 abusively. For example, let

$$f(X, Y) = X^2 + Y^2 - 1,$$

considered as a polynomial in $\mathbb{R}[X, Y]$. Since \mathbb{R} is an infinite field and since

$$f\left(\frac{1-t^2}{1+t^2}, \frac{2t}{1+t^2}\right) = \frac{(1-t^2)^2}{(1+t^2)^2} + \frac{(2t)^2}{(1+t^2)^2} - 1 = 0,$$

for every $t \in \mathbb{R}$, it would be tempting to say that $f = 0$. But what's wrong with the above reasoning is that there are no two infinite subsets R_1, R_2 of \mathbb{R} such that $f(\alpha_1, \alpha_2) = 0$ for all $(\alpha_1, \alpha_2) \in \mathbb{R}^2$. For every $\alpha_1 \in \mathbb{R}$, there are at most two $\alpha_2 \in \mathbb{R}$ such that $f(\alpha_1, \alpha_2) = 0$. What the example shows though, is that a nonnull polynomial $f \in A[X_1, \dots, X_n]$ where $n \geq 2$ can have an infinite number of zeros. This is in contrast with nonnull polynomials in one variables over an infinite field (which have a number of roots bounded by their degree).

We now look at polynomial interpolation.

20.7 Polynomial Interpolation (Lagrange, Newton, Hermite)

Let K be a field. First, we consider the following interpolation problem: Given a sequence $(\alpha_1, \dots, \alpha_{m+1})$ of pairwise distinct scalars in K and any sequence $(\beta_1, \dots, \beta_{m+1})$ of scalars in K , where the β_j are not necessarily distinct, find a polynomial $P(X)$ of degree $\leq m$ such that

$$P(\alpha_1) = \beta_1, \dots, P(\alpha_{m+1}) = \beta_{m+1}.$$

First, observe that if such a polynomial exists, then it is unique. Indeed, this is a consequence of Proposition 20.23. Thus, we just have to find any polynomial of degree $\leq m$. Consider the following so-called *Lagrange polynomials*:

$$L_i(X) = \frac{(X - \alpha_1) \cdots (X - \alpha_{i-1})(X - \alpha_{i+1}) \cdots (X - \alpha_{m+1})}{(\alpha_i - \alpha_1) \cdots (\alpha_i - \alpha_{i-1})(\alpha_i - \alpha_{i+1}) \cdots (\alpha_i - \alpha_{m+1})}.$$

Note that $L(\alpha_i) = 1$ and that $L(\alpha_j) = 0$, for all $j \neq i$. But then,

$$P(X) = \beta_1 L_1 + \cdots + \beta_{m+1} L_{m+1}$$

is the unique desired polynomial, since clearly, $P(\alpha_i) = \beta_i$. Such a polynomial is called a *Lagrange interpolant*. Also note that the polynomials (L_1, \dots, L_{m+1}) form a basis of the vector space of all polynomials of degree $\leq m$. Indeed, if we had

$$\lambda_1 L_1(X) + \cdots + \lambda_{m+1} L_{m+1}(X) = 0,$$

setting X to α_i , we would get $\lambda_i = 0$. Thus, the L_i are linearly independent, and by the previous argument, they are a set of generators. We call (L_1, \dots, L_{m+1}) the *Lagrange basis* (of order $m + 1$).

It is known from numerical analysis that from a computational point of view, the Lagrange basis is not very good. Newton proposed another solution, the method of divided differences.

Consider the polynomial $P(X)$ of degree $\leq m$, called the *Newton interpolant*,

$$P(X) = \lambda_0 + \lambda_1(X - \alpha_1) + \lambda_2(X - \alpha_1)(X - \alpha_2) + \cdots + \lambda_m(X - \alpha_1)(X - \alpha_2) \cdots (X - \alpha_m).$$

Then, the λ_i can be determined by successively setting X to, $\alpha_1, \alpha_2, \dots, \alpha_{m+1}$. More precisely, we define inductively the polynomials $Q(X)$ and $Q(\alpha_1, \dots, \alpha_i, X)$, for $1 \leq i \leq m$, as follows:

$$\begin{aligned}
 Q(X) &= P(X) \\
 Q_1(\alpha_1, X) &= \frac{Q(X) - Q(\alpha_1)}{X - \alpha_1} \\
 Q(\alpha_1, \alpha_2, X) &= \frac{Q(\alpha_1, X) - Q(\alpha_1, \alpha_2)}{X - \alpha_2} \\
 &\dots \\
 Q(\alpha_1, \dots, \alpha_i, X) &= \frac{Q(\alpha_1, \dots, \alpha_{i-1}, X) - Q(\alpha_1, \dots, \alpha_{i-1}, \alpha_i)}{X - \alpha_i}, \\
 &\dots \\
 Q(\alpha_1, \dots, \alpha_m, X) &= \frac{Q(\alpha_1, \dots, \alpha_{m-1}, X) - Q(\alpha_1, \dots, \alpha_{m-1}, \alpha_m)}{X - \alpha_m}.
 \end{aligned}$$

By induction on i , $1 \leq i \leq m-1$, it is easily verified that

$$\begin{aligned}
 Q(X) &= P(X), \\
 Q(\alpha_1, \dots, \alpha_i, X) &= \lambda_i + \lambda_{i+1}(X - \alpha_{i+1}) + \dots + \lambda_m(X - \alpha_{i+1}) \dots (X - \alpha_m), \\
 Q(\alpha_1, \dots, \alpha_m, X) &= \lambda_m.
 \end{aligned}$$

From the above expressions, it is clear that

$$\begin{aligned}
 \lambda_0 &= Q(\alpha_1), \\
 \lambda_i &= Q(\alpha_1, \dots, \alpha_i, \alpha_{i+1}), \\
 \lambda_m &= Q(\alpha_1, \dots, \alpha_m, \alpha_{m+1}).
 \end{aligned}$$

The expression $Q(\alpha_1, \alpha_2, \dots, \alpha_{i+1})$ is called the *i-th difference quotient*. Then, we can compute the λ_i in terms of $\beta_1 = P(\alpha_1), \dots, \beta_{m+1} = P(\alpha_{m+1})$, using the inductive formulae for the $Q(\alpha_1, \dots, \alpha_i, X)$ given above, initializing the $Q(\alpha_i)$ such that $Q(\alpha_i) = \beta_i$.

The above method is called the method of *divided differences* and it is due to Newton.

An astute observation may be used to optimize the computation. Observe that if $P_i(X)$ is the polynomial of degree $\leq i$ taking the values $\beta_1, \dots, \beta_{i+1}$ at the points $\alpha_1, \dots, \alpha_{i+1}$, then the coefficient of X^i in $P_i(X)$ is $Q(\alpha_1, \alpha_2, \dots, \alpha_{i+1})$, which is the value of λ_i in the Newton interpolant

$$P_i(X) = \lambda_0 + \lambda_1(X - \alpha_1) + \lambda_2(X - \alpha_1)(X - \alpha_2) + \dots + \lambda_i(X - \alpha_1)(X - \alpha_2) \dots (X - \alpha_i).$$

As a consequence, $Q(\alpha_1, \alpha_2, \dots, \alpha_{i+1})$ does not depend on the specific ordering of the α_j and there are better ways of computing it. For example, $Q(\alpha_1, \alpha_2, \dots, \alpha_{i+1})$ can be computed

wilder and wilder oscillations around the points $x = -5$ and $x = +5$. This phenomenon becomes quite noticeable beginning for degree 14, and gets worse and worse. For degree 22, things are quite bad! Equivalently, one may consider the function

$$f(x) = \frac{1}{1 + 25x^2},$$

in the interval $[-1, +1]$.

We now consider a more general interpolation problem which will lead to the Hermite polynomials.

We consider the following interpolation problem:

Given a sequence $(\alpha_1, \dots, \alpha_{m+1})$ of pairwise distinct scalars in K , integers n_1, \dots, n_{m+1} where $n_j \geq 0$, and $m + 1$ sequences $(\beta_j^0, \dots, \beta_j^{n_j})$ of scalars in K , letting

$$n = n_1 + \dots + n_{m+1} + m,$$

find a polynomial P of degree $\leq n$, such that

$$\begin{array}{lll} P(\alpha_1) = \beta_1^0, & \dots & P(\alpha_{m+1}) = \beta_{m+1}^0, \\ D^1 P(\alpha_1) = \beta_1^1, & \dots & D^1 P(\alpha_{m+1}) = \beta_{m+1}^1, \\ & \dots & \\ D^i P(\alpha_1) = \beta_1^i, & \dots & D^i P(\alpha_{m+1}) = \beta_{m+1}^i, \\ & \dots & \\ D^{n_1} P(\alpha_1) = \beta_1^{n_1}, & \dots & D^{n_{m+1}} P(\alpha_{m+1}) = \beta_{m+1}^{n_{m+1}}. \end{array}$$

Note that the above equations constitute $n + 1$ constraints, and thus, we can expect that there is a unique polynomial of degree $\leq n$ satisfying the above problem. This is indeed the case and such a polynomial is called a *Hermite polynomial*. We call the above problem the *Hermite interpolation problem*.

Proposition 20.28. *The Hermite interpolation problem has a unique solution of degree $\leq n$, where $n = n_1 + \dots + n_{m+1} + m$.*

Proof. First, we prove that the Hermite interpolation problem has at most one solution. Assume that P and Q are two distinct solutions of degree $\leq n$. Then, by Proposition 20.25 and the criterion following it, $P - Q$ has among its roots α_1 of multiplicity at least $n_1 + 1, \dots, \alpha_{m+1}$ of multiplicity at least $n_{m+1} + 1$. However, by Theorem 20.22, we should have

$$n_1 + 1 + \dots + n_{m+1} + 1 = n_1 + \dots + n_{m+1} + m + 1 \leq n,$$

which is a contradiction, since $n = n_1 + \dots + n_{m+1} + m$. Thus, $P = Q$. We are left with proving the existence of a Hermite interpolant. A quick way to do so is to use Proposition 5.13, which tells us that given a square matrix A over a field K , the following properties hold:

For every column vector B , there is a unique column vector X such that $AX = B$ iff the only solution to $AX = 0$ is the trivial vector $X = 0$ iff $D(A) \neq 0$.

If we let $P = y_0 + y_1X + \cdots + y_nX^n$, the Hermite interpolation problem yields a linear system of equations in the unknowns (y_0, \dots, y_n) with some associated $(n+1) \times (n+1)$ matrix A . Now, the system $AY = 0$ has a solution iff P has among its roots α_1 of multiplicity at least $n_1 + 1, \dots, \alpha_{m+1}$ of multiplicity at least $n_{m+1} + 1$. By the previous argument, since P has degree $\leq n$, we must have $P = 0$, that is, $Y = 0$. This concludes the proof. \square

Proposition 20.28 shows the existence of unique polynomials $H_j^i(X)$ of degree $\leq n$ such that $D^i H_j^i(\alpha_j) = 1$ and $D^k H_j^i(\alpha_l) = 0$, for $k \neq i$ or $l \neq j$, $1 \leq j, l \leq m+1$, $0 \leq i, k \leq n_j$. The polynomials H_j^i are called *Hermite basis polynomials*.

One problem with Proposition 20.28 is that it does not give an explicit way of computing the Hermite basis polynomials. We first show that this can be done explicitly in the special cases $n_1 = \dots = n_{m+1} = 1$, and $n_1 = \dots = n_{m+1} = 2$, and then suggest a method using a generalized Newton interpolant.

Assume that $n_1 = \dots = n_{m+1} = 1$. We try $H_j^0 = (a(X - \alpha_j) + b)L_j^2$, and $H_j^1 = (c(X - \alpha_j) + d)L_j^2$, where L_j is the Lagrange interpolant determined earlier. Since

$$DH_j^0 = aL_j^2 + 2(a(X - \alpha_j) + b)L_jDL_j,$$

requiring that $H_j^0(\alpha_j) = 1$, $H_j^0(\alpha_k) = 0$, $DH_j^0(\alpha_j) = 0$, and $DH_j^0(\alpha_k) = 0$, for $k \neq j$, implies $b = 1$ and $a = -2DL_j(\alpha_j)$. Similarly, from the requirements $H_j^1(\alpha_j) = 0$, $H_j^1(\alpha_k) = 0$, $DH_j^1(\alpha_j) = 1$, and $DH_j^1(\alpha_k) = 0$, $k \neq j$, we get $c = 1$ and $d = 0$.

Thus, we have the Hermite polynomials

$$H_j^0 = (1 - 2DL_j(\alpha_j)(X - \alpha_j))L_j^2, \quad H_j^1 = (X - \alpha_j)L_j^2.$$

In the special case where $m = 1$, $\alpha_1 = 0$, and $\alpha_2 = 1$, we leave as an exercise to show that the Hermite polynomials are

$$\begin{aligned} H_0^0 &= 2X^3 - 3X^2 + 1, \\ H_1^0 &= -2X^3 + 3X^2, \\ H_0^1 &= X^3 - 2X^2 + X, \\ H_1^1 &= X^3 - X^2. \end{aligned}$$

As a consequence, the polynomial P of degree 3 such that $P(0) = x_0$, $P(1) = x_1$, $P'(0) = m_0$, and $P'(1) = m_1$, can be written as

$$P(X) = x_0(2X^3 - 3X^2 + 1) + m_0(X^3 - 2X^2 + X) + m_1(X^3 - X^2) + x_1(-2X^3 + 3X^2).$$

If we want the polynomial P of degree 3 such that $P(a) = x_0$, $P(b) = x_1$, $P'(a) = m_0$, and $P'(b) = m_1$, where $b \neq a$, then we have

$$P(X) = x_0(2t^3 - 3t^2 + 1) + (b - a)m_0(t^3 - 2t^2 + t) + (b - a)m_1(t^3 - t^2) + x_1(-2t^3 + 3t^2),$$

where

$$t = \frac{X - a}{b - a}.$$

Observe the presence of the extra factor $(b - a)$ in front of m_0 and m_1 , the formula would be false otherwise!

We now consider the case where $n_1 = \dots = n_{m+1} = 2$. Let us try

$$H_j^i(X) = (a^i(X - \alpha_j)^2 + b^i(X - \alpha_j) + c^i)L_j^3,$$

where $0 \leq i \leq 2$. Sparing the readers some (tedious) computations, we find:

$$\begin{aligned} H_j^0(X) &= \left((6(DL_j(\alpha_j))^2 - \frac{3}{2}D^2L_j(\alpha_j))(X - \alpha_j)^2 - 3DL_j(\alpha_j)(X - \alpha_j) + 1 \right) L_j^3(X), \\ H_j^1(X) &= \left(9(DL_j(\alpha_j))^2(X - \alpha_j)^2 - 3DL_j(\alpha_j)(X - \alpha_j) \right) L_j^3(X), \\ H_j^2(X) &= \frac{1}{2}(X - \alpha_j)^2 L_j^3(X). \end{aligned}$$

Going back to the general problem, it seems to us that a kind of Newton interpolant will be more manageable. Let

$$\begin{aligned} P_0^0(X) &= 1, \\ P_j^0(X) &= (X - \alpha_1)^{n_1+1} \dots (X - \alpha_j)^{n_j+1}, \quad 1 \leq j \leq m \\ P_0^i(X) &= (X - \alpha_1)^i (X - \alpha_2)^{n_2+1} \dots (X - \alpha_{m+1})^{n_{m+1}+1}, \quad 1 \leq i \leq n_1, \\ P_j^i(X) &= (X - \alpha_1)^{n_1+1} \dots (X - \alpha_j)^{n_j+1} (X - \alpha_{j+1})^i (X - \alpha_{j+2})^{n_{j+2}+1} \dots (X - \alpha_{m+1})^{n_{m+1}+1}, \\ &\quad 1 \leq j \leq m-1, \quad 1 \leq i \leq n_{j+1}, \\ P_m^i(X) &= (X - \alpha_1)^{n_1+1} \dots (X - \alpha_m)^{n_m+1} (X - \alpha_{m+1})^i, \quad 1 \leq i \leq n_{m+1}, \end{aligned}$$

and let

$$P(X) = \sum_{j=0, i=0}^{j=m, i=n_{j+1}} \lambda_j^i P_j^i(X).$$

We can think of $P(X)$ as a generalized Newton interpolant. We can compute the derivatives $D^k P_j^i$, for $1 \leq k \leq n_{j+1}$, and if we look for the Hermite basis polynomials $H_j^i(X)$ such that $D^i H_j^i(\alpha_j) = 1$ and $D^k H_j^i(\alpha_l) = 0$, for $k \neq i$ or $l \neq j$, $1 \leq j, l \leq m+1$, $0 \leq i, k \leq n_j$, we find that we have to solve triangular systems of linear equations. Thus, as in the simple case $n_1 = \dots = n_{m+1} = 0$, we can solve successively for the λ_j^i . Obviously, the computations are quite formidable and we leave such considerations for further study.

Chapter 21

UFD's, Noetherian Rings, Hilbert's Basis Theorem

21.1 Unique Factorization Domains (Factorial Rings)

We saw in Section 20.5 that if K is a field, then every nonnull polynomial in $K[X]$ can be factored as a product of irreducible factors, and that such a factorization is essentially unique. The same property holds for the ring $K[X_1, \dots, X_n]$ where $n \geq 2$, but a different proof is needed.

The reason why unique factorization holds for $K[X_1, \dots, X_n]$ is that if A is an integral domain for which unique factorization holds in some suitable sense, then the property of unique factorization lifts to the polynomial ring $A[X]$. Such rings are called factorial rings, or unique factorization domains. The first step is to define the notion of irreducible element in an integral domain, and then to define a factorial ring. It will turn out that in a factorial ring, any nonnull element a is irreducible (or prime) iff the principal ideal (a) is a prime ideal.

Recall that given a ring A , a *unit* is any invertible element (w.r.t. multiplication). The set of units of A is denoted by A^* . It is a multiplicative subgroup of A , with identity 1. Also, given $a, b \in A$, recall that a *divides* b if $b = ac$ for some $c \in A$; equivalently, a divides b iff $(b) \subseteq (a)$. Any nonzero $a \in A$ is divisible by any unit u , since $a = u(u^{-1}a)$. The relation “ a divides b ,” often denoted by $a \mid b$, is reflexive and transitive, and thus, a preorder on $A - \{0\}$.

Definition 21.1. Let A be an integral domain. Some element $a \in A$ is *irreducible* if $a \neq 0$, $a \notin A^*$ (a is not a unit), and whenever $a = bc$, then either b or c is a unit (where $b, c \in A$). Equivalently, $a \in A$ is *reducible* if $a = 0$, or $a \in A^*$ (a is a unit), or $a = bc$ where $b, c \notin A^*$ (a, b are both noninvertible) and $b, c \neq 0$.

Observe that if $a \in A$ is irreducible and $u \in A$ is a unit, then ua is also irreducible. Generally, if $a \in A$, $a \neq 0$, and u is a unit, then a and ua are said to be *associated*. This is the equivalence relation on nonnull elements of A induced by the divisibility preorder.

The following simple proposition gives a sufficient condition for an element $a \in A$ to be irreducible.

Proposition 21.1. *Let A be an integral domain. For any $a \in A$ with $a \neq 0$, if the principal ideal (a) is a prime ideal, then a is irreducible.*

Proof. If (a) is prime, then $(a) \neq A$ and a is not a unit. Assume that $a = bc$. Then, $bc \in (a)$, and since (a) is prime, either $b \in (a)$ or $c \in (a)$. Consider the case where $b \in (a)$, the other case being similar. Then, $b = ax$ for some $x \in A$. As a consequence,

$$a = bc = axc,$$

and since A is an integral domain and $a \neq 0$, we get

$$1 = xc,$$

which proves that $c = x^{-1}$ is a unit. □

It should be noted that the converse of Proposition 21.1 is generally false. However, it holds for factorial rings, defined next.

Definition 21.2. A *factorial ring* or *unique factorization domain (UFD)* (or *unique factorization ring*) is an integral domain A such that the following two properties hold:

- (1) For every nonnull $a \in A$, if $a \notin A^*$ (a is not a unit), then a can be factored as a product

$$a = a_1 \cdots a_m$$

where each $a_i \in A$ is irreducible ($m \geq 1$).

- (2) For every nonnull $a \in A$, if $a \notin A^*$ (a is not a unit) and if

$$a = a_1 \cdots a_m = b_1 \cdots b_n$$

where $a_i \in A$ and $b_j \in A$ are irreducible, then $m = n$ and there is a permutation σ of $\{1, \dots, m\}$ and some units $u_1, \dots, u_m \in A^*$ such that $a_i = u_i b_{\sigma(i)}$ for all i , $1 \leq i \leq m$.

Example 21.1. The ring \mathbb{Z} of integers is a typical example of a UFD. Given a field K , the polynomial ring $K[X]$ is a UFD. More generally, we will show later that every PID is a UFD (see Theorem 21.12). Thus, in particular, $\mathbb{Z}[X]$ is a UFD. However, we leave as an exercise to prove that the ideal $(2X, X^2)$ generated by $2X$ and X^2 is not principal, and thus, $\mathbb{Z}[X]$ is not a PID.

First, we prove that condition (2) in Definition 21.2 is equivalent to the usual “Euclidean” condition.



There are integral domains that are not UFD's. For example, the subring $\mathbb{Z}[\sqrt{-5}]$ of \mathbb{C} consisting of the complex numbers of the form $a + bi\sqrt{5}$ where $a, b \in \mathbb{Z}$ is not a UFD. Indeed, we have

$$9 = 3 \cdot 3 = (2 + i\sqrt{5})(2 - i\sqrt{5}),$$

and it can be shown that 3, $2 + i\sqrt{5}$, and $2 - i\sqrt{5}$ are irreducible, and that the units are ± 1 . The uniqueness condition (2) fails and $\mathbb{Z}[\sqrt{-5}]$ is not a UFD.

Remark: For $d \in \mathbb{Z}$ with $d < 0$, it is known that the ring of integers of $\mathbb{Q}(\sqrt{d})$ is a UFD iff d is one of the nine primes, $d = -1, -2, -3, -7, -11, -19, -43, -67$ and -163 . This is a hard theorem that was conjectured by Gauss but not proved until 1966, independently by Stark and Baker. Heegner had published a proof of this result in 1952 but there was some doubt about its validity. After finding his proof, Stark reexamined Heegner's proof and concluded that it was essentially correct after all. In sharp contrast, when d is a positive integer, the problem of determining which of the rings of integers of $\mathbb{Q}(\sqrt{d})$ are UFD's, is still open. It can also be shown that if $d < 0$, then the ring $\mathbb{Z}[\sqrt{d}]$ is a UFD iff $d = -1$ or $d = -2$. If $d \equiv 1 \pmod{4}$, then $\mathbb{Z}[\sqrt{d}]$ is never a UFD. For more details about these remarkable results, see Stark [100] (Chapter 8).

Proposition 21.2. *Let A be an integral domain satisfying condition (1) in Definition 21.2. Then, condition (2) in Definition 21.2 is equivalent to the following condition:*

(2') *If $a \in A$ is irreducible and a divides the product bc , where $b, c \in A$ and $b, c \neq 0$, then either a divides b or a divides c .*

Proof. First, assume that (2) holds. Let $bc = ad$, where $d \in A$, $d \neq 0$. If b is a unit, then

$$c = adb^{-1},$$

and c is divisible by a . A similar argument applies to c . Thus, we may assume that b and c are not units. In view of (1), we can write

$$b = p_1 \cdots p_m \quad \text{and} \quad c = p_{m+1} \cdots p_{m+n},$$

where $p_i \in A$ is irreducible. Since $bc = ad$, a is irreducible, and b, c are not units, d cannot be a unit. In view of (1), we can write

$$d = q_1 \cdots q_r,$$

where $q_i \in A$ is irreducible. Thus,

$$p_1 \cdots p_m p_{m+1} \cdots p_{m+n} = a q_1 \cdots q_r,$$

where all the factors involved are irreducible. By (2), we must have

$$a = u_{i_0} p_{i_0}$$

for some unit $u_{i_0} \in A$ and some index i_0 , $1 \leq i_0 \leq m+n$. As a consequence, if $1 \leq i_0 \leq m$, then a divides b , and if $m+1 \leq i_0 \leq m+n$, then a divides c . This proves that (2') holds.

Let us now assume that (2') holds. Assume that

$$a = a_1 \cdots a_m = b_1 \cdots b_n,$$

where $a_i \in A$ and $b_j \in A$ are irreducible. Without loss of generality, we may assume that $m \leq n$. We proceed by induction on m . If $m = 1$,

$$a_1 = b_1 \cdots b_n,$$

and since a_1 is irreducible, $u = b_1 \cdots b_{i-1} b_{i+1} \cdots b_n$ must be a unit for some i , $1 \leq i \leq n$. Thus, (2) holds with $n = 1$ and $a_1 = b_i u$. Assume that $m > 1$ and that the induction hypothesis holds for $m-1$. Since

$$a_1 a_2 \cdots a_m = b_1 \cdots b_n,$$

a_1 divides $b_1 \cdots b_n$, and in view of (2'), a_1 divides some b_j . Since a_1 and b_j are irreducible, we must have $b_j = u_j a_1$, where $u_j \in A$ is a unit. Since A is an integral domain,

$$a_1 a_2 \cdots a_m = b_1 \cdots b_{j-1} u_j a_1 b_{j+1} \cdots b_n$$

implies that

$$a_2 \cdots a_m = (u_j b_1) \cdots b_{j-1} b_{j+1} \cdots b_n,$$

and by the induction hypothesis, $m-1 = n-1$ and $a_i = v_i b_{\tau(i)}$ for some units $v_i \in A$ and some bijection τ between $\{2, \dots, m\}$ and $\{1, \dots, j-1, j+1, \dots, n\}$. However, the bijection τ extends to a permutation σ of $\{1, \dots, m\}$ by letting $\sigma(1) = j$, and the result holds by letting $v_1 = u_j^{-1}$. \square

As a corollary of Proposition 21.2, we get the converse of Proposition 21.1.

Proposition 21.3. *Let A be a factorial ring. For any $a \in A$ with $a \neq 0$, the principal ideal (a) is a prime ideal iff a is irreducible.*

Proof. In view of Proposition 21.1, we just have to prove that if $a \in A$ is irreducible, then the principal ideal (a) is a prime ideal. Indeed, if $bc \in (a)$, then a divides bc , and by Proposition 21.2, property (2') implies that either a divides b or a divides c , that is, either $b \in (a)$ or $c \in (a)$, which means that (a) is prime. \square

Because Proposition 21.3 holds, in a UFD, an irreducible element is often called a *prime*.

In a UFD A , every nonzero element $a \in A$ that is not a unit can be expressed as a product $a = a_1 \cdots a_n$ of irreducible elements a_i , and by property (2), the number n of factors only depends on a , that is, it is the same for all factorizations into irreducible factors. We agree that this number is 0 for a unit.

Remark: If A is a UFD, we can state the factorization properties so that they also applies to units:

- (1) For every nonnull $a \in A$, a can be factored as a product

$$a = ua_1 \cdots a_m$$

where $u \in A^*$ (u is a unit) and each $a_i \in A$ is irreducible ($m \geq 0$).

- (2) For every nonnull $a \in A$, if

$$a = ua_1 \cdots a_m = vb_1 \cdots b_n$$

where $u, v \in A^*$ (u, v are units) and $a_i \in A$ and $b_j \in A$ are irreducible, then $m = n$, and if $m = n = 0$ then $u = v$, else if $m \geq 1$, then there is a permutation σ of $\{1, \dots, m\}$ and some units $u_1, \dots, u_m \in A^*$ such that $a_i = u_i b_{\sigma(i)}$ for all i , $1 \leq i \leq m$.

We are now ready to prove that if A is a UFD, then the polynomial ring $A[X]$ is also a UFD.

The fact that nonnull and nonunit polynomials in $A[X]$ factor as products of irreducible polynomials is rather easy to prove. First, observe that the units of $A[X]$ are just the units of A . If $f(X)$ is a polynomial of degree 0 that is not a unit, the fact that A is a UFD yields the desired factorization of $f(X)$. If $f(X)$ has degree $m > 0$ and $f(X)$ is reducible, then $f(X)$ factors as the product of two nonunit polynomials $g(X), h(X)$. Let f_m be the coefficient of degree m in f . We have

$$f(X) = g(X)h(X),$$

and if both $g(X)$ and $h(X)$ have degree strictly less than m , by induction, we get a factorization of $f(X)$ as a product of irreducible polynomials. Otherwise, either $g(X)$ or $h(X)$ is a constant. Consider the case where $g(X)$ is a constant, the other case being similar. Then, $g(X) = b$ is not a unit, and b factors as a product $b = b_1 \cdots b_n$ of irreducible elements b_i , where n only depends on b . Since

$$f_m = bh_m,$$

where h_m be the coefficient of degree m in h , we see that h_m is a product of p of the b_i 's, up to units, and thus, $p < m$. Again, we conclude by induction. More formally, we can proceed by induction on (m, n) , where m is the degree of $f(X)$ and n is the number of irreducible factors in f_m .

For the uniqueness of the factorization, by Proposition 21.2, it is enough to prove that condition (2') holds. This is a little more tricky. There are several proofs, but they all involve a pretty Lemma due to Gauss.

First, note the following trivial fact. Given a ring A , for any $a \in A$, $a \neq 0$, if a divides every coefficient of some nonnull polynomial $f(X) \in A[X]$, then a divides $f(X)$. If A is an integral domain, we get the following converse.

Proposition 21.4. *Let A be an integral domain. For any $a \in A$, $a \neq 0$, if a divides a nonnull polynomial $f(X) \in A[X]$, then a divides every coefficient of $f(X)$.*

Proof. Assume that $f(X) = ag(X)$, for some $g(X) \in A[X]$. Since $a \neq 0$ and A is an integral ring, $f(X)$ and $g(X)$ have the same degree m , and since for every i ($0 \leq i \leq m$) the coefficient of X^i in $f(X)$ is equal to the coefficient of X^i in $ag(x)$, we have $f_i = ag_i$, and whenever $f_i \neq 0$, we see that a divides f_i . \square

Lemma 21.5. (*Gauss's lemma*) *Let A be a UFD. For any $a \in A$, if a is irreducible and a divides the product $f(X)g(X)$ of two polynomials $f(X), g(X) \in A[X]$, then either a divides $f(X)$ or a divides $g(X)$.*

Proof. Let $f(X) = f_m X^m + \cdots + f_i X^i + \cdots + f_0$ and $g(X) = g_n X^n + \cdots + g_j X^j + \cdots + g_0$. Assume that a divides neither $f(X)$ nor $g(X)$. By the (easy) converse of Proposition 21.4, there is some i ($0 \leq i \leq m$) such that a does not divide f_i , and there is some j ($0 \leq j \leq n$) such that a does not divide g_j . Pick i and j minimal such that a does not divide f_i and a does not divide g_j . The coefficient c_{i+j} of X^{i+j} in $f(X)g(X)$ is

$$c_{i+j} = f_0 g_{i+j} + f_1 g_{i+j-1} + \cdots + f_i g_j + \cdots + f_{i+j} g_0$$

(letting $f_h = 0$ if $h > m$ and $g_k = 0$ if $k > n$). From the choice of i and j , a cannot divide $f_i g_j$, since a being irreducible, by (2') of Proposition 21.2, a would divide f_i or g_j . However, by the choice of i and j , a divides every other nonnull term in the sum for c_{i+j} , and since a is irreducible and divides $f(X)g(X)$, by Proposition 21.4, a divides c_{i+j} , which implies that a divides $f_i g_j$, a contradiction. Thus, either a divides $f(X)$ or a divides $g(X)$. \square

As a corollary, we get the following proposition.

Proposition 21.6. *Let A be a UFD. For any $a \in A$, $a \neq 0$, if a divides the product $f(X)g(X)$ of two polynomials $f(X), g(X) \in A[X]$ and $f(X)$ is irreducible and of degree at least 1, then a divides $g(X)$.*

Proof. The Proposition is trivial if a is a unit. Otherwise, $a = a_1 \cdots a_m$ where $a_i \in A$ is irreducible. Using induction and applying Lemma 21.5, we conclude that a divides $g(X)$. \square

We now show that Lemma 21.5 also applies to the case where a is an irreducible polynomial. This requires a little excursion involving the fraction field F of A .

Remark: If A is a UFD, it is possible to prove the uniqueness condition (2) for $A[X]$ directly without using the fraction field of A , see Malliavin [75], Chapter 3.

Given an integral domain A , we can construct a field F such that every element of F is of the form a/b , where $a, b \in A$, $b \neq 0$, using essentially the method for constructing the field \mathbb{Q} of rational numbers from the ring \mathbb{Z} of integers.

Proposition 21.7. *Let A be an integral domain.*

- (1) *There is a field F and an injective ring homomorphism $i: A \rightarrow F$ such that every element of F is of the form $i(a)i(b)^{-1}$, where $a, b \in A$, $b \neq 0$.*

- (2) For every field K and every injective ring homomorphism $h: A \rightarrow K$, there is a (unique) field homomorphism $\widehat{h}: F \rightarrow K$ such that

$$\widehat{h}(i(a)i(b)^{-1}) = h(a)h(b)^{-1}$$

for all $a, b \in A, b \neq 0$.

- (3) The field F in (1) is unique up to isomorphism.

Proof. (1) Consider the binary relation \simeq on $A \times (A - \{0\})$ defined as follows:

$$(a, b) \simeq (a', b') \quad \text{iff} \quad ab' = a'b.$$

It is easily seen that \simeq is an equivalence relation. Note that the fact that A is an integral domain is used to prove transitivity. The equivalence class of (a, b) is denoted by a/b . Clearly, $(0, b) \simeq (0, 1)$ for all $b \in A$, and we denote the class of $(0, 1)$ also by 0 . The equivalence class $a/1$ of $(a, 1)$ is also denoted by a . We define addition and multiplication on $A \times (A - \{0\})$ as follows:

$$\begin{aligned} (a, b) + (a', b') &= (ab' + a'b, bb'), \\ (a, b) \cdot (a', b') &= (aa', bb'). \end{aligned}$$

It is easily verified that \simeq is congruential w.r.t. $+$ and \cdot , which means that $+$ and \cdot are well-defined on equivalence classes modulo \simeq . When $a, b \neq 0$, the inverse of a/b is b/a , and it is easily verified that F is a field. The map $i: A \rightarrow F$ defined such that $i(a) = a/1$ is an injection of A into F and clearly

$$\frac{a}{b} = i(a)i(b)^{-1}.$$

- (2) Given an injective ring homomorphism $h: A \rightarrow K$ into a field K ,

$$\frac{a}{b} = \frac{a'}{b'} \quad \text{iff} \quad ab' = a'b,$$

which implies that

$$h(a)h(b') = h(a')h(b),$$

and since h is injective and $b, b' \neq 0$, we get

$$h(a)h(b)^{-1} = h(a')h(b')^{-1}.$$

Thus, there is a map $\widehat{h}: F \rightarrow K$ such that

$$\widehat{h}(a/b) = \widehat{h}(i(a)i(b)^{-1}) = h(a)h(b)^{-1}$$

for all $a, b \in A, b \neq 0$, and it is easily checked that \widehat{h} is a field homomorphism. The map \widehat{h} is clearly unique.

- (3) The uniqueness of F up to isomorphism follows from (2), and is left as an exercise. \square

The field F given by Proposition 21.7 is called the *fraction field* of A , and it is denoted by $\text{Frac}(A)$.

In particular, given an integral domain A , since $A[X_1, \dots, X_n]$ is also an integral domain, we can form the fraction field of the polynomial ring $A[X_1, \dots, X_n]$, denoted by $F(X_1, \dots, X_n)$, where $F = \text{Frac}(A)$ is the fraction field of A . It is also called the field of *rational functions* over F , although the terminology is a bit misleading, since elements of $F(X_1, \dots, X_n)$ only define functions when the dominator is nonnull.

We now have the following crucial lemma which shows that if a polynomial $f(X)$ is reducible over $F[X]$ where F is the fraction field of A , then $f(X)$ is already reducible over $A[X]$.

Lemma 21.8. *Let A be a UFD and let F be the fraction field of A . For any nonnull polynomial $f(X) \in A[X]$ of degree m , if $f(X)$ is not the product of two polynomials of degree strictly smaller than m , then $f(X)$ is irreducible in $F[X]$.*

Proof. Assume that $f(X)$ is reducible in $F[X]$ and that $f(X)$ is neither null nor a unit. Then,

$$f(X) = G(X)H(X),$$

where $G(X), H(X) \in F[X]$ are polynomials of degree $p, q \geq 1$. Let a be the product of the denominators of the coefficients of $G(X)$, and b the product of the denominators of the coefficients of $H(X)$. Then, $a, b \neq 0$, $g_1(X) = aG(X) \in A[X]$ has degree $p \geq 1$, $h_1(X) = bH(X) \in A[X]$ has degree $q \geq 1$, and

$$abf(X) = g_1(X)h_1(X).$$

Let $c = ab$. If c is a unit, then $f(X)$ is also reducible in $A[X]$. Otherwise, $c = c_1 \cdots c_n$, where $c_i \in A$ is irreducible. We now use induction on n to prove that

$$f(X) = g(X)h(X),$$

for some polynomials $g(X) \in A[X]$ of degree $p \geq 1$ and $h(X) \in A[X]$ of degree $q \geq 1$.

If $n = 1$, since $c = c_1$ is irreducible, by Lemma 21.5, either c divides $g_1(X)$ or c divides $h_1(X)$. Say that c divides $g_1(X)$, the other case being similar. Then, $g_1(X) = cg(X)$ for some $g(X) \in A[X]$ of degree $p \geq 1$, and since $A[X]$ is an integral ring, we get

$$f(X) = g(X)h_1(X),$$

showing that $f(X)$ is reducible in $A[X]$. If $n > 1$, since

$$c_1 \cdots c_n f(X) = g_1(X)h_1(X),$$

c_1 divides $g_1(X)h_1(X)$, and as above, either c_1 divides $g_1(X)$ or c divides $h_1(X)$. In either case, we get

$$c_2 \cdots c_n f(X) = g_2(X)h_2(X)$$

for some polynomials $g_2(X) \in A[X]$ of degree $p \geq 1$ and $h_2(X) \in A[X]$ of degree $q \geq 1$. By the induction hypothesis, we get

$$f(X) = g(X)h(X),$$

for some polynomials $g(X) \in A[X]$ of degree $p \geq 1$ and $h(X) \in A[X]$ of degree $q \geq 1$, showing that $f(X)$ is reducible in $A[X]$. \square

Finally, we can prove that (2') holds.

Lemma 21.9. *Let A be a UFD. Given any three nonnull polynomials $f(X), g(X), h(X) \in A[X]$, if $f(X)$ is irreducible and $f(X)$ divides the product $g(X)h(X)$, then either $f(X)$ divides $g(X)$ or $f(X)$ divides $h(X)$.*

Proof. If $f(X)$ has degree 0, then the result follows from Lemma 21.5. Thus, we may assume that the degree of $f(X)$ is $m \geq 1$. Let F be the fraction field of A . By Lemma 21.8, $f(X)$ is also irreducible in $F[X]$. Since $F[X]$ is a UFD (by Theorem 20.16), either $f(X)$ divides $g(X)$ or $f(X)$ divides $h(X)$, in $F[X]$. Assume that $f(X)$ divides $g(X)$, the other case being similar. Then,

$$g(X) = f(X)G(X),$$

for some $G(X) \in F[X]$. If a is the product the denominators of the coefficients of G , we have

$$ag(X) = q_1(X)f(X),$$

where $q_1(X) = aG(X) \in A[X]$. If a is a unit, we see that $f(X)$ divides $g(X)$. Otherwise, $a = a_1 \cdots a_n$, where $a_i \in A$ is irreducible. We prove by induction on n that

$$g(X) = q(X)f(X)$$

for some $q(X) \in A[X]$.

If $n = 1$, since $f(X)$ is irreducible and of degree $m \geq 1$ and

$$a_1g(X) = q_1(X)f(X),$$

by Lemma 21.5, a_1 divides $q_1(X)$. Thus, $q_1(X) = a_1q(X)$ where $q(X) \in A[X]$. Since $A[X]$ is an integral domain, we get

$$g(X) = q(X)f(X),$$

and $f(X)$ divides $g(X)$. If $n > 1$, from

$$a_1 \cdots a_n g(X) = q_1(X)f(X),$$

we note that a_1 divides $q_1(X)f(X)$, and as in the previous case, a_1 divides $q_1(X)$. Thus, $q_1(X) = a_1q_2(X)$ where $q_2(X) \in A[X]$, and we get

$$a_2 \cdots a_n g(X) = q_2(X)f(X).$$

By the induction hypothesis, we get

$$g(X) = q(X)f(X)$$

for some $q(X) \in A[X]$, and $f(X)$ divides $g(X)$. \square

We finally obtain the fact that $A[X]$ is a UFD when A is.

Theorem 21.10. *If A is a UFD then the polynomial ring $A[X]$ is also a UFD.*

Proof. As we said earlier, the factorization property (1) is easy to prove. Assume that $f(X)$ has degree m and that its coefficient f_m of degree m is the product of n irreducible elements (where $n = 0$ if f_m is a unit). We proceed by induction on the pair (m, n) , using the well-founded ordering on pairs, i.e.,

$$(m, n) \leq (m', n')$$

iff either $m < m'$, or $m = m'$ and $n < n'$. If $f(X)$ is a nonnull polynomial of degree 0 which is not a unit, then $f(X) \in A$, and $f(X) = f_m = a_1 \cdots a_n$ for some irreducible $a_i \in A$, since A is a UFD. If $f(X)$ has degree $m > 0$ and is reducible, then

$$f(X) = g(X)h(X),$$

where $g(X)$ and $h(X)$ have degree $p, q \leq m$ and are not units. If $p, q < m$, then $(p, n_1) < (m, n)$ and $(q, n_2) < (m, n)$, where n_1 is the number of irreducible factors in g_p and n_2 is the number of irreducible factors in h_q , and by the induction hypothesis, both $g(X)$ and $h(X)$ can be written as products of irreducible factors. If $p = 0$, then $g(X) = g_0$ is not a unit, and since

$$f_m = g_0 h_m,$$

h_m is a product of n_2 irreducible elements where $n_2 < n$. Since $(m, n_2) < (m, n)$, by the induction hypothesis, $h(X)$ can be written as products of irreducible polynomials. Since $g_0 \in A$ is not a unit, g_0 can also be factored as a product of irreducible elements. The case where $q = 0$ is similar.

Property (2') follows by Lemma 21.9. By Proposition 21.2, $A[X]$ is a UFD. \square

As a corollary of Theorem 21.10 and using induction, we note that for any field K , the polynomial ring $K[X_1, \dots, X_n]$ is a UFD.

For the sake of completeness, we shall prove that every PID is a UFD. First, we review the notion of gcd and the characterization of gcd's in a PID.

Given an integral domain A , for any two elements $a, b \in A$, $a, b \neq 0$, we say that $d \in A$ ($d \neq 0$) is a *greatest common divisor (gcd)* of a and b if

- (1) d divides both a and b .

(2) For any $h \in A$ ($h \neq 0$), if h divides both a and b , then h divides d .

We also say that a and b are *relatively prime* if 1 is a gcd of a and b .

Note that a and b are relatively prime iff every gcd of a and b is a unit. If A is a PID, then gcd's are characterized as follows.

Proposition 21.11. *Let A be a PID.*

(1) *For any $a, b, d \in A$ ($a, b, d \neq 0$), d is a gcd of a and b iff*

$$(d) = (a, b) = (a) + (b),$$

i.e., d generates the principal ideal generated by a and b .

(2) *(Bezout identity) Two nonnull elements $a, b \in A$ are relatively prime iff there are some $x, y \in A$ such that*

$$ax + by = 1.$$

Proof. (1) Recall that the ideal generated by a and b is the set

$$(a) + (b) = aA + bA = \{ax + by \mid x, y \in A\}.$$

First, assume that d is a gcd of a and b . If so, $a \in Ad$, $b \in Ad$, and thus, $(a) \subseteq (d)$ and $(b) \subseteq (d)$, so that

$$(a) + (b) \subseteq (d).$$

Since A is a PID, there is some $t \in A$, $t \neq 0$, such that

$$(a) + (b) = (t),$$

and thus, $(a) \subseteq (t)$ and $(b) \subseteq (t)$, which means that t divides both a and b . Since d is a gcd of a and b , t must divide d . But then,

$$(d) \subseteq (t) = (a) + (b),$$

and thus, $(d) = (a) + (b)$.

Assume now that

$$(d) = (a) + (b) = (a, b).$$

Since $(a) \subseteq (d)$ and $(b) \subseteq (d)$, d divides both a and b . Assume that t divides both a and b , so that $(a) \subseteq (t)$ and $(b) \subseteq (t)$. Then,

$$(d) = (a) + (b) \subseteq (t),$$

which means that t divides d , and d is indeed a gcd of a and b .

(2) By (1), if a and b are relatively prime, then

$$(1) = (a) + (b),$$

which yields the result. Conversely, if

$$ax + by = 1,$$

then

$$(1) = (a) + (b),$$

and 1 is a gcd of a and b . □

Given two nonnull elements $a, b \in A$, if a is an irreducible element and a does not divide b , then a and b are relatively prime. Indeed, if d is not a unit and d divides both a and b , then $a = dp$ and $b = dq$ where p must be a unit, so that

$$b = ap^{-1}q,$$

and a divides b , a contradiction.

Theorem 21.12. *Let A be ring. If A is a PID, then A is a UFD.*

Proof. First, we prove that every any nonnull element that is a not a unit can be factored as a product of irreducible elements. Let \mathcal{S} be the set of nontrivial principal ideals (a) such that $a \neq 0$ is not a unit and cannot be factored as a product of irreducible elements. Assume that \mathcal{S} is nonempty. We claim that every ascending chain in \mathcal{S} is finite. Otherwise, consider an infinite ascending chain

$$(a_1) \subset (a_2) \subset \cdots \subset (a_n) \subset \cdots .$$

It is immediately verified that

$$\bigcup_{n \geq 1} (a_n)$$

is an ideal in A . Since A is a PID,

$$\bigcup_{n \geq 1} (a_n) = (a)$$

for some $a \in A$. However, there must be some n such that $a \in (a_n)$, and thus,

$$(a_n) \subseteq (a) \subseteq (a_n),$$

and the chain stabilizes at (a_n) . As a consequence, for any ideal (d) such that

$$(a_n) \subset (d)$$

and $(a_n) \neq (d)$, d has the desired factorization. Observe that a_n is not irreducible, since $(a_n) \in \mathcal{S}$, and thus,

$$a_n = bc$$

for some $b, c \in A$, where neither b nor c is a unit. Then,

$$(a_n) \subseteq (b) \quad \text{and} \quad (a_n) \subseteq (c).$$

If $(a_n) = (b)$, then $b = a_n u$ for some $u \in A$, and then

$$a_n = a_n u c,$$

so that

$$1 = uc,$$

since A is an integral domain, and thus, c is a unit, a contradiction. Thus, $(a_n) \neq (b)$, and similarly, $(a_n) \neq (c)$. But then, both b and c factor as products of irreducible elements and so does $a_n = bc$, a contradiction. This implies that $\mathcal{S} = \emptyset$.

To prove the uniqueness of factorizations, we use Proposition 21.2. Assume that a is irreducible and that a divides bc . If a does not divide b , by a previous remark, a and b are relatively prime, and by Proposition 21.11, there are some $x, y \in A$ such that

$$ax + by = 1.$$

Thus,

$$acx + bcy = c,$$

and since a divides bc , we see that a must divide c , as desired. \square

Thus, we get another justification of the fact that \mathbb{Z} is a UFD and that if K is a field, then $K[X]$ is a UFD.

It should also be noted that in a UFD, gcd's of nonnull elements always exist. Indeed, this is trivial if a or b is a unit, and otherwise, we can write

$$a = p_1 \cdots p_m \quad \text{and} \quad b = q_1 \cdots q_n$$

where $p_i, q_j \in A$ are irreducible, and the product of the common factors of a and b is a gcd of a and b (it is 1 if there are no common factors).

We conclude this section on UFD's by proving a proposition characterizing when a UFD is a PID. The proof is nontrivial and makes use of Zorn's lemma (several times).

Proposition 21.13. *Let A be a ring that is a UFD, and not a field. Then, A is a PID iff every nonzero prime ideal is maximal.*

Proof. Assume that A is a PID that is not a field. Consider any nonzero prime ideal, (p) , and pick any proper ideal \mathfrak{A} in A such that

$$(p) \subseteq \mathfrak{A}.$$

Since A is a PID, the ideal \mathfrak{A} is a principal ideal, so $\mathfrak{A} = (q)$, and since \mathfrak{A} is a proper nonzero ideal, $q \neq 0$ and q is not a unit. Since

$$(p) \subseteq (q),$$

q divides p , and we have $p = qp_1$ for some $p_1 \in A$. Now, by Proposition 21.1, since $p \neq 0$ and (p) is a prime ideal, p is irreducible. But then, since $p = qp_1$ and p is irreducible, p_1 must be a unit (since q is not a unit), which implies that

$$(p) = (q);$$

that is, (p) is a maximal ideal.

Conversely, let us assume that every nonzero prime ideal is maximal. First, we prove that every prime ideal is principal. This is obvious for (0) . If \mathfrak{A} is a nonzero prime ideal, then, by hypothesis, it is maximal. Since $\mathfrak{A} \neq (0)$, there is some nonzero element $a \in \mathfrak{A}$. Since \mathfrak{A} is maximal, a is not a unit, and since A is a UFD, there is a factorization $a = a_1 \cdots a_n$ of a into irreducible elements. Since \mathfrak{A} is prime, we have $a_i \in \mathfrak{A}$ for some i . Now, by Proposition 21.3, since a_i is irreducible, the ideal (a_i) is prime, and so, by hypothesis, (a_i) is maximal. Since $(a_i) \subseteq \mathfrak{A}$ and (a_i) is maximal, we get $\mathfrak{A} = (a_i)$.

Next, assume that A is not a PID. Define the set, \mathcal{F} , by

$$\mathcal{F} = \{\mathfrak{A} \mid \mathfrak{A} \subseteq A, \mathfrak{A} \text{ is not a principal ideal}\}.$$

Since A is not a PID, the set \mathcal{F} is nonempty. Also, the reader will easily check that every chain in \mathcal{F} is bounded. Then, by Zorn's lemma (Lemma 31.1), the set \mathcal{F} has some maximal element, \mathfrak{A} . Clearly, $\mathfrak{A} \neq (0)$ is a proper ideal (since $A = (1)$), and \mathfrak{A} is not prime, since we just showed that prime ideals are principal. Then, by Theorem 31.3, there is some maximal ideal, \mathfrak{M} , so that $\mathfrak{A} \subset \mathfrak{M}$. However, a maximal ideal is prime, and we have shown that a prime ideal is principal. Thus,

$$\mathfrak{A} \subseteq (p),$$

for some $p \in A$ that is not a unit. Moreover, by Proposition 21.1, the element p is irreducible. Define

$$\mathfrak{B} = \{a \in A \mid pa \in \mathfrak{A}\}.$$

Clearly, $\mathfrak{A} = p\mathfrak{B}$, $\mathfrak{B} \neq (0)$, $\mathfrak{A} \subseteq \mathfrak{B}$, and \mathfrak{B} is a proper ideal. We claim that $\mathfrak{A} \neq \mathfrak{B}$. Indeed, if $\mathfrak{A} = \mathfrak{B}$ were true, then we would have $\mathfrak{A} = p\mathfrak{B} = \mathfrak{B}$, but this is impossible since p is irreducible, A is a UFD, and $\mathfrak{B} \neq (0)$ (we get $\mathfrak{B} = p^m\mathfrak{B}$ for all m , and every element of \mathfrak{B} would be a multiple of p^m for arbitrarily large m , contradicting the fact that A is a UFD). Thus, we have $\mathfrak{A} \subset \mathfrak{B}$, and since \mathfrak{A} is a maximal element of \mathcal{F} , we must have $\mathfrak{B} \notin \mathcal{F}$. However, $\mathfrak{B} \notin \mathcal{F}$ means that \mathfrak{B} is a principal ideal, and thus, $\mathfrak{A} = p\mathfrak{B}$ is also a principal ideal, a contradiction. \square

Observe that the above proof shows that Proposition 21.13 also holds under the assumption that every prime ideal is principal.

21.2 The Chinese Remainder Theorem

In this section, which is a bit of an interlude, we prove a basic result about quotients of commutative rings by products of ideals that are pairwise relatively prime. This result has applications in number theory and in the structure theorem for finitely generated modules over a PID, which will be presented later.

Given two ideals \mathfrak{a} and \mathfrak{b} of a ring A , we define the ideal $\mathfrak{a}\mathfrak{b}$ as the set of all finite sums of the form

$$a_1b_1 + \cdots + a_kb_k, \quad a_i \in \mathfrak{a}, b_i \in \mathfrak{b}.$$

The reader should check that $\mathfrak{a}\mathfrak{b}$ is indeed an ideal. Observe that $\mathfrak{a}\mathfrak{b} \subseteq \mathfrak{a}$ and $\mathfrak{a}\mathfrak{b} \subseteq \mathfrak{b}$, so that

$$\mathfrak{a}\mathfrak{b} \subseteq \mathfrak{a} \cap \mathfrak{b}.$$

In general, equality does not hold. However if

$$\mathfrak{a} + \mathfrak{b} = A,$$

then we have

$$\mathfrak{a}\mathfrak{b} = \mathfrak{a} \cap \mathfrak{b}.$$

This is because there is some $a \in \mathfrak{a}$ and some $b \in \mathfrak{b}$ such that

$$a + b = 1,$$

so for every $x \in \mathfrak{a} \cap \mathfrak{b}$, we have

$$x = xa + xb,$$

which shows that $x \in \mathfrak{a}\mathfrak{b}$. Ideals \mathfrak{a} and \mathfrak{b} of A that satisfy the condition $\mathfrak{a} + \mathfrak{b} = A$ are sometimes said to be *comaximal*.

We define the homomorphism $\varphi: A \rightarrow A/\mathfrak{a} \times A/\mathfrak{b}$ by

$$\varphi(x) = (\bar{x}_{\mathfrak{a}}, \bar{x}_{\mathfrak{b}}),$$

where $\bar{x}_{\mathfrak{a}}$ is the equivalence class of x modulo \mathfrak{a} (resp. $\bar{x}_{\mathfrak{b}}$ is the equivalence class of x modulo \mathfrak{b}). Recall that the ideal \mathfrak{a} defines the equivalence relation $\equiv_{\mathfrak{a}}$ on A given by

$$x \equiv_{\mathfrak{a}} y \quad \text{iff} \quad x - y \in \mathfrak{a},$$

and that A/\mathfrak{a} is the quotient ring of equivalence classes $\bar{x}_{\mathfrak{a}}$, where $x \in A$, and similarly for A/\mathfrak{b} . Sometimes, we also write $x \equiv y \pmod{\mathfrak{a}}$ for $x \equiv_{\mathfrak{a}} y$.

Clearly, the kernel of the homomorphism φ is $\mathfrak{a} \cap \mathfrak{b}$. If we assume that $\mathfrak{a} + \mathfrak{b} = A$, then $\text{Ker}(\varphi) = \mathfrak{a} \cap \mathfrak{b} = \mathfrak{a}\mathfrak{b}$, and because φ has a constant value on the equivalence classes modulo $\mathfrak{a}\mathfrak{b}$, the map φ induces a quotient homomorphism

$$\theta: A/\mathfrak{a}\mathfrak{b} \rightarrow A/\mathfrak{a} \times A/\mathfrak{b}.$$

Because $\text{Ker}(\varphi) = \mathfrak{a}\mathfrak{b}$, the homomorphism θ is injective. The Chinese Remainder Theorem says that θ is an isomorphism.

Theorem 21.14. *Given a commutative ring A , let \mathfrak{a} and \mathfrak{b} be any two ideals of A such that $\mathfrak{a} + \mathfrak{b} = A$. Then, the homomorphism $\theta: A/\mathfrak{a}\mathfrak{b} \rightarrow A/\mathfrak{a} \times A/\mathfrak{b}$ is an isomorphism.*

Proof. We already showed that θ is injective, so we need to prove that θ is surjective. We need to prove that for any $y, z \in A$, there is some $x \in A$ such that

$$\begin{aligned} x &\equiv y \pmod{\mathfrak{a}} \\ x &\equiv z \pmod{\mathfrak{b}}. \end{aligned}$$

Since $\mathfrak{a} + \mathfrak{b} = A$, there exist some $a \in \mathfrak{a}$ and some $b \in \mathfrak{b}$ such that

$$a + b = 1.$$

If we let

$$x = az + by,$$

then we have

$$x \equiv_{\mathfrak{a}} by \equiv_{\mathfrak{a}} (1 - a)y \equiv_{\mathfrak{a}} y - ay \equiv_{\mathfrak{a}} y,$$

and similarly

$$x \equiv_{\mathfrak{b}} az \equiv_{\mathfrak{b}} (1 - b)z \equiv_{\mathfrak{b}} z - bz \equiv_{\mathfrak{b}} z,$$

which shows that $x = az + by$ works. □

Theorem 21.14 can be generalized to any (finite) number of ideals.

Theorem 21.15. *(Chinese Remainder Theorem) Given a commutative ring A , let $\mathfrak{a}_1, \dots, \mathfrak{a}_n$ be any $n \geq 2$ ideals of A such that $\mathfrak{a}_i + \mathfrak{a}_j = A$ for all $i \neq j$. Then, the homomorphism $\theta: A/\mathfrak{a}_1 \cdots \mathfrak{a}_n \rightarrow A/\mathfrak{a}_1 \times \cdots \times A/\mathfrak{a}_n$ is an isomorphism.*

Proof. The map $\theta: A/\mathfrak{a}_1 \cap \cdots \cap \mathfrak{a}_n \rightarrow A/\mathfrak{a}_1 \times \cdots \times A/\mathfrak{a}_n$ is induced by the homomorphism $\varphi: A \rightarrow A/\mathfrak{a}_1 \times \cdots \times A/\mathfrak{a}_n$ given by

$$\varphi(x) = (\bar{x}_{\mathfrak{a}_1}, \dots, \bar{x}_{\mathfrak{a}_n}).$$

Clearly, $\text{Ker}(\varphi) = \mathfrak{a}_1 \cap \cdots \cap \mathfrak{a}_n$, so θ is well-defined and injective. We need to prove that

$$\mathfrak{a}_1 \cap \cdots \cap \mathfrak{a}_n = \mathfrak{a}_1 \cdots \mathfrak{a}_n$$

and that θ is surjective. We proceed by induction. The case $n = 2$ is Theorem 21.14. By induction, assume that

$$\mathfrak{a}_2 \cap \cdots \cap \mathfrak{a}_n = \mathfrak{a}_2 \cdots \mathfrak{a}_n.$$

We claim that

$$\mathfrak{a}_1 + \mathfrak{a}_2 \cdots \mathfrak{a}_n = A.$$

Indeed, since $\mathfrak{a}_1 + \mathfrak{a}_i = A$ for $i = 2, \dots, n$, there exist some $a_i \in \mathfrak{a}_1$ and some $b_i \in \mathfrak{a}_i$ such that

$$a_i + b_i = 1, \quad i = 2, \dots, n,$$

and by multiplying these equations, we get

$$a + b_2 \cdots b_n = 1,$$

where a is a sum of terms each containing some a_j as a factor, so $a \in \mathfrak{a}_1$ and $b_2 \cdots b_n \in \mathfrak{a}_2 \cdots \mathfrak{a}_n$, which shows that

$$\mathfrak{a}_1 + \mathfrak{a}_2 \cdots \mathfrak{a}_n = A,$$

as claimed. It follows that

$$\mathfrak{a}_1 \cap \mathfrak{a}_2 \cap \cdots \cap \mathfrak{a}_n = \mathfrak{a}_1 \cap (\mathfrak{a}_2 \cdots \mathfrak{a}_n) = \mathfrak{a}_1 \mathfrak{a}_2 \cdots \mathfrak{a}_n.$$

Let us now prove that θ is surjective by induction. The case $n = 2$ is Theorem 21.14. Let x_1, \dots, x_n be any $n \geq 3$ elements of A . First, applying Theorem 21.14 to \mathfrak{a}_1 and $\mathfrak{a}_2 \cdots \mathfrak{a}_n$, we can find $y_1 \in A$ such that

$$\begin{aligned} y_1 &\equiv 1 \pmod{\mathfrak{a}_1} \\ y_1 &\equiv 0 \pmod{\mathfrak{a}_2 \cdots \mathfrak{a}_n}. \end{aligned}$$

By the induction hypothesis, we can find $y_2, \dots, y_n \in A$ such that for all i, j with $2 \leq i, j \leq n$,

$$\begin{aligned} y_i &\equiv 1 \pmod{\mathfrak{a}_i} \\ y_i &\equiv 0 \pmod{\mathfrak{a}_j}, \quad j \neq i. \end{aligned}$$

We claim that

$$x = x_1 y_1 + x_2 y_2 + \cdots + x_n y_n$$

works. Indeed, using the above congruences, for $i = 2, \dots, n$, we get

$$x \equiv x_1 y_1 + x_i \pmod{\mathfrak{a}_i}, \tag{*}$$

but since $\mathfrak{a}_2 \cdots \mathfrak{a}_n \subseteq \mathfrak{a}_i$ for $i = 2, \dots, n$ and $y_1 \equiv 0 \pmod{\mathfrak{a}_2 \cdots \mathfrak{a}_n}$, we have

$$x_1 y_1 \equiv 0 \pmod{\mathfrak{a}_i}, \quad i = 2, \dots, n$$

and equation (*) reduces to

$$x \equiv x_i \pmod{\mathfrak{a}_i}, \quad i = 2, \dots, n.$$

For $i = 1$, we get

$$x \equiv x_1 \pmod{\mathfrak{a}_1},$$

therefore

$$x \equiv x_i \pmod{\mathfrak{a}_i}, \quad i = 1, \dots, n.$$

proving surjectivity. \square

The classical version of the Chinese Remainder Theorem is the case where $A = \mathbb{Z}$ and where the ideals \mathfrak{a}_i are defined by n pairwise relatively prime integers m_1, \dots, m_n . By the Bezout identity, since m_i and m_j are relatively prime whenever $i \neq j$, there exist some $u_i, u_j \in \mathbb{Z}$ such that $u_i m_i + u_j m_j = 1$, and so $m_i \mathbb{Z} + m_j \mathbb{Z} = \mathbb{Z}$. In this case, we get an isomorphism

$$\mathbb{Z}/(m_1 \cdots m_n)\mathbb{Z} \approx \prod_{i=1}^n \mathbb{Z}/m_i \mathbb{Z}.$$

In particular, if m is an integer greater than 1 and

$$m = \prod_i p_i^{r_i}$$

is its factorization into prime factors, then

$$\mathbb{Z}/m\mathbb{Z} \approx \prod_i \mathbb{Z}/p_i^{r_i}\mathbb{Z}.$$

In the previous situation where the integers m_1, \dots, m_n are pairwise relatively prime, if we write $m = m_1 \cdots m_n$ and $m'_i = m/m_i$ for $i = 1, \dots, n$, then m_i and m'_i are relatively prime, and so m'_i has an inverse modulo m_i . If t_i is such an inverse, so that

$$m'_i t_i \equiv 1 \pmod{m_i},$$

then it is not hard to show that for any $a_1, \dots, a_n \in \mathbb{Z}$,

$$x = a_1 t_1 m'_1 + \cdots + a_n t_n m'_n$$

satisfies the congruences

$$x \equiv a_i \pmod{m_i}, \quad i = 1, \dots, n.$$

Theorem 21.15 can be used to characterize rings isomorphic to finite products of quotient rings. Such rings play a role in the structure theorem for torsion modules over a PID.

Given n rings A_1, \dots, A_n , recall that the product ring $A = A_1 \times \cdots \times A_n$ is the ring in which addition and multiplication are defined componenwise. That is,

$$\begin{aligned} (a_1, \dots, a_n) + (b_1, \dots, b_n) &= (a_1 + b_1, \dots, a_n + b_n) \\ (a_1, \dots, a_n) \cdot (b_1, \dots, b_n) &= (a_1 b_1, \dots, a_n b_n). \end{aligned}$$

The additive identity is $0_A = (0, \dots, 0)$ and the multiplicative identity is $1_A = (1, \dots, 1)$. Then, for $i = 1, \dots, n$, we can define the element $e_i \in A$ as follows:

$$e_i = (0, \dots, 0, 1, 0, \dots, 0),$$

where the 1 occurs in position i . Observe that the following properties hold for all $i, j = 1, \dots, n$:

$$\begin{aligned} e_i^2 &= e_i \\ e_i e_j &= 0, \quad i \neq j \\ e_1 + \dots + e_n &= 1_A. \end{aligned}$$

Also, for any element $a = (a_1, \dots, a_n) \in A$, we have

$$e_i a = (0, \dots, 0, a_i, 0, \dots, 0) = pr_i(a),$$

where pr_i is the projection of A onto A_i . As a consequence

$$\text{Ker}(pr_i) = (1_A - e_i)A.$$

Definition 21.3. Given a commutative ring A , a *direct decomposition* of A is a sequence $(\mathfrak{b}_1, \dots, \mathfrak{b}_n)$ of ideals in A such that there is an isomorphism $A \approx A/\mathfrak{b}_1 \times \dots \times A/\mathfrak{b}_n$.

The following theorem gives useful conditions characterizing direct decompositions of a ring.

Theorem 21.16. Let A be a commutative ring and let $(\mathfrak{b}_1, \dots, \mathfrak{b}_n)$ be a sequence of ideals in A . The following conditions are equivalent:

- (a) The sequence $(\mathfrak{b}_1, \dots, \mathfrak{b}_n)$ is a direct decomposition of A .
- (b) There exist some elements e_1, \dots, e_n of A such that

$$\begin{aligned} e_i^2 &= e_i \\ e_i e_j &= 0, \quad i \neq j \\ e_1 + \dots + e_n &= 1_A, \end{aligned}$$

and $\mathfrak{b}_i = (1_A - e_i)A$, for $i, j = 1, \dots, n$.

- (c) We have $\mathfrak{b}_i + \mathfrak{b}_j = A$ for all $i \neq j$, and $\mathfrak{b}_1 \cdots \mathfrak{b}_n = (0)$.
- (d) We have $\mathfrak{b}_i + \mathfrak{b}_j = A$ for all $i \neq j$, and $\mathfrak{b}_1 \cap \dots \cap \mathfrak{b}_n = (0)$.

Proof. Assume (a). Since we have an isomorphism $A \approx A/\mathfrak{b}_1 \times \cdots \times A/\mathfrak{b}_n$, we may identify A with $A/\mathfrak{b}_1 \times \cdots \times A/\mathfrak{b}_n$, and \mathfrak{b}_i with $\text{Ker}(pr_i)$. Then, e_1, \dots, e_n are the elements defined just before Definition 21.3. As noted, $\mathfrak{b}_i = \text{Ker}(pr_i) = (1_A - e_i)A$. This proves (b).

Assume (b). Since $\mathfrak{b}_i = (1_A - e_i)A$ and A is a ring with unit 1_A , we have $1_A - e_i \in \mathfrak{b}_i$ for $i = 1, \dots, n$. For all $i \neq j$, we also have $e_i(1_A - e_j) = e_i - e_i e_j = e_i$, so (because \mathfrak{b}_j is an ideal), $e_i \in \mathfrak{b}_j$, and thus, $1_A = 1_A - e_i + e_i \in \mathfrak{b}_i + \mathfrak{b}_j$, which shows that $\mathfrak{b}_i + \mathfrak{b}_j = A$ for all $i \neq j$. Furthermore, for any $x_i \in A$, with $1 \leq i \leq n$, we have

$$\begin{aligned} \prod_{i=1}^n x_i (1_A - e_i) &= \left(\prod_{i=1}^n x_i \right) \prod_{i=1}^n (1_A - e_i) \\ &= \left(\prod_{i=1}^n x_i \right) (1_A - \sum_{i=1}^n e_i) \\ &= 0, \end{aligned}$$

which proves that $\mathfrak{b}_1 \cdots \mathfrak{b}_n = (0)$. Thus, (c) holds.

The equivalence of (c) and (d) follows from the proof of Theorem 21.15.

The fact that (c) implies (a) is an immediate consequence of Theorem 21.15. \square

21.3 Noetherian Rings and Hilbert's Basis Theorem

Given a (commutative) ring A (with unit element 1), an ideal $\mathfrak{A} \subseteq A$ is said to be *finitely generated* if there exists a finite set $\{a_1, \dots, a_n\}$ of elements from \mathfrak{A} so that

$$\mathfrak{A} = (a_1, \dots, a_n) = \{\lambda_1 a_1 + \cdots + \lambda_n a_n \mid \lambda_i \in A, 1 \leq i \leq n\}.$$

If K is a field, it turns out that every polynomial ideal \mathfrak{A} in $K[X_1, \dots, X_m]$ is finitely generated. This fact due to Hilbert and known as Hilbert's basis theorem, has very important consequences. For example, in algebraic geometry, one is interested in the zero locus of a set of polynomial equations, i.e., the set, $V(\mathcal{P})$, of n -tuples $(\lambda_1, \dots, \lambda_n) \in K^n$ so that

$$P_i(\lambda_1, \dots, \lambda_n) = 0$$

for all polynomials $P_i(X_1, \dots, X_n)$ in some given family, $\mathcal{P} = (P_i)_{i \in I}$. However, it is clear that

$$V(\mathcal{P}) = V(\mathfrak{A}),$$

where \mathfrak{A} is the ideal generated by \mathcal{P} . Then, Hilbert's basis theorem says that $V(\mathfrak{A})$ is actually defined by a *finite* number of polynomials (any set of generators of \mathfrak{A}), even if \mathcal{P} is infinite.

The property that every ideal in a ring is finitely generated is equivalent to other natural properties, one of which is the so-called *ascending chain condition*, abbreviated *a.c.c.* Before proving Hilbert's basis theorem, we explore the equivalence of these conditions.

Definition 21.4. Let A be a commutative ring with unit 1. We say that A satisfies the *ascending chain condition*, for short, the *a.c.c.*, if for every ascending chain of ideals

$$\mathfrak{A}_1 \subseteq \mathfrak{A}_2 \subseteq \cdots \subseteq \mathfrak{A}_i \subseteq \cdots,$$

there is some integer $n \geq 1$ so that

$$\mathfrak{A}_i = \mathfrak{A}_n \quad \text{for all } i \geq n + 1.$$

We say that A satisfies the *maximum condition* if every nonempty collection C of ideals in A has a maximal element, i.e., there is some ideal $\mathfrak{A} \in C$ which is not contained in any other ideal in C .

Proposition 21.17. *A ring A satisfies the a.c.c if and only if it satisfies the maximum condition.*

Proof. Suppose that A does not satisfy the a.c.c. Then, there is an infinite strictly ascending sequence of ideals

$$\mathfrak{A}_1 \subset \mathfrak{A}_2 \subset \cdots \subset \mathfrak{A}_i \subset \cdots,$$

and the collection $C = \{\mathfrak{A}_i\}$ has no maximal element.

Conversely, assume that A satisfies the a.c.c. Let C be a nonempty collection of ideals. Since C is nonempty, we may pick some ideal \mathfrak{A}_1 in C . If \mathfrak{A}_1 is not maximal, then there is some ideal \mathfrak{A}_2 in C so that

$$\mathfrak{A}_1 \subset \mathfrak{A}_2.$$

Using this process, if C has no maximal element, we can define by induction an infinite strictly increasing sequence

$$\mathfrak{A}_1 \subset \mathfrak{A}_2 \subset \cdots \subset \mathfrak{A}_i \subset \cdots.$$

However, the a.c.c. implies that such a sequence cannot exist. Therefore, C has a maximal element. \square

Having shown that the a.c.c. condition is equivalent to the maximal condition, we now prove that the a.c.c. condition is equivalent to the fact that every ideal is finitely generated.

Proposition 21.18. *A ring A satisfies the a.c.c if and only if every ideal is finitely generated.*

Proof. Assume that every ideal is finitely generated. Consider an ascending sequence of ideals

$$\mathfrak{A}_1 \subseteq \mathfrak{A}_2 \subseteq \cdots \subseteq \mathfrak{A}_i \subseteq \cdots.$$

Observe that $\mathfrak{A} = \bigcup_i \mathfrak{A}_i$ is also an ideal. By hypothesis, \mathfrak{A} has a finite generating set $\{a_1, \dots, a_n\}$. By definition of \mathfrak{A} , each a_i belongs to some \mathfrak{A}_{j_i} , and since the \mathfrak{A}_i form an ascending chain, there is some m so that $a_i \in \mathfrak{A}_m$ for $i = 1, \dots, n$. But then,

$$\mathfrak{A} = \mathfrak{A}_m$$

for all $i \geq m + 1$, and the a.c.c. holds.

Conversely, assume that the a.c.c. holds. Let \mathfrak{A} be any ideal in A and consider the family C of subideals of \mathfrak{A} that are finitely generated. The family C is nonempty, since (0) is a subideal of \mathfrak{A} . By Proposition 21.17, the family C has some maximal element, say \mathfrak{B} . For any $a \in \mathfrak{A}$, the ideal $\mathfrak{B} + (a)$ (where $\mathfrak{B} + (a) = \{b + \lambda a \mid b \in \mathfrak{B}, \lambda \in A\}$) is also finitely generated (since \mathfrak{B} is finitely generated), and by maximality, we have

$$\mathfrak{B} = \mathfrak{B} + (a).$$

So, we get $a \in \mathfrak{B}$ for all $a \in \mathfrak{A}$, and thus, $\mathfrak{A} = \mathfrak{B}$, and \mathfrak{A} is finitely generated. \square

Definition 21.5. A commutative ring A (with unit 1) is called *noetherian* if it satisfies the a.c.c. condition. A *noetherian domain* is a noetherian ring that is also a domain.

By Proposition 21.17 and Proposition 21.18, a noetherian ring can also be defined as a ring that either satisfies the maximal property or such that every ideal is finitely generated. The proof of Hilbert's basis theorem will make use the following lemma:

Lemma 21.19. *Let A be a (commutative) ring. For every ideal \mathfrak{A} in $A[X]$, for every $i \geq 0$, let $L_i(\mathfrak{A})$ denote the set of elements of A consisting of 0 and of the coefficients of X^i in all the polynomials $f(X) \in \mathfrak{A}$ which are of degree i . Then, the $L_i(\mathfrak{A})$'s form an ascending chain of ideals in A . Furthermore, if \mathfrak{B} is any ideal of $A[X]$ so that $\mathfrak{A} \subseteq \mathfrak{B}$ and if $L_i(\mathfrak{A}) = L_i(\mathfrak{B})$ for all $i \geq 0$, then $\mathfrak{A} = \mathfrak{B}$.*

Proof. That $L_i(\mathfrak{A})$ is an ideal and that $L_i(\mathfrak{A}) \subseteq L_{i+1}(\mathfrak{A})$ follows from the fact that if $f(X) \in \mathfrak{A}$ and $g(X) \in \mathfrak{A}$, then $f(X) + g(X)$, $\lambda f(X)$, and $Xf(X)$ all belong to \mathfrak{A} . Now, let $g(X)$ be any polynomial in \mathfrak{B} , and assume that $g(X)$ has degree n . Since $L_n(\mathfrak{A}) = L_n(\mathfrak{B})$, there is some polynomial $f_n(X)$ in \mathfrak{A} , of degree n , so that $g(X) - f_n(X)$ is of degree at most $n - 1$. Now, since $\mathfrak{A} \subseteq \mathfrak{B}$, the polynomial $g(X) - f_n(X)$ belongs to \mathfrak{B} . Using this process, we can define by induction a sequence of polynomials $f_{n+i}(X) \in \mathfrak{A}$, so that each $f_{n+i}(X)$ is either zero or has degree $n - i$, and

$$g(X) - (f_n(X) + f_{n+1}(X) + \cdots + f_{n+i}(X))$$

is of degree at most $n - i - 1$. Note that this last polynomial must be zero when $i = n$, and thus, $g(X) \in \mathfrak{A}$. \square

We now prove Hilbert's basis theorem. The proof is substantially Hilbert's original proof. A slightly shorter proof can be given but it is not as transparent as Hilbert's proof (see the remark just after the proof of Theorem 21.20, and Zariski and Samuel [117], Chapter IV, Section 1, Theorem 1).

Theorem 21.20. (*Hilbert's basis theorem*) *If A is a noetherian ring, then $A[X]$ is also a noetherian ring.*

Proof. Let \mathfrak{A} be any ideal in $A[X]$, and denote by \mathcal{L} the set of elements of A consisting of 0 and of all the coefficients of the highest degree terms of all the polynomials in \mathfrak{A} . Observe that

$$\mathcal{L} = \bigcup_i L_i(\mathfrak{A}).$$

Thus, \mathcal{L} is an ideal in A (this can also be proved directly). Since A is noetherian, \mathcal{L} is finitely generated, and let $\{a_1, \dots, a_n\}$ be a set of generators of \mathcal{L} . Let $f_1(X), \dots, f_n(X)$ be polynomials in \mathfrak{A} having respectively a_1, \dots, a_n as highest degree term coefficients. These polynomials generate an ideal \mathfrak{B} . Let q be the maximum of the degrees of the $f_i(X)$'s. Now, pick any polynomial $g(X) \in \mathfrak{A}$ of degree $d \geq q$, and let aX^d be its term of highest degree. Since $a \in \mathcal{L}$, we have

$$a = \lambda_1 a_1 + \dots + \lambda_n a_n,$$

for some $\lambda_i \in A$. Consider the polynomial

$$g_1(X) = \sum_{i=1}^n \lambda_i f_i(X) X^{d-d_i},$$

where d_i is the degree of $f_i(X)$. Now, $g(X) - g_1(X)$ is a polynomial in \mathfrak{A} of degree at most $d - 1$. By repeating this procedure, we get a sequence of polynomials $g_i(X)$ in \mathfrak{B} , having strictly decreasing degrees, and such that the polynomial

$$g(X) - (g_1(X) + \dots + g_i(X))$$

is of degree at most $d - i$. This polynomial must be of degree at most $q - 1$ as soon as $i = d - q + 1$. Thus, we proved that every polynomial in \mathfrak{A} of degree $d \geq q$ belongs to \mathfrak{B} .

It remains to take care of the polynomials in \mathfrak{A} of degree at most $q - 1$. Since A is noetherian, each ideal $L_i(\mathfrak{A})$ is finitely generated, and let $\{a_{i1}, \dots, a_{in_i}\}$ be a set of generators for $L_i(\mathfrak{A})$ (for $i = 0, \dots, q - 1$). Let $f_{ij}(X)$ be a polynomial in \mathfrak{A} having $a_{ij}X^i$ as its highest degree term. Given any polynomial $g(X) \in \mathfrak{A}$ of degree $d \leq q - 1$, if we denote its term of highest degree by aX^d , then, as in the previous argument, we can write

$$a = \lambda_1 a_{d1} + \dots + \lambda_{n_d} a_{dn_d},$$

and we define

$$g_1(X) = \sum_{i=1}^{n_d} \lambda_i f_{di}(X) X^{d-d_i},$$

where d_i is the degree of $f_{di}(X)$. Then, $g(X) - g_1(X)$ is a polynomial in \mathfrak{A} of degree at most $d - 1$, and by repeating this procedure at most q times, we get an element of \mathfrak{A} of degree 0, and the latter is a linear combination of the f_{0i} 's. This proves that every polynomial in \mathfrak{A} of degree at most $q - 1$ is a combination of the polynomials $f_{ij}(X)$, for $0 \leq i \leq q - 1$ and $1 \leq j \leq n_i$. Therefore, \mathfrak{A} is generated by the $f_k(X)$'s and the $f_{ij}(X)$'s, a finite number of polynomials. \square

Remark: Only a small part of Lemma 21.19 was used in the above proof, namely, the fact that $L_i(\mathfrak{A})$ is an ideal. A shorter proof of Theorem 21.21 making full use of Lemma 21.19 can be given as follows:

Proof. (Second proof) Let $(\mathfrak{A}_i)_{i \geq 1}$ be an ascending sequence of ideals in $A[X]$. Consider the doubly indexed family $(L_i(\mathfrak{A}_j))$ of ideals in A . Since A is noetherian, by the maximal property, this family has a maximal element $L_p(\mathfrak{A}_q)$. Since the $L_i(\mathfrak{A}_j)$'s form an ascending sequence when either i or j is fixed, we have $L_i(\mathfrak{A}_j) = L_p(\mathfrak{A}_q)$ for all i and j with $i \geq p$ and $j \geq q$, and thus, $L_i(\mathfrak{A}_q) = L_i(\mathfrak{A}_j)$ for all i and j with $i \geq p$ and $j \geq q$. On the other hand, for any fixed i , the a.c.c. shows that there exists some integer $n(i)$ so that $L_i(\mathfrak{A}_j) = L_i(\mathfrak{A}_{n(i)})$ for all $j \geq n(i)$. Since $L_i(\mathfrak{A}_q) = L_i(\mathfrak{A}_j)$ when $i \geq p$ and $j \geq q$, we may take $n(i) = q$ if $i \geq p$. This shows that there is some n_0 so that $n(i) \leq n_0$ for all $i \geq 0$, and thus, we have $L_i(\mathfrak{A}_j) = L_i(\mathfrak{A}_{n(0)})$ for every i and for every $j \geq n(0)$. By Lemma 21.19, we get $\mathfrak{A}_j = \mathfrak{A}_{n(0)}$ for every $j \geq n(0)$, establishing the fact that $A[X]$ satisfies the a.c.c. \square

Using induction, we immediately obtain the following important result.

Corollary 21.21. *If A is a noetherian ring, then $A[X_1, \dots, X_n]$ is also a noetherian ring.*

Since a field K is obviously noetherian (since it has only two ideals, (0) and K), we also have:

Corollary 21.22. *If K is a field, then $K[X_1, \dots, X_n]$ is a noetherian ring.*

21.4 Futher Readings

The material of this Chapter is thoroughly covered in Lang [67], Artin [4], Mac Lane and Birkhoff [73], Bourbaki [14, 15], Malliavin [75], Zariski and Samuel [117], and Van Der Waerden [112].

Chapter 22

Annihilating Polynomials and the Primary Decomposition

22.1 Annihilating Polynomials and the Minimal Polynomial

In Section 5.7, we explained that if $f: E \rightarrow E$ is a linear map on a K -vector space E , then for any polynomial $p(X) = a_0X^d + a_1X^{d-1} + \cdots + a_d$ with coefficients in the field K , we can define the *linear map* $p(f): E \rightarrow E$ by

$$p(f) = a_0f^d + a_1f^{d-1} + \cdots + a_d\text{id},$$

where $f^k = f \circ \cdots \circ f$, the k -fold composition of f with itself. Note that

$$p(f)(u) = a_0f^d(u) + a_1f^{d-1}(u) + \cdots + a_du,$$

for every vector $u \in E$. Then, we showed that if E is finite-dimensional and if $\chi_f(X) = \det(XI - f)$ is the characteristic polynomial of f , by the Cayley–Hamilton Theorem, we have

$$\chi_f(f) = 0.$$

This fact suggests looking at the set of all polynomials $p(X)$ such that

$$p(f) = 0.$$

We say that the polynomial $p(X)$ *annihilates* f . It is easy to check that the set $\text{Ann}(f)$ of polynomials that annihilate f is an ideal. Furthermore, when E is finite-dimensional, the Cayley–Hamilton Theorem implies that $\text{Ann}(f)$ is not the zero ideal. Therefore, by Proposition 20.9, there is a unique monic polynomial m_f that generates $\text{Ann}(f)$. Results from Chapter 20, especially about gcd's of polynomials, will come handy.

Definition 22.1. If $f: E \rightarrow E$, is linear map on a finite-dimensional vector space E , the unique monic polynomial $m_f(X)$ that generates the ideal $\text{Ann}(f)$ of polynomials which annihilate f (the *annihilator* of f) is called the *minimal polynomial* of f .

The minimal polynomial m_f of f is the monic polynomial of smallest degree that annihilates f . Thus, the minimal polynomial divides the characteristic polynomial χ_f , and $\deg(m_f) \geq 1$. For simplicity of notation, we often write m instead of m_f .

If A is any $n \times n$ matrix, the set $\text{Ann}(A)$ of polynomials that annihilate A is the set of polynomials

$$p(X) = a_0X^d + a_1X^{d-1} + \cdots + a_{d-1}X + a_d$$

such that

$$a_0A^d + a_1A^{d-1} + \cdots + a_{d-1}A + a_dI = 0.$$

It is clear that $\text{Ann}(A)$ is a nonzero ideal and its unique monic generator is called the *minimal polynomial* of A . We check immediately that if Q is an invertible matrix, then A and $Q^{-1}AQ$ have the same minimal polynomial. Also, if A is the matrix of f with respect to some basis, then f and A have the same minimal polynomial.

The zeros (in K) of the minimal polynomial of f and the eigenvalues of f (in K) are intimately related.

Proposition 22.1. *Let $f: E \rightarrow E$ be a linear map on some finite-dimensional vector space E . Then, $\lambda \in K$ is a zero of the minimal polynomial $m_f(X)$ of f iff λ is an eigenvalue of f iff λ is a zero of $\chi_f(X)$. Therefore, the minimal and the characteristic polynomials have the same zeros (in K), except for multiplicities.*

Proof. First, assume that $m(\lambda) = 0$ (with $\lambda \in K$, and writing m instead of m_f). If so, using polynomial division, m can be factored as

$$m = (X - \lambda)q,$$

with $\deg(q) < \deg(m)$. Since m is the minimal polynomial, $q(f) \neq 0$, so there is some nonzero vector $v \in E$ such that $u = q(f)(v) \neq 0$. But then, because m is the minimal polynomial,

$$\begin{aligned} 0 &= m(f)(v) \\ &= (f - \lambda \text{id})(q(f)(v)) \\ &= (f - \lambda \text{id})(u), \end{aligned}$$

which shows that λ is an eigenvalue of f .

Conversely, assume that $\lambda \in K$ is an eigenvalue of f . This means that for some $u \neq 0$, we have $f(u) = \lambda u$. Now, it is easy to show that

$$m(f)(u) = m(\lambda)u,$$

and since m is the minimal polynomial of f , we have $m(f)(u) = 0$, so $m(\lambda)u = 0$, and since $u \neq 0$, we must have $m(\lambda) = 0$. \square

If we assume that f is diagonalizable, then its eigenvalues are all in K , and if $\lambda_1, \dots, \lambda_k$ are the distinct eigenvalues of f , then by Proposition 22.1, the minimal polynomial m of f must be a product of powers of the polynomials $(X - \lambda_i)$. Actually, we claim that

$$m = (X - \lambda_1) \cdots (X - \lambda_k).$$

For this, we just have to show that m annihilates f . However, for any eigenvector u of f , one of the linear maps $f - \lambda_i \text{id}$ sends u to 0, so

$$m(f)(u) = (f - \lambda_1 \text{id}) \circ \cdots \circ (f - \lambda_k \text{id})(u) = 0.$$

Since E is spanned by the eigenvectors of f , we conclude that

$$m(f) = 0.$$

Therefore, if a linear map is diagonalizable, then its minimal polynomial is a product of distinct factors of degree 1. It turns out that the converse is true, but this will take a little work to establish it.

22.2 Minimal Polynomials of Diagonalizable Linear Maps

In this section, we prove that if the minimal polynomial m_f of a linear map f is of the form

$$m_f = (X - \lambda_1) \cdots (X - \lambda_k)$$

for distinct scalars $\lambda_1, \dots, \lambda_k \in K$, then f is diagonalizable. This is a powerful result that has a number of implications. We need a few properties of invariant subspaces.

Given a linear map $f: E \rightarrow E$, recall that a subspace W of E is *invariant under f* if $f(u) \in W$ for all $u \in W$.

Proposition 22.2. *Let W be a subspace of E invariant under the linear map $f: E \rightarrow E$ (where E is finite-dimensional). Then, the minimal polynomial of the restriction $f|_W$ of f to W divides the minimal polynomial of f , and the characteristic polynomial of $f|_W$ divides the characteristic polynomial of f .*

Sketch of proof. The key ingredient is that we can pick a basis (e_1, \dots, e_n) of E in which (e_1, \dots, e_k) is a basis of W . Then, the matrix of f over this basis is a block matrix of the form

$$A = \begin{pmatrix} B & C \\ 0 & D \end{pmatrix},$$

where B is a $k \times k$ matrix, D is a $(n - k) \times (n - k)$ matrix, and C is a $k \times (n - k)$ matrix. Then

$$\det(XI - A) = \det(XI - B) \det(XI - D),$$

which implies the statement about the characteristic polynomials. Furthermore,

$$A^i = \begin{pmatrix} B^i & C_i \\ 0 & D^i \end{pmatrix},$$

for some $k \times (n - k)$ matrix C_i . It follows that any polynomial which annihilates A also annihilates B and D . So, the minimal polynomial of B divides the minimal polynomial of A . \square

For the next step, there are at least two ways to proceed. We can use an old-fashion argument using Lagrange interpolants, or use a slight generalization of the notion of annihilator. We pick the second method because it illustrates nicely the power of principal ideals.

What we need is the notion of conductor (also called transporter).

Definition 22.2. Let $f: E \rightarrow E$ be a linear map on a finite-dimensional vector space E , let W be an invariant subspace of f , and let u be any vector in E . The set $S_f(u, W)$ consisting of all polynomials $q \in K[X]$ such that $q(f)(u) \in W$ is called the *f-conductor of u into W* .

Observe that the minimal polynomial m_f of f always belongs to $S_f(u, W)$, so this is a nontrivial set. Also, if $W = (0)$, then $S_f(u, (0))$ is just the annihilator of f . The crucial property of $S_f(u, W)$ is that it is an ideal.

Proposition 22.3. *If W is an invariant subspace for f , then for each $u \in E$, the f -conductor $S_f(u, W)$ is an ideal in $K[X]$.*

We leave the proof as a simple exercise, using the fact that if W invariant under f , then W is invariant under every polynomial $q(f)$ in f .

Since $S_f(u, W)$ is an ideal, it is generated by a unique monic polynomial q of smallest degree, and because the minimal polynomial m_f of f is in $S_f(u, W)$, the polynomial q divides m .

Proposition 22.4. *Let $f: E \rightarrow E$ be a linear map on a finite-dimensional space E , and assume that the minimal polynomial m of f is of the form*

$$m = (X - \lambda_1)^{r_1} \cdots (X - \lambda_k)^{r_k},$$

where the eigenvalues $\lambda_1, \dots, \lambda_k$ of f belong to K . If W is a proper subspace of E which is invariant under f , then there is a vector $u \in E$ with the following properties:

- (a) $u \notin W$;
- (b) $(f - \lambda \text{id})(u) \in W$, for some eigenvalue λ of f .

Proof. Observe that (a) and (b) together assert that the f -conductor of u into W is a polynomial of the form $X - \lambda_i$. Pick any vector $v \in E$ not in W , and let g be the conductor of v into W . Since g divides m and $v \notin W$, the polynomial g is not a constant, and thus it is of the form

$$g = (X - \lambda_1)^{s_1} \cdots (X - \lambda_k)^{s_k},$$

with at least some $s_i > 0$. Choose some index j such that $s_j > 0$. Then $X - \lambda_j$ is a factor of g , so we can write

$$g = (X - \lambda_j)q.$$

By definition of g , the vector $u = q(f)(v)$ cannot be in W , since otherwise g would not be of minimal degree. However,

$$\begin{aligned} (f - \lambda_j \text{id})(u) &= (f - \lambda_j \text{id})(q(f)(v)) \\ &= g(f)(v) \end{aligned}$$

is in W , which concludes the proof. \square

We can now prove the main result of this section.

Theorem 22.5. *Let $f: E \rightarrow E$ be a linear map on a finite-dimensional space E . Then f is diagonalizable iff its minimal polynomial m is of the form*

$$m = (X - \lambda_1) \cdots (X - \lambda_k),$$

where $\lambda_1, \dots, \lambda_k$ are distinct elements of K .

Proof. We already showed in Section 22.2 that if f is diagonalizable, then its minimal polynomial is of the above form (where $\lambda_1, \dots, \lambda_k$ are the distinct eigenvalues of f).

For the converse, let W be the subspace spanned by all the eigenvectors of f . If $W \neq E$, since W is invariant under f , by Proposition 22.4, there is some vector $u \notin W$ such that for some λ_j , we have

$$(f - \lambda_j \text{id})(u) \in W.$$

Let $v = (f - \lambda_j \text{id})(u) \in W$. Since $v \in W$, we can write

$$v = w_1 + \cdots + w_k$$

where $f(w_i) = \lambda_i w_i$ (either $w_i = 0$ or w_i is an eigenvector for λ_i), and so, for every polynomial h , we have

$$h(f)(v) = h(\lambda_1)w_1 + \cdots + h(\lambda_k)w_k,$$

which shows that $h(f)(v) \in W$ for every polynomial h . We can write

$$m = (X - \lambda_j)q$$

for some polynomial q , and also

$$q - q(\lambda_j) = p(X - \lambda_j)$$

for some polynomial p . We know that $p(f)(v) \in W$, and since m is the minimal polynomial of f , we have

$$0 = m(f)(u) = (f - \lambda_j \text{id})(q(f)(u)),$$

which implies that $q(f)(u) \in W$ (either $q(f)(u) = 0$, or it is an eigenvector associated with λ_j). However,

$$q(f)(u) - q(\lambda_j)u = p(f)((f - \lambda_j \text{id})(u)) = p(f)(v),$$

and since $p(f)(v) \in W$ and $q(f)(u) \in W$, we conclude that $q(\lambda_j)u \in W$. But, $u \notin W$, which implies that $q(\lambda_j) = 0$, so λ_j is a double root of m , a contradiction. Therefore, we must have $W = E$. \square

Remark: Proposition 22.4 can be used to give a quick proof of Theorem 8.4.

Using Theorem 22.5, we can give a short proof about commuting diagonalizable linear maps. If \mathcal{F} is a family of linear maps on a vector space E , we say that \mathcal{F} is a *commuting family* iff $f \circ g = g \circ f$ for all $f, g \in \mathcal{F}$.

Proposition 22.6. *Let \mathcal{F} be a finite commuting family of diagonalizable linear maps on a vector space E . There exists a basis of E such that every linear map in \mathcal{F} is represented in that basis by a diagonal matrix.*

Proof. We proceed by induction on $n = \dim(E)$. If $n = 1$, there is nothing to prove. If $n > 1$, there are two cases. If all linear maps in \mathcal{F} are of the form λid for some $\lambda \in K$, then the proposition holds trivially. In the second case, let $f \in \mathcal{F}$ be some linear map in \mathcal{F} which is not a scalar multiple of the identity. In this case, f has at least two distinct eigenvalues $\lambda_1, \dots, \lambda_k$, and because f is diagonalizable, E is the direct sum of the corresponding eigenspaces $E_{\lambda_1}, \dots, E_{\lambda_k}$. For every index i , the eigenspace E_{λ_i} is invariant under f and under every other linear map g in \mathcal{F} , since for any $g \in \mathcal{F}$ and any $u \in E_{\lambda_i}$, because f and g commute, we have

$$f(g(u)) = g(f(u)) = g(\lambda_i u) = \lambda_i g(u)$$

so $g(u) \in E_{\lambda_i}$. Let \mathcal{F}_i be the family obtained by restricting each $f \in \mathcal{F}$ to E_{λ_i} . By proposition 22.2, the minimal polynomial of every linear map $f|_{E_{\lambda_i}}$ in \mathcal{F}_i divides the minimal polynomial m_f of f , and since f is diagonalizable, m_f is a product of distinct linear factors, so the minimal polynomial of $f|_{E_{\lambda_i}}$ is also a product of distinct linear factors. By Theorem 22.5, the linear map $f|_{E_{\lambda_i}}$ is diagonalizable. Since $k > 1$, we have $\dim(E_{\lambda_i}) < \dim(E)$ for $i = 1, \dots, k$, and by the induction hypothesis, for each i there is a basis of E_{λ_i} over which $f|_{E_{\lambda_i}}$ is represented by a diagonal matrix. Since the above argument holds for all i , by combining the bases of the E_{λ_i} , we obtain a basis of E such that the matrix of every linear map $f \in \mathcal{F}$ is represented by a diagonal matrix. \square

Remark: Proposition 22.6 also holds for infinite commuting families \mathcal{F} of diagonalizable linear maps, because E being finite dimensional, there is a finite subfamily of linearly independent linear maps in \mathcal{F} spanning \mathcal{F} .

There is also an analogous result for commuting families of linear maps represented by upper triangular matrices. To prove this, we need the following proposition.

Proposition 22.7. *Let \mathcal{F} be a nonempty finite commuting family of triangulable linear maps on a finite-dimensional vector space E . Let W be a proper subspace of E which is invariant under \mathcal{F} . Then there exists a vector $u \in E$ such that:*

1. $u \notin W$.
2. For every $f \in \mathcal{F}$, the vector $f(u)$ belongs to the subspace $W \oplus Ku$ spanned by W and u .

Proof. By renaming the elements of \mathcal{F} if necessary, we may assume that (f_1, \dots, f_r) is a basis of the subspace of $\text{End}(E)$ spanned by \mathcal{F} . We prove by induction on r that there exists some vector $u \in E$ such that

1. $u \notin W$.
2. $(f_i - \alpha_i \text{id})(u) \in W$ for $i = 1, \dots, r$, for some scalars $\alpha_i \in K$.

Consider the base case $r = 1$. Since f_1 is triangulable, its eigenvalues all belong to K since they are the diagonal entries of the triangular matrix associated with f_1 (this is the easy direction of Theorem 8.4), so the minimal polynomial of f_1 is of the form

$$m = (X - \lambda_1)^{r_1} \cdots (X - \lambda_k)^{r_k},$$

where the eigenvalues $\lambda_1, \dots, \lambda_k$ of f_1 belong to K . We conclude by applying Proposition 22.4.

Next, assume that $r \geq 2$ and that the induction hypothesis holds for f_1, \dots, f_{r-1} . Thus, there is a vector $u_{r-1} \in E$ such that

1. $u_{r-1} \notin W$.
2. $(f_i - \alpha_i \text{id})(u_{r-1}) \in W$ for $i = 1, \dots, r-1$, for some scalars $\alpha_i \in K$.

Let

$$V_{r-1} = \{w \in E \mid (f_i - \alpha_i \text{id})(w) \in W, i = 1, \dots, r-1\}.$$

Clearly, $W \subseteq V_{r-1}$ and $u_{r-1} \in V_{r-1}$. We claim that V_{r-1} is invariant under \mathcal{F} . This is because, for any $v \in V_{r-1}$ and any $f \in \mathcal{F}$, since f and f_i commute, we have

$$(f_i - \alpha_i \text{id})(f(v)) = f(f_i - \alpha_i \text{id})(v), \quad 1 \leq i \leq r-1.$$

Now, $(f_i - \alpha_i \text{id})(v) \in W$ because $v \in V_{r-1}$, and W is invariant under \mathcal{F} so $f(f_i - \alpha_i \text{id})(v) \in W$, that is, $(f_i - \alpha_i \text{id})(f(v)) \in W$.

Consider the restriction g_r of f_r to V_{r-1} . The minimal polynomial of g_r divides the minimal polynomial of f_r , and since f_r is triangulable, just as we saw for f_1 , the minimal polynomial of f_r is of the form

$$m = (X - \lambda_1)^{r_1} \cdots (X - \lambda_k)^{r_k},$$

where the eigenvalues $\lambda_1, \dots, \lambda_k$ of f_r belong to K , so the minimal polynomial of g_r is of the same form. By Proposition 22.4, there is some vector $u_r \in V_{r-1}$ such that

1. $u_r \notin W$.
2. $(g_r - \alpha_r \text{id})(u_r) \in W$ for some scalars $\alpha_r \in K$.

Now, since $u_r \in V_{r-1}$, we have $(f_i - \alpha_i \text{id})(u_r) \in W$ for $i = 1, \dots, r-1$, so $(f_i - \alpha_i \text{id})(u_r) \in W$ for $i = 1, \dots, r$ (since g_r is the restriction of f_r), which concludes the proof of the induction step. Finally, since every $f \in \mathcal{F}$ is the linear combination of (f_1, \dots, f_r) , condition (2) of the inductive claim implies condition (2) of the proposition. \square

We can now prove the following result.

Proposition 22.8. *Let \mathcal{F} be a nonempty finite commuting family of triangulable linear maps on a finite-dimensional vector space E . There exists a basis of E such that every linear map in \mathcal{F} is represented in that basis by an upper triangular matrix.*

Proof. Let $n = \dim(E)$. We construct inductively a basis (u_1, \dots, u_n) of E such that if W_i is the subspace spanned by (u_1, \dots, u_i) , then for every $f \in \mathcal{F}$,

$$f(u_i) = a_{1i}^f u_1 + \cdots + a_{ii}^f u_i,$$

for some $a_{ij}^f \in K$; that is, $f(u_i)$ belongs to the subspace W_i .

We begin by applying Proposition 22.7 to the subspace $W_0 = (0)$ to get u_1 so that for all $f \in \mathcal{F}$,

$$f(u_1) = \alpha_1^f u_1.$$

For the induction step, since W_i is invariant under \mathcal{F} , we apply Proposition 22.7 to the subspace W_i , to get $u_{i+1} \in E$ such that

1. $u_{i+1} \notin W_i$.
2. For every $f \in \mathcal{F}$, the vector $f(u_{i+1})$ belongs to the subspace spanned by W_i and u_{i+1} .

Condition (1) implies that $(u_1, \dots, u_i, u_{i+1})$ is linearly independent, and condition (2) means that for every $f \in \mathcal{F}$,

$$f(u_{i+1}) = a_{1i+1}^f u_1 + \cdots + a_{i+1,i+1}^f u_{i+1},$$

for some $a_{i+1,j}^f \in K$, establishing the induction step. After n steps, each $f \in \mathcal{F}$ is represented by an upper triangular matrix. \square

Observe that if \mathcal{F} consists of a single linear map f and if the minimal polynomial of f is of the form

$$m = (X - \lambda_1)^{r_1} \cdots (X - \lambda_k)^{r_k},$$

with all $\lambda_i \in K$, using Proposition 22.4 instead of Proposition 22.7, the proof of Proposition 22.8 yields another proof of Theorem 8.4.

22.3 The Primary Decomposition Theorem

If $f: E \rightarrow E$ is a linear map and $\lambda \in K$ is an eigenvalue of f , recall that the eigenspace E_λ associated with λ is the kernel of the linear map $\lambda \text{id} - f$. If all the eigenvalues $\lambda_1, \dots, \lambda_k$ of f are in K , it may happen that

$$E = E_{\lambda_1} \oplus \cdots \oplus E_{\lambda_k},$$

but in general there are not enough eigenvectors to span E . What if we generalize the notion of eigenvector and look for (nonzero) vectors u such that

$$(\lambda \text{id} - f)^r(u) = 0, \quad \text{for some } r \geq 1?$$

Then, it turns out that if the minimal polynomial of f is of the form

$$m = (X - \lambda_1)^{r_1} \cdots (X - \lambda_k)^{r_k},$$

then $r = r_i$ does the job for λ_i ; that is, if we let

$$W_i = \text{Ker}(\lambda_i \text{id} - f)^{r_i},$$

then

$$E = W_1 \oplus \cdots \oplus W_k.$$

This result is very nice but seems to require that the eigenvalues of f all belong to K . Actually, it is a special case of a more general result involving the factorization of the minimal polynomial m into its irreducible monic factors (See Theorem 20.16),

$$m = p_1^{r_1} \cdots p_k^{r_k},$$

where the p_i are distinct irreducible monic polynomials over K .

Theorem 22.9. (*Primary Decomposition Theorem*) Let $f: E \rightarrow E$ be a linear map on the finite-dimensional vector space E over the field K . Write the minimal polynomial m of f as

$$m = p_1^{r_1} \cdots p_k^{r_k},$$

where the p_i are distinct irreducible monic polynomials over K , and the r_i are positive integers. Let

$$W_i = \text{Ker}(p_i^{r_i}(f)), \quad i = 1, \dots, k.$$

Then

(a) $E = W_1 \oplus \cdots \oplus W_k$.

(b) Each W_i is invariant under f .

(c) The minimal polynomial of the restriction $f|_{W_i}$ of f to W_i is $p_i^{r_i}$.

Proof. The trick is to construct projections π_i using the polynomials $p_j^{r_j}$ so that the range of π_i is equal to W_i . Let

$$g_i = m/p_i^{r_i} = \prod_{j \neq i} p_j^{r_j}.$$

Note that

$$p_i^{r_i} g_i = m.$$

Since p_1, \dots, p_k are irreducible and distinct, they are relatively prime. Then, using Proposition 20.13, it is easy to show that g_1, \dots, g_k are relatively prime. Otherwise, some irreducible polynomial p would divide all of g_1, \dots, g_k , so by Proposition 20.13 it would be equal to one of the irreducible factors p_i . But, that p_i is missing from g_i , a contradiction. Therefore, by Proposition 20.14, there exist some polynomials h_1, \dots, h_k such that

$$g_1 h_1 + \cdots + g_k h_k = 1.$$

Let $q_i = g_i h_i$ and let $\pi_i = q_i(f) = g_i(f) h_i(f)$. We have

$$q_1 + \cdots + q_k = 1,$$

and since m divides $q_i q_j$ for $i \neq j$, we get

$$\begin{aligned} \pi_1 + \cdots + \pi_k &= \text{id} \\ \pi_i \pi_j &= 0, \quad i \neq j. \end{aligned}$$

(We implicitly used the fact that if p, q are two polynomials, the linear maps $p(f) \circ q(f)$ and $q(f) \circ p(f)$ are the same since $p(f)$ and $q(f)$ are polynomials in the powers of f , which commute.) Composing the first equation with π_i and using the second equation, we get

$$\pi_i^2 = \pi_i.$$

Therefore, the π_i are projections, and E is the direct sum of the images of the π_i . Indeed, every $u \in E$ can be expressed as

$$u = \pi_1(u) + \cdots + \pi_k(u).$$

Also, if

$$\pi_1(u) + \cdots + \pi_k(u) = 0,$$

then by applying π_i we get

$$0 = \pi_i^2(u) = \pi_i(u), \quad i = 1, \dots, k.$$

To finish proving (a), we need to show that

$$W_i = \text{Ker}(p_i^{r_i}(f)) = \pi_i(E).$$

If $v \in \pi_i(E)$, then $v = \pi_i(u)$ for some $u \in E$, so

$$\begin{aligned} p_i^{r_i}(f)(v) &= p_i^{r_i}(f)(\pi_i(u)) \\ &= p_i^{r_i}(f)g_i(f)h_i(f)(u) \\ &= h_i(f)p_i^{r_i}(f)g_i(f)(u) \\ &= h_i(f)m(f)(u) = 0, \end{aligned}$$

because m is the minimal polynomial of f . Therefore, $v \in W_i$.

Conversely, assume that $v \in W_i = \text{Ker}(p_i^{r_i}(f))$. If $j \neq i$, then $g_j h_j$ is divisible by $p_i^{r_i}$, so

$$g_j(f)h_j(f)(v) = \pi_j(v) = 0, \quad j \neq i.$$

Then, since $\pi_1 + \cdots + \pi_k = \text{id}$, we have $v = \pi_i v$, which shows that v is in the range of π_i . Therefore, $W_i = \text{Im}(\pi_i)$, and this finishes the proof of (a).

If $p_i^{r_i}(f)(u) = 0$, then $p_i^{r_i}(f)(f(u)) = f(p_i^{r_i}(f)(u)) = 0$, so (b) holds.

If we write $f_i = f|_{W_i}$, then $p_i^{r_i}(f_i) = 0$, because $p_i^{r_i}(f) = 0$ on W_i (its kernel). Therefore, the minimal polynomial of f_i divides $p_i^{r_i}$. Conversely, let q be any polynomial such that $q(f_i) = 0$ (on W_i). Since $m = p_i^{r_i}g_i$, the fact that $m(f)(u) = 0$ for all $u \in E$ shows that

$$p_i^{r_i}(f)(g_i(f)(u)) = 0, \quad u \in E,$$

and thus $\text{Im}(g_i(f)) \subseteq \text{Ker}(p_i^{r_i}(f)) = W_i$. Consequently, since $q(f)$ is zero on W_i ,

$$q(f)g_i(f) = 0 \quad \text{for all } u \in E.$$

But then, qg_i is divisible by the minimal polynomial $m = p_i^{r_i}g_i$ of f , and since $p_i^{r_i}$ and g_i are relatively prime, by Euclid's Proposition, $p_i^{r_i}$ must divide q . This finishes the proof that the minimal polynomial of f_i is $p_i^{r_i}$, which is (c). \square

If all the eigenvalues of f belong to the field K , we obtain the following result.

Theorem 22.10. (Primary Decomposition Theorem, Version 2) *Let $f: E \rightarrow E$ be a linear map on the finite-dimensional vector space E over the field K . If all the eigenvalues $\lambda_1, \dots, \lambda_k$ of f belong to K , write*

$$m = (X - \lambda_1)^{r_1} \cdots (X - \lambda_k)^{r_k}$$

for the minimal polynomial of f ,

$$\chi_f = (X - \lambda_1)^{n_1} \cdots (X - \lambda_k)^{n_k}$$

for the characteristic polynomial of f , with $1 \leq r_i \leq n_i$, and let

$$W_i = \text{Ker}(\lambda_i \text{id} - f)^{r_i}, \quad i = 1, \dots, k.$$

Then

- (a) $E = W_1 \oplus \cdots \oplus W_k$.
- (b) Each W_i is invariant under f .
- (c) $\dim(W_i) = n_i$.
- (d) The minimal polynomial of the restriction $f|_{W_i}$ of f to W_i is $(X - \lambda_i)^{r_i}$.

Proof. Parts (a), (b) and (d) have already been proved in Theorem 22.10, so it remains to prove (c). Since W_i is invariant under f , let f_i be the restriction of f to W_i . The characteristic polynomial χ_{f_i} of f_i divides $\chi(f)$, and since $\chi(f)$ has all its roots in K , so does $\chi_i(f)$. By Theorem 8.4, there is a basis of W_i in which f_i is represented by an upper triangular matrix, and since $(\lambda_i \text{id} - f)^{r_i} = 0$, the diagonal entries of this matrix are equal to λ_i . Consequently,

$$\chi_{f_i} = (X - \lambda_i)^{\dim(W_i)},$$

and since χ_{f_i} divides $\chi(f)$, we conclude that

$$\dim(W_i) \leq n_i, \quad i = 1, \dots, k.$$

Because E is the direct sum of the W_i , we have $\dim(W_1) + \cdots + \dim(W_k) = n$, and since $n_1 + \cdots + n_k = n$, we must have

$$\dim(W_i) = n_i, \quad i = 1, \dots, k,$$

proving (c). □

Definition 22.3. If $\lambda \in K$ is an eigenvalue of f , we define a *generalized eigenvector* of f as a nonzero vector $u \in E$ such that

$$(\lambda \text{id} - f)^r(u) = 0, \quad \text{for some } r \geq 1.$$

The *index* of λ is defined as the smallest $r \geq 1$ such that

$$\text{Ker}(\lambda \text{id} - f)^r = \text{Ker}(\lambda \text{id} - f)^{r+1}.$$

It is clear that $\text{Ker}(\lambda \text{id} - f)^i \subseteq \text{Ker}(\lambda \text{id} - f)^{i+1}$ for all $i \geq 1$. By Theorem 22.10(d), if $\lambda = \lambda_i$, the index of λ_i is equal to r_i .

Another important consequence of Theorem 22.10 is that f can be written as the sum of a diagonalizable and a nilpotent linear map (which commute). If we write

$$D = \lambda_1 \pi_1 + \cdots + \lambda_k \pi_k,$$

where π_i is the projection from E onto the subspace W_i defined in the proof of Theorem 22.9, since

$$\pi_1 + \cdots + \pi_k = \text{id},$$

we have

$$f = f\pi_1 + \cdots + f\pi_k,$$

and so we get

$$f - D = (f - \lambda_1 \text{id})\pi_1 + \cdots + (f - \lambda_k \text{id})\pi_k.$$

Since the π_i are polynomials in f , they commute with f , and if we write $N = f - D$, using the properties of the π_i , we get

$$N^r = (f - \lambda_1 \text{id})^r \pi_1 + \cdots + (f - \lambda_k \text{id})^r \pi_k.$$

Therefore, if $r = \max\{r_i\}$, we have $(f - \lambda_k \text{id})^r = 0$ for $i = 1, \dots, k$, which implies that

$$N^r = 0.$$

A linear map $g: E \rightarrow E$ is said to be *nilpotent* if there is some positive integer r such that $g^r = 0$.

Since N is a polynomial in f , it commutes with f , and thus with D . From

$$D = \lambda_1 \pi_1 + \cdots + \lambda_k \pi_k,$$

and

$$\pi_1 + \cdots + \pi_k = \text{id},$$

we see that

$$\begin{aligned} D - \lambda_i \text{id} &= \lambda_1 \pi_1 + \cdots + \lambda_k \pi_k - \lambda_i (\pi_1 + \cdots + \pi_k) \\ &= (\lambda_1 - \lambda_i) \pi_1 + \cdots + (\lambda_{i-1} - \lambda_i) \pi_{i-1} + (\lambda_{i+1} - \lambda_i) \pi_{i+1} + \cdots + (\lambda_k - \lambda_i) \pi_k. \end{aligned}$$

Since the projections π_j with $j \neq i$ vanish on W_i , the above equation implies that $D - \lambda_i \text{id}$ vanishes on W_i and that $(D - \lambda_j \text{id})(W_i) \subseteq W_i$, and thus that the minimal polynomial of D is

$$(X - \lambda_1) \cdots (X - \lambda_k).$$

Since the λ_i are distinct, by Theorem 22.5, the linear map D is diagonalizable, so we have shown that when all the eigenvalues of f belong to K , there exist a diagonalizable linear map D and a nilpotent linear map N , such that

$$\begin{aligned} f &= D + N \\ DN &= ND, \end{aligned}$$

and N and D are polynomials in f .

A decomposition of f as above is called a *Jordan decomposition*. In fact, we can prove more: The maps D and N are uniquely determined by f .

Theorem 22.11. (*Jordan Decomposition*) Let $f: E \rightarrow E$ be a linear map on the finite-dimensional vector space E over the field K . If all the eigenvalues $\lambda_1, \dots, \lambda_k$ of f belong to K , then there exist a diagonalizable linear map D and a nilpotent linear map N such that

$$\begin{aligned} f &= D + N \\ DN &= ND. \end{aligned}$$

Furthermore, D and N are uniquely determined by the above equations and they are polynomials in f .

Proof. We already proved the existence part. Suppose we also have $f = D' + N'$, with $D'N' = N'D'$, where D' is diagonalizable, N' is nilpotent, and both are polynomials in f . We need to prove that $D = D'$ and $N = N'$.

Since D' and N' commute with one another and $f = D' + N'$, we see that D' and N' commute with f . Then, D' and N' commute with any polynomial in f ; hence they commute with D and N . From

$$D + N = D' + N',$$

we get

$$D - D' = N' - N,$$

and D, D', N, N' commute with one another. Since D and D' are both diagonalizable and commute, by Proposition 22.6, they are simultaneously diagonalizable, so $D - D'$ is diagonalizable. Since N and N' commute, by the binomial formula, for any $r \geq 1$,

$$(N' - N)^r = \sum_{j=0}^r (-1)^j \binom{r}{j} (N')^{r-j} N^j.$$

Since both N and N' are nilpotent, we have $N^{r_1} = 0$ and $(N')^{r_2} = 0$, for some $r_1, r_2 > 0$, so for $r \geq r_1 + r_2$, the right-hand side of the above expression is zero, which shows that $N' - N$ is nilpotent. (In fact, it is easy that $r_1 = r_2 = n$ works). It follows that $D - D' = N' - N$ is both diagonalizable and nilpotent. Clearly, the minimal polynomial of a nilpotent linear map is of the form X^r for some $r > 0$ (and $r \leq \dim(E)$). But $D - D'$ is diagonalizable, so its minimal polynomial has simple roots, which means that $r = 1$. Therefore, the minimal polynomial of $D - D'$ is X , which says that $D - D' = 0$, and then $N = N'$. \square

If K is an algebraically closed field, then Theorem 22.11 holds. This is the case when $K = \mathbb{C}$. This theorem reduces the study of linear maps (from E to itself) to the study of nilpotent operators. There is a special normal form for such operators which is discussed in the next section.

22.4 Nilpotent Linear Maps and Jordan Form

This section is devoted to a normal form for nilpotent maps. We follow Godement's exposition [47]. Let $f: E \rightarrow E$ be a nilpotent linear map on a finite-dimensional vector space over a field K , and assume that f is not the zero map. Then, there is a smallest positive integer $r \geq 1$ such $f^r \neq 0$ and $f^{r+1} = 0$. Clearly, the polynomial X^{r+1} annihilates f , and it is the minimal polynomial of f since $f^r \neq 0$. It follows that $r + 1 \leq n = \dim(E)$. Let us define the subspaces N_i by

$$N_i = \text{Ker}(f^i), \quad i \geq 0.$$

Note that $N_0 = (0)$, $N_1 = \text{Ker}(f)$, and $N_{r+1} = E$. Also, it is obvious that

$$N_i \subseteq N_{i+1}, \quad i \geq 0.$$

Proposition 22.12. *Given a nilpotent linear map f with $f^r \neq 0$ and $f^{r+1} = 0$ as above, the inclusions in the following sequence are strict:*

$$(0) = N_0 \subset N_1 \subset \cdots \subset N_r \subset N_{r+1} = E.$$

Proof. We proceed by contradiction. Assume that $N_i = N_{i+1}$ for some i with $0 \leq i \leq r$. Since $f^{r+1} = 0$, for every $u \in E$, we have

$$0 = f^{r+1}(u) = f^{i+1}(f^{r-i}(u)),$$

which shows that $f^{r-i}(u) \in N_{i+1}$. Since $N_i = N_{i+1}$, we get $f^{r-i}(u) \in N_i$, and thus $f^r(u) = 0$. Since this holds for all $u \in E$, we see that $f^r = 0$, a contradiction. \square

Proposition 22.13. *Given a nilpotent linear map f with $f^r \neq 0$ and $f^{r+1} = 0$, for any integer i with $1 \leq i \leq r$, for any subspace U of E , if $U \cap N_i = (0)$, then $f(U) \cap N_{i-1} = (0)$, and the restriction of f to U is an isomorphism onto $f(U)$.*

Proof. Pick $v \in f(U) \cap N_{i-1}$. We have $v = f(u)$ for some $u \in U$ and $f^{i-1}(v) = 0$, which means that $f^i(u) = 0$. Then, $u \in U \cap N_i$, so $u = 0$ since $U \cap N_i = (0)$, and $v = f(u) = 0$. Therefore, $f(U) \cap N_{i-1} = (0)$. The restriction of f to U is obviously surjective on $f(U)$. Suppose that $f(u) = 0$ for some $u \in U$. Then $u \in U \cap N_1 \subseteq U \cap N_i = (0)$ (since $i \geq 1$), so $u = 0$, which proves that f is also injective on U . \square

Proposition 22.14. *Given a nilpotent linear map f with $f^r \neq 0$ and $f^{r+1} = 0$, there exists a sequence of subspace U_1, \dots, U_{r+1} of E with the following properties:*

(1) $N_i = N_{i-1} \oplus U_i$, for $i = 1, \dots, r+1$.

(2) We have $f(U_i) \subseteq U_{i-1}$, and the restriction of f to U_i is an injection, for $i = 2, \dots, r+1$.

Proof. We proceed inductively, by defining the sequence U_{r+1}, U_r, \dots, U_1 . We pick U_{r+1} to be any supplement of N_r in $N_{r+1} = E$, so that

$$E = N_{r+1} = N_r \oplus U_{r+1}.$$

Since $f^{r+1} = 0$ and $N_r = \text{Ker}(f^r)$, we have $f(U_{r+1}) \subseteq N_r$, and by Proposition 22.13, as $U_{r+1} \cap N_r = (0)$, we have $f(U_{r+1}) \cap N_{r-1} = (0)$. As a consequence, we can pick a supplement U_r of N_{r-1} in N_r so that $f(U_{r+1}) \subseteq U_r$. We have

$$N_r = N_{r-1} \oplus U_r \quad \text{and} \quad f(U_{r+1}) \subseteq U_r.$$

By Proposition 22.13, f is an injection from U_{r+1} to U_r . Assume inductively that U_{r+1}, \dots, U_i have been defined for $i \geq 2$ and that they satisfy (1) and (2). Since

$$N_i = N_{i-1} \oplus U_i,$$

we have $U_i \subseteq N_i$, so $f^{i-1}(f(U_i)) = f^i(U_i) = (0)$, which implies that $f(U_i) \subseteq N_{i-1}$. Also, since $U_i \cap N_{i-1} = (0)$, by Proposition 22.13, we have $f(U_i) \cap N_{i-2} = (0)$. It follows that there is a supplement U_{i-1} of N_{i-2} in N_{i-1} that contains $f(U_i)$. We have

$$N_{i-1} = N_{i-2} \oplus U_{i-1} \quad \text{and} \quad f(U_i) \subseteq U_{i-1}.$$

The fact that f is an injection from U_i into U_{i-1} follows from Proposition 22.13. Therefore, the induction step is proved. The construction stops when $i = 1$. \square

Because $N_0 = (0)$ and $N_{r+1} = E$, we see that E is the direct sum of the U_i :

$$E = U_1 \oplus \cdots \oplus U_{r+1},$$

with $f(U_i) \subseteq U_{i-1}$, and f an injection from U_i to U_{i-1} , for $i = r+1, \dots, 2$. By a clever choice of bases in the U_i , we obtain the following nice theorem.

Theorem 22.15. *For any nilpotent linear map $f: E \rightarrow E$ on a finite-dimensional vector space E of dimension n over a field K , there is a basis of E such that the matrix N of f is of the form*

$$N = \begin{pmatrix} 0 & \nu_1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \nu_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 0 & \nu_n \\ 0 & 0 & 0 & \cdots & 0 & 0 \end{pmatrix},$$

where $\nu_i = 1$ or $\nu_i = 0$.

Proof. First, apply Proposition 22.14 to obtain a direct sum $E = \bigoplus_{i=1}^{r+1} U_i$. Then, we define a basis of E inductively as follows. First, we choose a basis

$$e_1^{r+1}, \dots, e_{n_{r+1}}^{r+1}$$

of U_{r+1} . Next, for $i = r, \dots, 2$, given the basis

$$e_1^i, \dots, e_{n_i}^i$$

of U_i , since f is injective on U_i and $f(U_i) \subseteq U_{i-1}$, the vectors $f(e_1^i), \dots, f(e_{n_i}^i)$ are linearly independent, so we define a basis of U_{i-1} by completing $f(e_1^i), \dots, f(e_{n_i}^i)$ to a basis in U_{i-1} :

$$e_1^{i-1}, \dots, e_{n_i}^{i-1}, e_{n_i+1}^{i-1}, \dots, e_{n_{i-1}}^{i-1}$$

with

$$e_j^{i-1} = f(e_j^i), \quad j = 1, \dots, n_i.$$

Since $U_1 = N_1 = \text{Ker}(f)$, we have

$$f(e_j^1) = 0, \quad j = 1, \dots, n_1.$$

These basis vectors can be arranged as the rows of the following matrix:

$$\begin{pmatrix} e_1^{r+1} & \cdots & e_{n_{r+1}}^{r+1} & & & & & & & & \\ \vdots & & \vdots & & & & & & & & \\ e_1^r & \cdots & e_{n_r}^r & e_{n_r+1}^r & \cdots & e_{n_{r-1}}^r & & & & & \\ \vdots & & \vdots & \vdots & & \vdots & & & & & \\ e_1^{r-1} & \cdots & e_{n_{r-1}}^{r-1} & e_{n_{r-1}+1}^{r-1} & \cdots & e_{n_{r-2}}^{r-1} & e_{n_{r-2}+1}^{r-1} & \cdots & e_{n_{r-3}}^{r-1} & & \\ \vdots & & \vdots & \vdots & & \vdots & \vdots & & \vdots & & \\ \vdots & & \vdots & \vdots & & \vdots & \vdots & & \vdots & & \\ e_1^1 & \cdots & e_{n_1}^1 & e_{n_1+1}^1 & \cdots & e_{n_2}^1 & e_{n_2+1}^1 & \cdots & e_{n_3}^1 & \cdots & \cdots & e_{n_1}^1 \end{pmatrix}$$

Finally, we define the basis (e_1, \dots, e_n) by listing each column of the above matrix from the bottom-up, starting with column one, then column two, *etc.* This means that we list the vectors e_j^i in the following order:

For $j = 1, \dots, n_{r+1}$, list e_j^1, \dots, e_j^{r+1} ;

In general, for $i = r, \dots, 1$,

for $j = n_{i+1} + 1, \dots, n_i$, list e_j^1, \dots, e_j^i .

Then, because $f(e_j^1) = 0$ and $e_j^{i-1} = f(e_j^i)$ for $i \geq 2$, either

$$f(e_i) = 0 \quad \text{or} \quad f(e_i) = e_{i-1},$$

which proves the theorem. □

As an application of Theorem 22.15, we obtain the *Jordan form* of a linear map.

Definition 22.4. A *Jordan block* is an $r \times r$ matrix $J_r(\lambda)$, of the form

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

where $\lambda \in K$, with $J_1(\lambda) = (\lambda)$ if $r = 1$. A *Jordan matrix*, J , is an $n \times n$ block diagonal matrix of the form

$$J = \begin{pmatrix} J_{r_1}(\lambda_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & J_{r_m}(\lambda_m) \end{pmatrix},$$

where each $J_{r_k}(\lambda_k)$ is a Jordan block associated with some $\lambda_k \in K$, and with $r_1 + \cdots + r_m = n$.

To simplify notation, we often write $J(\lambda)$ for $J_r(\lambda)$. Here is an example of a Jordan matrix with four blocks:

$$J = \begin{pmatrix} \lambda & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \lambda & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \lambda & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \lambda & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \lambda & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \lambda & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \mu & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \mu \end{pmatrix}.$$

Theorem 22.16. (*Jordan form*) Let E be a vector space of dimension n over a field K and let $f: E \rightarrow E$ be a linear map. The following properties are equivalent:

- (1) The eigenvalues of f all belong to K (i.e. the roots of the characteristic polynomial χ_f all belong to K).
- (2) There is a basis of E in which the matrix of f is a Jordan matrix.

Proof. Assume (1). First we apply Theorem 22.10, and we get a direct sum $E = \bigoplus_{j=1}^k W_k$, such that the restriction of $g_i = f - \lambda_j \text{id}$ to W_i is nilpotent. By Theorem 22.15, there is a basis of W_i such that the matrix of the restriction of g_i is of the form

$$G_i = \begin{pmatrix} 0 & \nu_1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \nu_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 0 & \nu_{n_i} \\ 0 & 0 & 0 & \cdots & 0 & 0 \end{pmatrix},$$

where $\nu_i = 1$ or $\nu_i = 0$. Furthermore, over any basis, $\lambda_i \text{id}$ is represented by the diagonal matrix D_i with λ_i on the diagonal. Then, it is clear that we can split $D_i + G_i$ into Jordan blocks by forming a Jordan block for every uninterrupted chain of 1s. By Putting the bases of the W_i together, we obtain a matrix in Jordan form for f .

Now, assume (2). If f can be represented by a Jordan matrix, it is obvious that the diagonal entries are the eigenvalues of f , so they all belong to K . \square

Observe that Theorem 22.16 applies if $K = \mathbb{C}$. It turns out that there are uniqueness properties of the Jordan blocks. There are also other fundamental normal forms for linear maps, such as the rational canonical form, but to prove these results, it is better to develop more powerful machinery about finitely generated modules over a PID. To accomplish this most effectively, we need some basic knowledge about tensor products.

Chapter 23

Tensor Algebras, Symmetric Algebras and Exterior Algebras

23.1 Tensors Products

We begin by defining tensor products of vector spaces over a field and then we investigate some basic properties of these tensors, in particular the existence of bases and duality. After this, we investigate special kinds of tensors, namely symmetric tensors and skew-symmetric tensors. Tensor products of modules over a commutative ring with identity will be discussed very briefly. They show up naturally when we consider the space of sections of a tensor product of vector bundles.

Given a linear map $f: E \rightarrow F$, we know that if we have a basis $(u_i)_{i \in I}$ for E , then f is completely determined by its values $f(u_i)$ on the basis vectors. For a multilinear map $f: E^n \rightarrow F$, we don't know if there is such a nice property but it would certainly be very useful.

In many respects, tensor products allow us to define multilinear maps in terms of their action on a suitable basis. The crucial idea is to *linearize*, that is, to create a new vector space $E^{\otimes n}$ such that the multilinear map $f: E^n \rightarrow F$ is turned into a *linear map* $f_{\otimes}: E^{\otimes n} \rightarrow F$ which is equivalent to f in a strong sense. If in addition, f is symmetric, then we can define a symmetric tensor power $\text{Sym}^n(E)$, and every symmetric multilinear map $f: E^n \rightarrow F$ is turned into a *linear map* $f_{\odot}: \text{Sym}^n(E) \rightarrow F$ which is equivalent to f in a strong sense. Similarly, if f is alternating, then we can define a skew-symmetric tensor power $\bigwedge^n(E)$, and every alternating multilinear map is turned into a *linear map* $f_{\wedge}: \bigwedge^n(E) \rightarrow F$ which is equivalent to f in a strong sense.

Tensor products can be defined in various ways, some more abstract than others. We tried to stay down to earth, without excess!

Let K be a given field, and let E_1, \dots, E_n be $n \geq 2$ given vector spaces. For any vector space F , recall that a map $f: E_1 \times \dots \times E_n \rightarrow F$ is *multilinear* iff it is linear in each of its

argument; that is,

$$\begin{aligned} f(u_1, \dots, u_{i_1}, v + w, u_{i+1}, \dots, u_n) &= f(u_1, \dots, u_{i_1}, v, u_{i+1}, \dots, u_n) \\ &\quad + f(u_1, \dots, u_{i_1}, w, u_{i+1}, \dots, u_n) \\ f(u_1, \dots, u_{i_1}, \lambda v, u_{i+1}, \dots, u_n) &= \lambda f(u_1, \dots, u_{i_1}, v, u_{i+1}, \dots, u_n), \end{aligned}$$

for all $u_j \in E_j$ ($j \neq i$), all $v, w \in E_i$ and all $\lambda \in K$, for $i = 1 \dots, n$.

The set of multilinear maps as above forms a vector space denoted $L(E_1, \dots, E_n; F)$ or $\text{Hom}(E_1, \dots, E_n; F)$. When $n = 1$, we have the vector space of linear maps $L(E, F)$ (also denoted $\text{Hom}(E, F)$). (To be very precise, we write $\text{Hom}_K(E_1, \dots, E_n; F)$ and $\text{Hom}_K(E, F)$.) As usual, the *dual space* E^* of E is defined by $E^* = \text{Hom}(E, K)$.

Before proceeding any further, we recall a basic fact about pairings. We will use this fact to deal with dual spaces of tensors.

Definition 23.1. Given two vector spaces E and F , a map $\langle -, - \rangle: E \times F \rightarrow K$ is a *nondegenerate pairing* iff it is bilinear and iff $\langle u, v \rangle = 0$ for all $v \in F$ implies $u = 0$, and $\langle u, v \rangle = 0$ for all $u \in E$ implies $v = 0$. A nondegenerate pairing induces two linear maps $\varphi: E \rightarrow F^*$ and $\psi: F \rightarrow E^*$ defined by

$$\begin{aligned} \varphi(u)(y) &= \langle u, y \rangle \\ \psi(v)(x) &= \langle x, v \rangle, \end{aligned}$$

for all $u, x \in E$ and all $v, y \in F$.

Proposition 23.1. *For every nondegenerate pairing $\langle -, - \rangle: E \times F \rightarrow K$, the induced maps $\varphi: E \rightarrow F^*$ and $\psi: F \rightarrow E^*$ are linear and injective. Furthermore, if E and F are finite dimensional, then $\varphi: E \rightarrow F^*$ and $\psi: F \rightarrow E^*$ are bijective.*

Proof. The maps $\varphi: E \rightarrow F^*$ and $\psi: F \rightarrow E^*$ are linear because $u, v \mapsto \langle u, v \rangle$ is bilinear. Assume that $\varphi(u) = 0$. This means that $\varphi(u)(y) = \langle u, y \rangle = 0$ for all $y \in F$, and as our pairing is nondegenerate, we must have $u = 0$. Similarly, ψ is injective. If E and F are finite dimensional, then $\dim(E) = \dim(E^*)$ and $\dim(F) = \dim(F^*)$. However, the injectivity of φ and ψ implies that $\dim(E) \leq \dim(F^*)$ and $\dim(F) \leq \dim(E^*)$. Consequently $\dim(E) \leq \dim(F)$ and $\dim(F) \leq \dim(E)$, so $\dim(E) = \dim(F)$. Therefore, $\dim(E) = \dim(F^*)$ and φ is bijective (and similarly $\dim(F) = \dim(E^*)$ and ψ is bijective). \square

Proposition 23.1 shows that when E and F are finite dimensional, a nondegenerate pairing induces *canonical isomorphisms* $\varphi: E \rightarrow F^*$ and $\psi: F \rightarrow E^*$; that is, isomorphisms that do not depend on the choice of bases. An important special case is the case where $E = F$ and we have an inner product (a symmetric, positive definite bilinear form) on E .

Remark: When we use the term “canonical isomorphism,” we mean that such an isomorphism is defined independently of any choice of bases. For example, if E is a finite dimensional vector space and (e_1, \dots, e_n) is any basis of E , we have the dual basis (e_1^*, \dots, e_n^*) of

E^* (where, $e_i^*(e_j) = \delta_{ij}$), and thus the map $e_i \mapsto e_i^*$ is an isomorphism between E and E^* . This isomorphism is *not* canonical.

On the other hand, if $\langle -, - \rangle$ is an inner product on E , then Proposition 23.1 shows that the nondegenerate pairing $\langle -, - \rangle$ on $E \times E$ induces a canonical isomorphism between E and E^* . This isomorphism is often denoted $\flat: E \rightarrow E^*$, and we usually write u^\flat for $\flat(u)$, with $u \in E$.

Given any basis, (e_1, \dots, e_n) , of E (not necessarily orthonormal), if we let $g_{ij} = \langle e_i, e_j \rangle$, then for every $u = \sum_{j=1}^n u_j e_j$, since $u^\flat(v) = \langle u, v \rangle$ for all $v \in V$, we have

$$u^\flat(e_i) = \langle u, e_i \rangle = \left\langle \sum_{j=1}^n u_j e_j, e_i \right\rangle = \sum_{j=1}^n u_j \langle e_j, e_i \rangle = \sum_{j=1}^n g_{ij} u_j,$$

so we get

$$u^\flat = \sum_{i=1}^n \omega_i e_i^*, \quad \text{with} \quad \omega_i = \sum_{j=1}^n g_{ij} u_j.$$

If we use the convention that coordinates of vectors are written using superscripts ($u = \sum_{i=1}^n u^i e_i$) and coordinates of one-forms (covectors) are written using subscripts ($\omega = \sum_{i=1}^n \omega_i e_i^*$), then the map \flat has the effect of lowering (flattening!) indices. The inverse of \flat is denoted $\sharp: E^* \rightarrow E$. If we write $\omega \in E^*$ as $\omega = \sum_{i=1}^n \omega_i e_i^*$ and $\omega^\sharp \in E$ as $\omega^\sharp = \sum_{j=1}^n (\omega^\sharp)^j e_j$, since

$$\omega_i = \omega(e_i) = \langle \omega^\sharp, e_i \rangle = \sum_{j=1}^n (\omega^\sharp)^j g_{ij}, \quad 1 \leq i \leq n,$$

we get

$$(\omega^\sharp)^i = \sum_{j=1}^n g^{ij} \omega_j,$$

where (g^{ij}) is the inverse of the matrix (g_{ij}) .

The inner product $\langle -, - \rangle$ on E induces an inner product on E^* also denoted $\langle -, - \rangle$, and given by

$$\langle \omega_1, \omega_2 \rangle = \langle \omega_1^\sharp, \omega_2^\sharp \rangle, \quad \text{for all } \omega_1, \omega_2 \in E^*.$$

Then, it is obvious that

$$\langle u, v \rangle = \langle u^\flat, v^\flat \rangle, \quad \text{for all } u, v \in E.$$

If (e_1, \dots, e_n) is a basis of E and $g_{ij} = \langle e_i, e_j \rangle$, as

$$(e_i^*)^\sharp = \sum_{k=1}^n g^{ik} e_k,$$

an easy computation shows that

$$\langle e_i^*, e_j^* \rangle = \langle (e_i^*)^\sharp, (e_j^*)^\sharp \rangle = g^{ij};$$

that is, in the basis (e_1^*, \dots, e_n^*) , the inner product on E^* is represented by the matrix (g^{ij}) , the inverse of the matrix (g_{ij}) .

The inner product on a finite vector space also yields a natural isomorphism between the space $\text{Hom}(E, E; K)$ of bilinear forms on E , and the space $\text{Hom}(E, E)$ of linear maps from E to itself. Using this isomorphism, we can define the trace of a bilinear form in an intrinsic manner. This technique is used in differential geometry, for example, to define the divergence of a differential one-form.

Proposition 23.2. *If $\langle -, - \rangle$ is an inner product on a finite vector space E (over a field, K), then for every bilinear form $f: E \times E \rightarrow K$, there is a unique linear map $f^\sharp: E \rightarrow E$ such that*

$$f(u, v) = \langle f^\sharp(u), v \rangle, \quad \text{for all } u, v \in E.$$

The map $f \mapsto f^\sharp$ is a linear isomorphism between $\text{Hom}(E, E; K)$ and $\text{Hom}(E, E)$.

Proof. For every $g \in \text{Hom}(E, E)$, the map given by

$$f(u, v) = \langle g(u), v \rangle, \quad u, v \in E,$$

is clearly bilinear. It is also clear that the above defines a linear map from $\text{Hom}(E, E)$ to $\text{Hom}(E, E; K)$. This map is injective, because if $f(u, v) = 0$ for all $u, v \in E$, as $\langle -, - \rangle$ is an inner product, we get $g(u) = 0$ for all $u \in E$. Furthermore, both spaces $\text{Hom}(E, E)$ and $\text{Hom}(E, E; K)$ have the same dimension, so our linear map is an isomorphism. \square

If (e_1, \dots, e_n) is an orthonormal basis of E , then we check immediately that the trace of a linear map g (which is independent of the choice of a basis) is given by

$$\text{tr}(g) = \sum_{i=1}^n \langle g(e_i), e_i \rangle,$$

where $n = \dim(E)$. We define the *trace of the bilinear form f* by

$$\text{tr}(f) = \text{tr}(f^\sharp).$$

From Proposition 23.2, $\text{tr}(f)$ is given by

$$\text{tr}(f) = \sum_{i=1}^n f(e_i, e_i),$$

for any orthonormal basis (e_1, \dots, e_n) of E . We can also check directly that the above expression is independent of the choice of an orthonormal basis.

We will also need the following Proposition to show that various families are linearly independent.

Proposition 23.3. *Let E and F be two nontrivial vector spaces and let $(u_i)_{i \in I}$ be any family of vectors $u_i \in E$. The family $(u_i)_{i \in I}$ is linearly independent iff for every family $(v_i)_{i \in I}$ of vectors $v_i \in F$, there is some linear map $f: E \rightarrow F$ so that $f(u_i) = v_i$ for all $i \in I$.*

Proof. Left as an exercise. \square

First, we define tensor products, and then we prove their existence and uniqueness up to isomorphism.

Definition 23.2. A *tensor product* of $n \geq 2$ vector spaces E_1, \dots, E_n is a vector space T together with a multilinear map $\varphi: E_1 \times \dots \times E_n \rightarrow T$, such that for every vector space F and for every multilinear map $f: E_1 \times \dots \times E_n \rightarrow F$, there is a unique linear map $f_\otimes: T \rightarrow F$ with

$$f(u_1, \dots, u_n) = f_\otimes(\varphi(u_1, \dots, u_n)),$$

for all $u_1 \in E_1, \dots, u_n \in E_n$, or for short

$$f = f_\otimes \circ \varphi.$$

Equivalently, there is a unique linear map f_\otimes such that the following diagram commutes:

$$\begin{array}{ccc} E_1 \times \dots \times E_n & \xrightarrow{\varphi} & T \\ & \searrow f & \downarrow f_\otimes \\ & & F \end{array}$$

First, we show that any two tensor products (T_1, φ_1) and (T_2, φ_2) for E_1, \dots, E_n , are isomorphic.

Proposition 23.4. *Given any two tensor products (T_1, φ_1) and (T_2, φ_2) for E_1, \dots, E_n , there is an isomorphism $h: T_1 \rightarrow T_2$ such that*

$$\varphi_2 = h \circ \varphi_1.$$

Proof. Focusing on (T_1, φ_1) , we have a multilinear map $\varphi_2: E_1 \times \dots \times E_n \rightarrow T_2$, and thus there is a unique linear map $(\varphi_2)_\otimes: T_1 \rightarrow T_2$ with

$$\varphi_2 = (\varphi_2)_\otimes \circ \varphi_1.$$

Similarly, focusing now on (T_2, φ_2) , we have a multilinear map $\varphi_1: E_1 \times \dots \times E_n \rightarrow T_1$, and thus there is a unique linear map $(\varphi_1)_\otimes: T_2 \rightarrow T_1$ with

$$\varphi_1 = (\varphi_1)_\otimes \circ \varphi_2.$$

But then, we get

$$\varphi_1 = (\varphi_1)_\otimes \circ (\varphi_2)_\otimes \circ \varphi_1,$$

and

$$\varphi_2 = (\varphi_2)_\otimes \circ (\varphi_1)_\otimes \circ \varphi_2.$$

On the other hand, focusing on (T_1, φ_1) , we have a multilinear map $\varphi_1: E_1 \times \cdots \times E_n \rightarrow T_1$, but the unique linear map $h: T_1 \rightarrow T_1$ with

$$\varphi_1 = h \circ \varphi_1$$

is $h = \text{id}$, and since $(\varphi_1)_\otimes \circ (\varphi_2)_\otimes$ is linear, as a composition of linear maps, we must have

$$(\varphi_1)_\otimes \circ (\varphi_2)_\otimes = \text{id}.$$

Similarly, we must have

$$(\varphi_2)_\otimes \circ (\varphi_1)_\otimes = \text{id}.$$

This shows that $(\varphi_1)_\otimes$ and $(\varphi_2)_\otimes$ are inverse linear maps, and thus, $(\varphi_2)_\otimes: T_1 \rightarrow T_2$ is an isomorphism between T_1 and T_2 . \square

Now that we have shown that tensor products are unique up to isomorphism, we give a construction that produces one.

Theorem 23.5. *Given $n \geq 2$ vector spaces E_1, \dots, E_n , a tensor product $(E_1 \otimes \cdots \otimes E_n, \varphi)$ for E_1, \dots, E_n can be constructed. Furthermore, denoting $\varphi(u_1, \dots, u_n)$ as $u_1 \otimes \cdots \otimes u_n$, the tensor product $E_1 \otimes \cdots \otimes E_n$ is generated by the vectors $u_1 \otimes \cdots \otimes u_n$, where $u_1 \in E_1, \dots, u_n \in E_n$, and for every multilinear map $f: E_1 \times \cdots \times E_n \rightarrow F$, the unique linear map $f_\otimes: E_1 \otimes \cdots \otimes E_n \rightarrow F$ such that $f = f_\otimes \circ \varphi$ is defined by*

$$f_\otimes(u_1 \otimes \cdots \otimes u_n) = f(u_1, \dots, u_n)$$

on the generators $u_1 \otimes \cdots \otimes u_n$ of $E_1 \otimes \cdots \otimes E_n$.

Proof. Given any set I viewed as an index set, let $K^{(I)}$ be the set of all functions $f: I \rightarrow K$ such that $f(i) \neq 0$ only for finitely many $i \in I$. As usual, denote such a function by $(f_i)_{i \in I}$; it is a family of finite support. We make $K^{(I)}$ into a vector space by defining addition and scalar multiplication by

$$\begin{aligned} (f_i) + (g_i) &= (f_i + g_i) \\ \lambda(f_i) &= (\lambda f_i). \end{aligned}$$

The family $(e_i)_{i \in I}$ is defined such that $(e_i)_j = 0$ if $j \neq i$ and $(e_i)_i = 1$. It is a basis of the vector space $K^{(I)}$, so that every $w \in K^{(I)}$ can be uniquely written as a finite linear combination of the e_i . There is also an injection $\iota: I \rightarrow K^{(I)}$ such that $\iota(i) = e_i$ for every

$i \in I$. Furthermore, it is easy to show that for any vector space F , and for any function $f: I \rightarrow F$, there is a unique linear map $\bar{f}: K^{(I)} \rightarrow F$ such that $f = \bar{f} \circ \iota$, as in the following diagram:

$$\begin{array}{ccc} I & \xrightarrow{\iota} & K^{(I)} \\ & \searrow f & \downarrow \bar{f} \\ & & F \end{array}$$

This shows that $K^{(I)}$ is the *free vector space generated by I* . Now, apply this construction to the cartesian product $I = E_1 \times \cdots \times E_n$, obtaining the free vector space $M = K^{(I)}$ on $I = E_1 \times \cdots \times E_n$. Since every e_i is uniquely associated with some n -tuple $i = (u_1, \dots, u_n) \in E_1 \times \cdots \times E_n$, we denote e_i by (u_1, \dots, u_n) .

Next, let N be the subspace of M generated by the vectors of the following type:

$$\begin{aligned} & (u_1, \dots, u_i + v_i, \dots, u_n) - (u_1, \dots, u_i, \dots, u_n) - (u_1, \dots, v_i, \dots, u_n), \\ & (u_1, \dots, \lambda u_i, \dots, u_n) - \lambda(u_1, \dots, u_i, \dots, u_n). \end{aligned}$$

We let $E_1 \otimes \cdots \otimes E_n$ be the quotient M/N of the free vector space M by N , $\pi: M \rightarrow M/N$ be the quotient map, and set

$$\varphi = \pi \circ \iota.$$

By construction, φ is multilinear, and since π is surjective and the $\iota(i) = e_i$ generate M , since i is of the form $i = (u_1, \dots, u_n) \in E_1 \times \cdots \times E_n$, the $\varphi(u_1, \dots, u_n)$ generate M/N . Thus, if we denote $\varphi(u_1, \dots, u_n)$ as $u_1 \otimes \cdots \otimes u_n$, the tensor product $E_1 \otimes \cdots \otimes E_n$ is generated by the vectors $u_1 \otimes \cdots \otimes u_n$, where $u_1 \in E_1, \dots, u_n \in E_n$.

For every multilinear map $f: E_1 \times \cdots \times E_n \rightarrow F$, if a linear map $f_\otimes: E_1 \otimes \cdots \otimes E_n \rightarrow F$ exists such that $f = f_\otimes \circ \varphi$, since the vectors $u_1 \otimes \cdots \otimes u_n$ generate $E_1 \otimes \cdots \otimes E_n$, the map f_\otimes is uniquely defined by

$$f_\otimes(u_1 \otimes \cdots \otimes u_n) = f(u_1, \dots, u_n).$$

On the other hand, because $M = K^{(E_1 \times \cdots \times E_n)}$ is free on $I = E_1 \times \cdots \times E_n$, there is a unique linear map $\bar{f}: K^{(E_1 \times \cdots \times E_n)} \rightarrow F$, such that

$$f = \bar{f} \circ \iota,$$

as in the diagram below:

$$\begin{array}{ccc} E_1 \times \cdots \times E_n & \xrightarrow{\iota} & K^{(E_1 \times \cdots \times E_n)} \\ & \searrow f & \downarrow \bar{f} \\ & & F \end{array}$$

Because f is multilinear, note that we must have $\bar{f}(w) = 0$ for every $w \in N$. But then, $\bar{f}: M \rightarrow F$ induces a linear map $h: M/N \rightarrow F$ such that

$$f = h \circ \pi \circ \iota,$$

by defining $h([z]) = \bar{f}(z)$ for every $z \in M$, where $[z]$ denotes the equivalence class in M/N of $z \in M$:

$$\begin{array}{ccc} E_1 \times \cdots \times E_n & \xrightarrow{\pi \circ \iota} & K^{(E_1 \times \cdots \times E_n)} / N \\ & \searrow f & \downarrow h \\ & & F \end{array}$$

Indeed, the fact that \bar{f} vanishes on N insures that h is well defined on M/N , and it is clearly linear by definition. However, we showed that such a linear map h is unique, and thus it agrees with the linear map f_{\otimes} defined by

$$f_{\otimes}(u_1 \otimes \cdots \otimes u_n) = f(u_1, \dots, u_n)$$

on the generators of $E_1 \otimes \cdots \otimes E_n$. □

What is important about Theorem 23.5 is not so much the construction itself but the fact that it produces a tensor product with the universal mapping property with respect to multilinear maps. Indeed, Theorem 23.5 yields a canonical isomorphism

$$L(E_1 \otimes \cdots \otimes E_n, F) \cong L(E_1, \dots, E_n; F)$$

between the vector space of linear maps $L(E_1 \otimes \cdots \otimes E_n, F)$, and the vector space of multilinear maps $\mathcal{L}(E_1, \dots, E_n; F)$, *via* the linear map $- \circ \varphi$ defined by

$$h \mapsto h \circ \varphi,$$

where $h \in L(E_1 \otimes \cdots \otimes E_n, F)$. Indeed, $h \circ \varphi$ is clearly multilinear, and since by Theorem 23.5, for every multilinear map $f \in \mathcal{L}(E_1, \dots, E_n; F)$, there is a unique linear map $f_{\otimes} \in L(E_1 \otimes \cdots \otimes E_n, F)$ such that $f = f_{\otimes} \circ \varphi$, the map $- \circ \varphi$ is bijective. As a matter of fact, its inverse is the map

$$f \mapsto f_{\otimes}.$$

Using the “Hom” notation, the above canonical isomorphism is written

$$\text{Hom}(E_1 \otimes \cdots \otimes E_n, F) \cong \text{Hom}(E_1, \dots, E_n; F).$$

Remarks:

- (1) To be very precise, since the tensor product depends on the field K , we should subscript the symbol \otimes with K and write

$$E_1 \otimes_K \cdots \otimes_K E_n.$$

However, we often omit the subscript K unless confusion may arise.

- (2) For $F = K$, the base field, we obtain a canonical isomorphism between the vector space $L(E_1 \otimes \cdots \otimes E_n, K)$, and the vector space of multilinear forms $\mathcal{L}(E_1, \dots, E_n; K)$. However, $L(E_1 \otimes \cdots \otimes E_n, K)$ is the dual space $(E_1 \otimes \cdots \otimes E_n)^*$, and thus the vector space of multilinear forms $\mathcal{L}(E_1, \dots, E_n; K)$ is canonically isomorphic to $(E_1 \otimes \cdots \otimes E_n)^*$. We write

$$L(E_1, \dots, E_n; K) \cong (E_1 \otimes \cdots \otimes E_n)^*.$$

The fact that the map $\varphi: E_1 \times \cdots \times E_n \rightarrow E_1 \otimes \cdots \otimes E_n$ is multilinear, can also be expressed as follows:

$$\begin{aligned} u_1 \otimes \cdots \otimes (v_i + w_i) \otimes \cdots \otimes u_n &= (u_1 \otimes \cdots \otimes v_i \otimes \cdots \otimes u_n) + (u_1 \otimes \cdots \otimes w_i \otimes \cdots \otimes u_n), \\ u_1 \otimes \cdots \otimes (\lambda u_i) \otimes \cdots \otimes u_n &= \lambda(u_1 \otimes \cdots \otimes u_i \otimes \cdots \otimes u_n). \end{aligned}$$

Of course, this is just what we wanted! Tensors in $E_1 \otimes \cdots \otimes E_n$ are also called *n-tensors*, and tensors of the form $u_1 \otimes \cdots \otimes u_n$, where $u_i \in E_i$ are called *simple (or indecomposable) n-tensors*. Those *n-tensors* that are not simple are often called *compound n-tensors*.

Not only do tensor products act on spaces, but they also act on linear maps (they are functors). Given two linear maps $f: E \rightarrow E'$ and $g: F \rightarrow F'$, we can define $h: E \otimes F \rightarrow E' \otimes F'$ by

$$h(u, v) = f(u) \otimes g(v).$$

It is immediately verified that h is bilinear, and thus it induces a unique linear map

$$f \otimes g: E \otimes F \rightarrow E' \otimes F'$$

such that

$$(f \otimes g)(u \otimes v) = f(u) \otimes g(v).$$

If we also have linear maps $f': E' \rightarrow E''$ and $g': F' \rightarrow F''$, we can easily verify that the linear maps $(f' \circ f) \otimes (g' \circ g)$ and $(f' \otimes g') \circ (f \otimes g)$ agree on all vectors of the form $u \otimes v \in E \otimes F$. Since these vectors generate $E \otimes F$, we conclude that

$$(f' \circ f) \otimes (g' \circ g) = (f' \otimes g') \circ (f \otimes g).$$

The generalization to the tensor product $f_1 \otimes \cdots \otimes f_n$ of $n \geq 3$ linear maps $f_i: E_i \rightarrow F_i$ is immediate, and left to the reader.

23.2 Bases of Tensor Products

We showed that $E_1 \otimes \cdots \otimes E_n$ is generated by the vectors of the form $u_1 \otimes \cdots \otimes u_n$. However, these vectors are not linearly independent. This situation can be fixed when considering bases, which is the object of the next proposition.

Proposition 23.6. *Given $n \geq 2$ vector spaces E_1, \dots, E_n , if $(u_i^k)_{i \in I_k}$ is a basis for E_k , $1 \leq k \leq n$, then the family of vectors*

$$(u_{i_1}^1 \otimes \cdots \otimes u_{i_n}^n)_{(i_1, \dots, i_n) \in I_1 \times \dots \times I_n}$$

is a basis of the tensor product $E_1 \otimes \cdots \otimes E_n$.

Proof. For each k , $1 \leq k \leq n$, every $v^k \in E_k$ can be written uniquely as

$$v^k = \sum_{j \in I_k} v_j^k u_j^k,$$

for some family of scalars $(v_j^k)_{j \in I_k}$. Let F be any nontrivial vector space. We show that for every family

$$(w_{i_1, \dots, i_n})_{(i_1, \dots, i_n) \in I_1 \times \dots \times I_n},$$

of vectors in F , there is some linear map $h: E_1 \otimes \cdots \otimes E_n \rightarrow F$ such that

$$h(u_{i_1}^1 \otimes \cdots \otimes u_{i_n}^n) = w_{i_1, \dots, i_n}.$$

Then, by Proposition 23.3, it follows that

$$(u_{i_1}^1 \otimes \cdots \otimes u_{i_n}^n)_{(i_1, \dots, i_n) \in I_1 \times \dots \times I_n}$$

is linearly independent. However, since $(u_i^k)_{i \in I_k}$ is a basis for E_k , the $u_{i_1}^1 \otimes \cdots \otimes u_{i_n}^n$ also generate $E_1 \otimes \cdots \otimes E_n$, and thus, they form a basis of $E_1 \otimes \cdots \otimes E_n$.

We define the function $f: E_1 \times \cdots \times E_n \rightarrow F$ as follows:

$$f\left(\sum_{j_1 \in I_1} v_{j_1}^1 u_{j_1}^1, \dots, \sum_{j_n \in I_n} v_{j_n}^n u_{j_n}^n\right) = \sum_{j_1 \in I_1, \dots, j_n \in I_n} v_{j_1}^1 \cdots v_{j_n}^n w_{j_1, \dots, j_n}.$$

It is immediately verified that f is multilinear. By the universal mapping property of the tensor product, the linear map $f_\otimes: E_1 \otimes \cdots \otimes E_n \rightarrow F$ such that $f = f_\otimes \circ \varphi$, is the desired map h . \square

In particular, when each I_k is finite and of size $m_k = \dim(E_k)$, we see that the dimension of the tensor product $E_1 \otimes \cdots \otimes E_n$ is $m_1 \cdots m_n$. As a corollary of Proposition 23.6, if $(u_i^k)_{i \in I_k}$ is a basis for E_k , $1 \leq k \leq n$, then every tensor $z \in E_1 \otimes \cdots \otimes E_n$ can be written in a unique way as

$$z = \sum_{(i_1, \dots, i_n) \in I_1 \times \dots \times I_n} \lambda_{i_1, \dots, i_n} u_{i_1}^1 \otimes \cdots \otimes u_{i_n}^n,$$

for some unique family of scalars $\lambda_{i_1, \dots, i_n} \in K$, all zero except for a finite number.

23.3 Some Useful Isomorphisms for Tensor Products

Proposition 23.7. *Given 3 vector spaces E, F, G , there exists unique canonical isomorphisms*

- (1) $E \otimes F \simeq F \otimes E$
- (2) $(E \otimes F) \otimes G \simeq E \otimes (F \otimes G) \simeq E \otimes F \otimes G$
- (3) $(E \oplus F) \otimes G \simeq (E \otimes G) \oplus (F \otimes G)$
- (4) $K \otimes E \simeq E$

such that respectively

- (a) $u \otimes v \mapsto v \otimes u$
- (b) $(u \otimes v) \otimes w \mapsto u \otimes (v \otimes w) \mapsto u \otimes v \otimes w$
- (c) $(u, v) \otimes w \mapsto (u \otimes w, v \otimes w)$
- (d) $\lambda \otimes u \mapsto \lambda u$.

Proof. These isomorphisms are proved using the universal mapping property of tensor products. We illustrate the proof method on (2). Fix some $w \in G$. The map

$$(u, v) \mapsto u \otimes v \otimes w$$

from $E \times F$ to $E \otimes F \otimes G$ is bilinear, and thus there is a linear map $f_w: E \otimes F \rightarrow E \otimes F \otimes G$ such that $f_w(u \otimes v) = u \otimes v \otimes w$.

Next, consider the map

$$(z, w) \mapsto f_w(z),$$

from $(E \otimes F) \times G$ into $E \otimes F \otimes G$. It is easily seen to be bilinear, and thus it induces a linear map

$$f: (E \otimes F) \otimes G \rightarrow E \otimes F \otimes G$$

such that $f((u \otimes v) \otimes w) = u \otimes v \otimes w$.

Also consider the map

$$(u, v, w) \mapsto (u \otimes v) \otimes w$$

from $E \times F \times G$ to $(E \otimes F) \otimes G$. It is trilinear, and thus there is a linear map

$$g: E \otimes F \otimes G \rightarrow (E \otimes F) \otimes G$$

such that $g(u \otimes v \otimes w) = (u \otimes v) \otimes w$. Clearly, $f \circ g$ and $g \circ f$ are identity maps, and thus f and g are isomorphisms. The other cases are similar. \square

Given any three vector spaces, E, F, G , we have the canonical isomorphism

$$\text{Hom}(E, F; G) \cong \text{Hom}(E, \text{Hom}(F, G)).$$

Indeed, any bilinear map $f: E \times F \rightarrow G$ gives the linear map $\varphi(f) \in \text{Hom}(E, \text{Hom}(F, G))$, where $\varphi(f)(u)$ is the linear map in $\text{Hom}(F, G)$ given by

$$\varphi(f)(u)(v) = f(u, v).$$

Conversely, given a linear map $g \in \text{Hom}(E, \text{Hom}(F, G))$, we get the bilinear map $\psi(g)$ given by

$$\psi(g)(u, v) = g(u)(v),$$

and it is clear that φ and ψ are mutual inverses. Consequently, we have the important corollary:

Proposition 23.8. *For any three vector spaces, E, F, G , we have the canonical isomorphism*

$$\text{Hom}(E \otimes F, G) \cong \text{Hom}(E, \text{Hom}(F, G)).$$

23.4 Duality for Tensor Products

In this section, all vector spaces are assumed to have finite dimension. Let us now see how tensor products behave under duality. For this, we define a pairing between $E_1^* \otimes \cdots \otimes E_n^*$ and $E_1 \otimes \cdots \otimes E_n$ as follows: For any fixed $(v_1^*, \dots, v_n^*) \in E_1^* \times \cdots \times E_n^*$, we have the multilinear map

$$l_{v_1^*, \dots, v_n^*}: (u_1, \dots, u_n) \mapsto v_1^*(u_1) \cdots v_n^*(u_n)$$

from $E_1 \times \cdots \times E_n$ to K . The map $l_{v_1^*, \dots, v_n^*}$ extends uniquely to a linear map $L_{v_1^*, \dots, v_n^*}: E_1 \otimes \cdots \otimes E_n \rightarrow K$. We also have the multilinear map

$$(v_1^*, \dots, v_n^*) \mapsto L_{v_1^*, \dots, v_n^*}$$

from $E_1^* \times \cdots \times E_n^*$ to $\text{Hom}(E_1 \otimes \cdots \otimes E_n, K)$, which extends to a linear map L from $E_1^* \otimes \cdots \otimes E_n^*$ to $\text{Hom}(E_1 \otimes \cdots \otimes E_n, K)$. However, in view of the isomorphism

$$\text{Hom}(U \otimes V, W) \cong \text{Hom}(U, \text{Hom}(V, W)),$$

we can view L as a linear map

$$L: (E_1^* \otimes \cdots \otimes E_n^*) \otimes (E_1 \otimes \cdots \otimes E_n) \rightarrow K,$$

which corresponds to a bilinear map

$$(E_1^* \otimes \cdots \otimes E_n^*) \times (E_1 \otimes \cdots \otimes E_n) \rightarrow K,$$

via the isomorphism $(U \otimes V)^* \cong L(U, V; K)$. It is easy to check that this bilinear map is nondegenerate, and thus by Proposition 23.1, we have a canonical isomorphism

$$(E_1 \otimes \cdots \otimes E_n)^* \cong E_1^* \otimes \cdots \otimes E_n^*.$$

This, together with the isomorphism $L(E_1, \dots, E_n; K) \cong (E_1 \otimes \cdots \otimes E_n)^*$ yields a canonical isomorphism

$$\mu: E_1^* \otimes \cdots \otimes E_n^* \cong \text{Hom}(E_1, \dots, E_n; K).$$

Remark: The isomorphism $\mu: E_1^* \otimes \cdots \otimes E_n^* \cong \text{Hom}(E_1, \dots, E_n; K)$ can be described explicitly as the linear extension of the map given by

$$\mu(v_1^* \otimes \cdots \otimes v_n^*)(u_1, \dots, u_n) = v_1^*(u_1) \cdots v_n^*(u_n).$$

We prove another useful canonical isomorphism that allows us to treat linear maps as tensors.

Let E and F be two vector spaces and let $\alpha: E^* \times F \rightarrow \text{Hom}(E, F)$ be the map defined such that

$$\alpha(u^*, f)(x) = u^*(x)f,$$

for all $u^* \in E^*$, $f \in F$, and $x \in E$. This map is clearly bilinear, and thus it induces a linear map

$$\alpha_\otimes: E^* \otimes F \rightarrow \text{Hom}(E, F)$$

such that

$$\alpha_\otimes(u^* \otimes f)(x) = u^*(x)f.$$

Proposition 23.9. *If E and F are vector spaces, then the following properties hold:*

- (1) *The linear map $\alpha_\otimes: E^* \otimes F \rightarrow \text{Hom}(E, F)$ is injective.*
- (2) *If E is finite-dimensional, then $\alpha_\otimes: E^* \otimes F \rightarrow \text{Hom}(E, F)$ is a canonical isomorphism.*
- (3) *If F is finite-dimensional, then $\alpha_\otimes: E^* \otimes F \rightarrow \text{Hom}(E, F)$ is a canonical isomorphism.*

Proof. (1) Let $(e_i^*)_{i \in I}$ be a basis of E^* and let $(f_j)_{j \in J}$ be a basis of F . Then, we know that $(e_i^* \otimes f_j)_{i \in I, j \in J}$ is a basis of $E^* \otimes F$. To prove that α_\otimes is injective, let us show that its kernel is reduced to (0). For any vector

$$\omega = \sum_{i \in I', j \in J'} \lambda_{ij} e_i^* \otimes f_j$$

in $E^* \otimes F$, with I' and J' some finite sets, assume that $\alpha_\otimes(\omega) = 0$. This means that for every $x \in E$, we have $\alpha_\otimes(\omega)(x) = 0$; that is,

$$\sum_{i \in I', j \in J'} \alpha_\otimes(\lambda_{ij} e_i^* \otimes f_j)(x) = \sum_{j \in J'} \left(\sum_{i \in I'} \lambda_{ij} e_i^*(x) \right) f_j = 0.$$

Since $(f_j)_{j \in J}$ is a basis of F , for every $j \in J'$, we must have

$$\sum_{i \in I'} \lambda_{ij} e_i^*(x) = 0, \quad \text{for all } x \in E.$$

But, then $(e_i^*)_{i \in I'}$ would be linearly dependent, contradicting the fact that $(e_i^*)_{i \in I}$ is a basis of E^* , so we must have

$$\lambda_{ij} = 0, \quad \text{for all } i \in I' \text{ and all } j \in J',$$

which shows that $\omega = 0$. Therefore, α_\otimes is injective.

(2) Let $(e_j)_{1 \leq j \leq n}$ be a finite basis of E , and as usual, let $e_j^* \in E^*$ be the linear form defined by

$$e_j^*(e_k) = \delta_{j,k},$$

where $\delta_{j,k} = 1$ iff $j = k$ and 0 otherwise. We know that $(e_j^*)_{1 \leq j \leq n}$ is a basis of E^* (this is where we use the finite dimension of E). Now, for any linear map $f \in \text{Hom}(E, F)$, for every $x = x_1 e_1 + \cdots + x_n e_n \in E$, we have

$$f(x) = f(x_1 e_1 + \cdots + x_n e_n) = x_1 f(e_1) + \cdots + x_n f(e_n) = e_1^*(x) f(e_1) + \cdots + e_n^*(x) f(e_n).$$

Consequently, every linear map $f \in \text{Hom}(E, F)$ can be expressed as

$$f(x) = e_1^*(x) f_1 + \cdots + e_n^*(x) f_n,$$

for some $f_i \in F$. Furthermore, if we apply f to e_i , we get $f(e_i) = f_i$, so the f_i are unique. Observe that

$$(\alpha_\otimes(e_1^* \otimes f_1 + \cdots + e_n^* \otimes f_n))(x) = \sum_{i=1}^n (\alpha_\otimes(e_i^* \otimes f_i))(x) = \sum_{i=1}^n e_i^*(x) f_i.$$

Thus, α_\otimes is surjective, so α_\otimes is a bijection.

(3) Let (f_1, \dots, f_m) be a finite basis of F . It is easy to show that $\text{Hom}(E, F)$ is isomorphic to $(E^*)^m$, and similarly that $E^* \otimes F$ is isomorphic to $(E^*)^m$. The details are left as an exercise. \square

Note that in Proposition 23.9, we have an isomorphism if either E or F has finite dimension. In view of the canonical isomorphism

$$\text{Hom}(E_1, \dots, E_n; F) \cong \text{Hom}(E_1 \otimes \cdots \otimes E_n, F)$$

and the canonical isomorphism $(E_1 \otimes \cdots \otimes E_n)^* \cong E_1^* \otimes \cdots \otimes E_n^*$, where the E_i 's are finite-dimensional, Proposition 23.9 yields the canonical isomorphism

$$\text{Hom}(E_1, \dots, E_n; F) \cong E_1^* \otimes \cdots \otimes E_n^* \otimes F.$$

23.5 Tensor Algebras

The tensor product

$$\underbrace{V \otimes \cdots \otimes V}_m$$

is also denoted as

$$\bigotimes^m V \quad \text{or} \quad V^{\otimes m}$$

and is called the m -th tensor power of V (with $V^{\otimes 1} = V$, and $V^{\otimes 0} = K$). We can pack all the tensor powers of V into the “big” vector space

$$T(V) = \bigoplus_{m \geq 0} V^{\otimes m},$$

also denoted $T^\bullet(V)$ to avoid confusion with the tangent bundle. This is an interesting object because we can define a multiplication operation on it which makes it into an *algebra* called the *tensor algebra of V* . When V is of finite dimension n , this space corresponds to the algebra of polynomials with coefficients in K in n noncommuting variables.

Let us recall the definition of an algebra over a field. Let K denote any (commutative) field, although for our purposes, we may assume that $K = \mathbb{R}$ (and occasionally, $K = \mathbb{C}$). Since we will only be dealing with associative algebras with a multiplicative unit, we only define algebras of this kind.

Definition 23.3. Given a field K , a K -algebra is a K -vector space A together with a bilinear operation $\cdot : A \times A \rightarrow A$, called *multiplication*, which makes A into a ring with unity 1 (or 1_A , when we want to be very precise). This means that \cdot is associative and that there is a multiplicative identity element 1 so that $1 \cdot a = a \cdot 1 = a$, for all $a \in A$. Given two K -algebras A and B , a K -algebra homomorphism $h : A \rightarrow B$ is a linear map that is also a ring homomorphism, with $h(1_A) = 1_B$.

For example, the ring $M_n(K)$ of all $n \times n$ matrices over a field K is a K -algebra.

There is an obvious notion of *ideal* of a K -algebra: An ideal $\mathfrak{A} \subseteq A$ is a linear subspace of A that is also a two-sided ideal with respect to multiplication in A . If the field K is understood, we usually simply say an algebra instead of a K -algebra.

We would like to define a multiplication operation on $T(V)$ which makes it into a K -algebra. As

$$T(V) = \bigoplus_{i \geq 0} V^{\otimes i},$$

for every $i \geq 0$, there is a natural injection $\iota_n : V^{\otimes n} \rightarrow T(V)$, and in particular, an injection $\iota_0 : K \rightarrow T(V)$. The multiplicative unit $\mathbf{1}$ of $T(V)$ is the image $\iota_0(1)$ in $T(V)$ of the unit 1 of the field K . Since every $v \in T(V)$ can be expressed as a finite sum

$$v = \iota_{n_1}(v_1) + \cdots + \iota_{n_k}(v_k),$$

where $v_i \in V^{\otimes n_i}$ and the n_i are natural numbers with $n_i \neq n_j$ if $i \neq j$, to define multiplication in $T(V)$, using bilinearity, it is enough to define multiplication operations $\cdot: V^{\otimes m} \times V^{\otimes n} \rightarrow V^{\otimes(m+n)}$, which, using the isomorphisms $V^{\otimes n} \cong \iota_n(V^{\otimes n})$, yield multiplication operations $\cdot: \iota_m(V^{\otimes m}) \times \iota_n(V^{\otimes n}) \rightarrow \iota_{m+n}(V^{\otimes(m+n)})$. More precisely, we use the canonical isomorphism

$$V^{\otimes m} \otimes V^{\otimes n} \cong V^{\otimes(m+n)}$$

which defines a bilinear operation

$$V^{\otimes m} \times V^{\otimes n} \longrightarrow V^{\otimes(m+n)},$$

which is taken as the multiplication operation. The isomorphism $V^{\otimes m} \otimes V^{\otimes n} \cong V^{\otimes(m+n)}$ can be established by proving the isomorphisms

$$\begin{aligned} V^{\otimes m} \otimes V^{\otimes n} &\cong V^{\otimes m} \otimes \underbrace{V \otimes \cdots \otimes V}_n \\ V^{\otimes m} \otimes \underbrace{V \otimes \cdots \otimes V}_n &\cong V^{\otimes(m+n)}, \end{aligned}$$

which can be shown using methods similar to those used to prove associativity. Of course, the multiplication $V^{\otimes m} \times V^{\otimes n} \rightarrow V^{\otimes(m+n)}$ is defined so that

$$(v_1 \otimes \cdots \otimes v_m) \cdot (w_1 \otimes \cdots \otimes w_n) = v_1 \otimes \cdots \otimes v_m \otimes w_1 \otimes \cdots \otimes w_n.$$

(This has to be made rigorous by using isomorphisms involving the associativity of tensor products, for details, see Atiyah and Macdonald [5].)

Remark: It is important to note that multiplication in $T(V)$ is **not** commutative. Also, in all rigor, the unit $\mathbf{1}$ of $T(V)$ is **not equal** to 1, the unit of the field K . However, in view of the injection $\iota_0: K \rightarrow T(V)$, for the sake of notational simplicity, we will denote $\mathbf{1}$ by 1. More generally, in view of the injections $\iota_n: V^{\otimes n} \rightarrow T(V)$, we identify elements of $V^{\otimes n}$ with their images in $T(V)$.

The algebra $T(V)$ satisfies a universal mapping property which shows that it is unique up to isomorphism. For simplicity of notation, let $i: V \rightarrow T(V)$ be the natural injection of V into $T(V)$.

Proposition 23.10. *Given any K -algebra A , for any linear map $f: V \rightarrow A$, there is a unique K -algebra homomorphism $\bar{f}: T(V) \rightarrow A$ so that*

$$f = \bar{f} \circ i,$$

as in the diagram below:

$$\begin{array}{ccc} V & \xrightarrow{i} & T(V) \\ & \searrow f & \downarrow \bar{f} \\ & & A \end{array}$$

Proof. Left as an exercise (use Theorem 23.5). \square

Most algebras of interest arise as well-chosen quotients of the tensor algebra $T(V)$. This is true for the *exterior algebra* $\bigwedge(V)$ (also called *Grassmann algebra*), where we take the quotient of $T(V)$ modulo the ideal generated by all elements of the form $v \otimes v$, where $v \in V$, and for the *symmetric algebra* $\text{Sym}(V)$, where we take the quotient of $T(V)$ modulo the ideal generated by all elements of the form $v \otimes w - w \otimes v$, where $v, w \in V$.

Algebras such as $T(V)$ are graded, in the sense that there is a sequence of subspaces $V^{\otimes n} \subseteq T(V)$ such that

$$T(V) = \bigoplus_{k \geq 0} V^{\otimes k},$$

and the multiplication \otimes behaves well w.r.t. the grading, *i.e.*, $\otimes: V^{\otimes m} \times V^{\otimes n} \rightarrow V^{\otimes(m+n)}$. Generally, a K -algebra E is said to be a *graded algebra* iff there is a sequence of subspaces $E^n \subseteq E$ such that

$$E = \bigoplus_{k \geq 0} E^k,$$

(with $E^0 = K$) and the multiplication \cdot respects the grading; that is, $\cdot: E^m \times E^n \rightarrow E^{m+n}$. Elements in E^n are called *homogeneous elements of rank (or degree) n* .

In differential geometry and in physics it is necessary to consider slightly more general tensors.

Definition 23.4. Given a vector space V , for any pair of nonnegative integers (r, s) , the *tensor space* $T^{r,s}(V)$ of *type* (r, s) is the tensor product

$$T^{r,s}(V) = V^{\otimes r} \otimes (V^*)^{\otimes s} = \underbrace{V \otimes \cdots \otimes V}_r \otimes \underbrace{V^* \otimes \cdots \otimes V^*}_s,$$

with $T^{0,0}(V) = K$. We also define the *tensor algebra* $T^{\bullet,\bullet}(V)$ as the coproduct

$$T^{\bullet,\bullet}(V) = \bigoplus_{r,s \geq 0} T^{r,s}(V).$$

Tensors in $T^{r,s}(V)$ are called *homogeneous of degree* (r, s) .

Note that tensors in $T^{r,0}(V)$ are just our “old tensors” in $V^{\otimes r}$. We make $T^{\bullet,\bullet}(V)$ into an algebra by defining multiplication operations

$$T^{r_1,s_1}(V) \times T^{r_2,s_2}(V) \longrightarrow T^{r_1+r_2,s_1+s_2}(V)$$

in the usual way, namely: For $u = u_1 \otimes \cdots \otimes u_{r_1} \otimes u_1^* \otimes \cdots \otimes u_{s_1}^*$ and $v = v_1 \otimes \cdots \otimes v_{r_2} \otimes v_1^* \otimes \cdots \otimes v_{s_2}^*$, let

$$u \otimes v = u_1 \otimes \cdots \otimes u_{r_1} \otimes v_1 \otimes \cdots \otimes v_{r_2} \otimes u_1^* \otimes \cdots \otimes u_{s_1}^* \otimes v_1^* \otimes \cdots \otimes v_{s_2}^*.$$

Denote by $\text{Hom}(V^r, (V^*)^s; W)$ the vector space of all multilinear maps from $V^r \times (V^*)^s$ to W . Then, we have the universal mapping property which asserts that there is a canonical isomorphism

$$\text{Hom}(T^{r,s}(V), W) \cong \text{Hom}(V^r, (V^*)^s; W).$$

In particular,

$$(T^{r,s}(V))^* \cong \text{Hom}(V^r, (V^*)^s; K).$$

For finite dimensional vector spaces, the duality of Section 23.4 is also easily extended to the tensor spaces $T^{r,s}(V)$. We define the pairing

$$T^{r,s}(V^*) \times T^{r,s}(V) \longrightarrow K$$

as follows: If

$$v^* = v_1^* \otimes \cdots \otimes v_r^* \otimes u_{r+1} \otimes \cdots \otimes u_{r+s} \in T^{r,s}(V^*)$$

and

$$u = u_1 \otimes \cdots \otimes u_r \otimes v_{r+1}^* \otimes \cdots \otimes v_{r+s}^* \in T^{r,s}(V),$$

then

$$(v^*, u) = v_1^*(u_1) \cdots v_{r+s}^*(u_{r+s}).$$

This is a nondegenerate pairing, and thus we get a canonical isomorphism

$$(T^{r,s}(V))^* \cong T^{r,s}(V^*).$$

Consequently, we get a canonical isomorphism

$$T^{r,s}(V^*) \cong \text{Hom}(V^r, (V^*)^s; K).$$

Remark: The tensor spaces, $T^{r,s}(V)$ are also denoted $T_s^r(V)$. A tensor $\alpha \in T^{r,s}(V)$ is said to be *contravariant* in the first r arguments and *covariant* in the last s arguments. This terminology refers to the way tensors behave under coordinate changes. Given a basis (e_1, \dots, e_n) of V , if (e_1^*, \dots, e_n^*) denotes the dual basis, then every tensor $\alpha \in T^{r,s}(V)$ is given by an expression of the form

$$\alpha = \sum_{\substack{i_1, \dots, i_r \\ j_1, \dots, j_s}} a_{j_1, \dots, j_s}^{i_1, \dots, i_r} e_{i_1} \otimes \cdots \otimes e_{i_r} \otimes e_{j_1}^* \otimes \cdots \otimes e_{j_s}^*.$$

The tradition in classical tensor notation is to use lower indices on vectors and upper indices on linear forms and in accordance to *Einstein summation convention* (or *Einstein notation*) the position of the indices on the coefficients is reversed. *Einstein summation convention* is to assume that a summation is performed for all values of every index that appears simultaneously once as an upper index and once as a lower index. According to this convention, the tensor α above is written

$$\alpha = a_{j_1, \dots, j_s}^{i_1, \dots, i_r} e_{i_1} \otimes \cdots \otimes e_{i_r} \otimes e^{j_1} \otimes \cdots \otimes e^{j_s}.$$

An older view of tensors is that they are multidimensional arrays of coefficients,

$$(a_{j_1, \dots, j_s}^{i_1, \dots, i_r}),$$

subject to the rules for changes of bases.

Another operation on general tensors, contraction, is useful in differential geometry.

Definition 23.5. For all $r, s \geq 1$, the *contraction* $c_{i,j}: T^{r,s}(V) \rightarrow T^{r-1,s-1}(V)$, with $1 \leq i \leq r$ and $1 \leq j \leq s$, is the linear map defined on generators by

$$\begin{aligned} c_{i,j}(u_1 \otimes \cdots \otimes u_r \otimes v_1^* \otimes \cdots \otimes v_s^*) \\ = v_j^*(u_i) u_1 \otimes \cdots \otimes \widehat{u_i} \otimes \cdots \otimes u_r \otimes v_1^* \otimes \cdots \otimes \widehat{v_j^*} \otimes \cdots \otimes v_s^*, \end{aligned}$$

where the hat over an argument means that it should be omitted.

Let us figure out what is $c_{1,1}: T^{1,1}(V) \rightarrow \mathbb{R}$, that is $c_{1,1}: V \otimes V^* \rightarrow \mathbb{R}$. If (e_1, \dots, e_n) is a basis of V and (e_1^*, \dots, e_n^*) is the dual basis, every $h \in V \otimes V^* \cong \text{Hom}(V, V)$ can be expressed as

$$h = \sum_{i,j=1}^n a_{ij} e_i \otimes e_j^*.$$

As

$$c_{1,1}(e_i \otimes e_j^*) = \delta_{i,j},$$

we get

$$c_{1,1}(h) = \sum_{i=1}^n a_{ii} = \text{tr}(h),$$

where $\text{tr}(h)$ is the *trace* of h , where h is viewed as the linear map given by the matrix, (a_{ij}) . Actually, since $c_{1,1}$ is defined independently of any basis, $c_{1,1}$ provides an intrinsic definition of the trace of a linear map $h \in \text{Hom}(V, V)$.

Remark: Using the Einstein summation convention, if

$$\alpha = a_{j_1, \dots, j_s}^{i_1, \dots, i_r} e_{i_1} \otimes \cdots \otimes e_{i_r} \otimes e^{j_1} \otimes \cdots \otimes e^{j_s},$$

then

$$c_{k,l}(\alpha) = a_{j_1, \dots, j_{l-1}, i, j_{l+1}, \dots, j_s}^{i_1, \dots, i_{k-1}, i, i_{k+1}, \dots, i_r} e_{i_1} \otimes \cdots \otimes \widehat{e_{i_k}} \otimes \cdots \otimes e_{i_r} \otimes e^{j_1} \otimes \cdots \otimes \widehat{e^{j_l}} \otimes \cdots \otimes e^{j_s}.$$

If E and F are two K -algebras, we know that their tensor product $E \otimes F$ exists as a vector space. We can make $E \otimes F$ into an algebra as well. Indeed, we have the multilinear map

$$E \times F \times E \times F \longrightarrow E \otimes F$$

given by $(a, b, c, d) \mapsto (ac) \otimes (bd)$, where ac is the product of a and c in E and bd is the product of b and d in F . By the universal mapping property, we get a linear map,

$$E \otimes F \otimes E \otimes F \longrightarrow E \otimes F.$$

Using the isomorphism

$$E \otimes F \otimes E \otimes F \cong (E \otimes F) \otimes (E \otimes F),$$

we get a linear map

$$(E \otimes F) \otimes (E \otimes F) \longrightarrow E \otimes F,$$

and thus a bilinear map,

$$(E \otimes F) \times (E \otimes F) \longrightarrow E \otimes F$$

which is our multiplication operation in $E \otimes F$. This multiplication is determined by

$$(a \otimes b) \cdot (c \otimes d) = (ac) \otimes (bd).$$

One immediately checks that $E \otimes F$ with this multiplication is a K -algebra.

We now turn to symmetric tensors.

23.6 Symmetric Tensor Powers

Our goal is to come up with a notion of tensor product that will allow us to treat symmetric multilinear maps as linear maps. First, note that we have to restrict ourselves to a single vector space E , rather than n vector spaces E_1, \dots, E_n , so that symmetry makes sense. Recall that a multilinear map $f: E^n \rightarrow F$ is *symmetric* iff

$$f(u_{\sigma(1)}, \dots, u_{\sigma(n)}) = f(u_1, \dots, u_n),$$

for all $u_i \in E$ and all permutations, $\sigma: \{1, \dots, n\} \rightarrow \{1, \dots, n\}$. The group of permutations on $\{1, \dots, n\}$ (the *symmetric group*) is denoted \mathfrak{S}_n . The vector space of all symmetric multilinear maps $f: E^n \rightarrow F$ is denoted by $S^n(E; F)$. Note that $S^1(E; F) = \text{Hom}(E, F)$.

We could proceed directly as in Theorem 23.5 and construct symmetric tensor products from scratch. However, since we already have the notion of a tensor product, there is a more economical method. First, we define symmetric tensor powers.

Definition 23.6. An n -th *symmetric tensor power* of a vector space E , where $n \geq 1$, is a vector space S together with a symmetric multilinear map $\varphi: E^n \rightarrow S$ such that, for every vector space F and for every symmetric multilinear map $f: E^n \rightarrow F$, there is a unique linear map $f_{\odot}: S \rightarrow F$, with

$$f(u_1, \dots, u_n) = f_{\odot}(\varphi(u_1, \dots, u_n)),$$

for all $u_1, \dots, u_n \in E$, or for short

$$f = f_{\odot} \circ \varphi.$$

Equivalently, there is a unique linear map f_{\odot} such that the following diagram commutes:

$$\begin{array}{ccc} E^n & \xrightarrow{\varphi} & S \\ & \searrow f & \downarrow f_{\odot} \\ & & F \end{array}$$

First, we show that any two symmetric n -th tensor powers (S_1, φ_1) and (S_2, φ_2) for E are isomorphic.

Proposition 23.11. *Given any two symmetric n -th tensor powers (S_1, φ_1) and (S_2, φ_2) for E , there is an isomorphism $h: S_1 \rightarrow S_2$ such that*

$$\varphi_2 = h \circ \varphi_1.$$

Proof. Replace tensor product by n -th symmetric tensor power in the proof of Proposition 23.4. \square

We now give a construction that produces a symmetric n -th tensor power of a vector space E .

Theorem 23.12. *Given a vector space E , a symmetric n -th tensor power $(\text{Sym}^n(E), \varphi)$ for E can be constructed ($n \geq 1$). Furthermore, denoting $\varphi(u_1, \dots, u_n)$ as $u_1 \odot \dots \odot u_n$, the symmetric tensor power $\text{Sym}^n(E)$ is generated by the vectors $u_1 \odot \dots \odot u_n$, where $u_1, \dots, u_n \in E$, and for every symmetric multilinear map $f: E^n \rightarrow F$, the unique linear map $f_{\odot}: \text{Sym}^n(E) \rightarrow F$ such that $f = f_{\odot} \circ \varphi$ is defined by*

$$f_{\odot}(u_1 \odot \dots \odot u_n) = f(u_1, \dots, u_n)$$

on the generators $u_1 \odot \dots \odot u_n$ of $\text{Sym}^n(E)$.

Proof. The tensor power $E^{\otimes n}$ is too big, and thus we define an appropriate quotient. Let C be the subspace of $E^{\otimes n}$ generated by the vectors of the form

$$u_1 \otimes \dots \otimes u_n - u_{\sigma(1)} \otimes \dots \otimes u_{\sigma(n)},$$

for all $u_i \in E$, and all permutations $\sigma: \{1, \dots, n\} \rightarrow \{1, \dots, n\}$. We claim that the quotient space $(E^{\otimes n})/C$ does the job.

Let $p: E^{\otimes n} \rightarrow (E^{\otimes n})/C$ be the quotient map. Let $\varphi: E^n \rightarrow (E^{\otimes n})/C$ be the map

$$(u_1, \dots, u_n) \mapsto p(u_1 \otimes \dots \otimes u_n),$$

or equivalently, $\varphi = p \circ \varphi_0$, where $\varphi_0(u_1, \dots, u_n) = u_1 \otimes \cdots \otimes u_n$.

Let us denote $\varphi(u_1, \dots, u_n)$ as $u_1 \odot \cdots \odot u_n$. It is clear that φ is symmetric. Since the vectors $u_1 \otimes \cdots \otimes u_n$ generate $E^{\otimes n}$, and p is surjective, the vectors $u_1 \odot \cdots \odot u_n$ generate $(E^{\otimes n})/C$.

Given any symmetric multilinear map $f: E^n \rightarrow F$, there is a linear map $f_{\otimes}: E^{\otimes n} \rightarrow F$ such that $f = f_{\otimes} \circ \varphi_0$, as in the diagram below:

$$\begin{array}{ccc} E^n & \xrightarrow{\varphi_0} & E^{\otimes n} \\ & \searrow f & \downarrow f_{\otimes} \\ & & F \end{array}$$

However, since f is symmetric, we have $f_{\otimes}(z) = 0$ for every $z \in E^{\otimes n}$. Thus, we get an induced linear map $h: (E^{\otimes n})/C \rightarrow F$ such that $h([z]) = f_{\otimes}(z)$, where $[z]$ is the equivalence class in $(E^{\otimes n})/C$ of $z \in E^{\otimes n}$:

$$\begin{array}{ccc} E^n & \xrightarrow{p \circ \varphi_0} & (E^{\otimes n})/C \\ & \searrow f & \downarrow h \\ & & F \end{array}$$

However, if a linear map $f_{\odot}: (E^{\otimes n})/C \rightarrow F$ exists, since the vectors $u_1 \odot \cdots \odot u_n$ generate $(E^{\otimes n})/C$, we must have

$$f_{\odot}(u_1 \odot \cdots \odot u_n) = f(u_1, \dots, u_n),$$

which shows that h and f_{\odot} agree. Thus, $\text{Sym}^n(E) = (E^{\otimes n})/C$ and φ constitute a symmetric n -th tensor power of E . \square

Again, the actual construction is not important. What is important is that the symmetric n -th power has the universal mapping property with respect to symmetric multilinear maps.

Remark: The notation \odot for the commutative multiplication of symmetric tensor powers is not standard. Another notation commonly used is \cdot . We often abbreviate “symmetric tensor power” as “symmetric power.” The symmetric power $\text{Sym}^n(E)$ is also denoted $\text{Sym}^n E$ or $S(E)$. To be consistent with the use of \odot , we could have used the notation $\odot^n E$. Clearly, $\text{Sym}^1(E) \cong E$ and it is convenient to set $\text{Sym}^0(E) = K$.

The fact that the map $\varphi: E^n \rightarrow \text{Sym}^n(E)$ is symmetric and multilinear can also be expressed as follows:

$$\begin{aligned} u_1 \odot \cdots \odot (v_i + w_i) \odot \cdots \odot u_n &= (u_1 \odot \cdots \odot v_i \odot \cdots \odot u_n) + (u_1 \odot \cdots \odot w_i \odot \cdots \odot u_n), \\ u_1 \odot \cdots \odot (\lambda u_i) \odot \cdots \odot u_n &= \lambda(u_1 \odot \cdots \odot u_i \odot \cdots \odot u_n), \\ u_{\sigma(1)} \odot \cdots \odot u_{\sigma(n)} &= u_1 \odot \cdots \odot u_n, \end{aligned}$$

for all permutations $\sigma \in \mathfrak{S}_n$.

The last identity shows that the “operation” \odot is commutative. Thus, we can view the symmetric tensor $u_1 \odot \cdots \odot u_n$ as a multiset.

Theorem 23.12 yields a canonical isomorphism

$$\text{Hom}(\text{Sym}^n(E), F) \cong \mathcal{S}(E^n; F),$$

between the vector space of linear maps $\text{Hom}(\text{Sym}^n(E), F)$, and the vector space of symmetric multilinear maps $\mathcal{S}(E^n; F)$, via the linear map $- \circ \varphi$ defined by

$$h \mapsto h \circ \varphi,$$

where $h \in \text{Hom}(\text{Sym}^n(E), F)$. Indeed, $h \circ \varphi$ is clearly symmetric multilinear, and since by Theorem 23.12, for every symmetric multilinear map $f \in \mathcal{S}(E^n; F)$, there is a unique linear map $f_\odot \in \text{Hom}(\text{Sym}^n(E), F)$ such that $f = f_\odot \circ \varphi$, the map $- \circ \varphi$ is bijective. As a matter of fact, its inverse is the map

$$f \mapsto f_\odot.$$

In particular, when $F = K$, we get a canonical isomorphism

$$(\text{Sym}^n(E))^* \cong \mathcal{S}^n(E; K).$$

Symmetric tensors in $\text{Sym}^n(E)$ are also called *symmetric n -tensors*, and tensors of the form $u_1 \odot \cdots \odot u_n$, where $u_i \in E$, are called *simple (or decomposable) symmetric n -tensors*. Those symmetric n -tensors that are not simple are often called *compound symmetric n -tensors*.

Given two linear maps $f: E \rightarrow E'$ and $g: E \rightarrow E'$, we can define $h: E \times E \rightarrow \text{Sym}^2(E')$ by

$$h(u, v) = f(u) \odot g(v).$$

It is immediately verified that h is symmetric bilinear, and thus it induces a unique linear map

$$f \odot g: \text{Sym}^2(E) \rightarrow \text{Sym}^2(E'),$$

such that

$$(f \odot g)(u \odot v) = f(u) \odot g(v).$$

If we also have linear maps $f': E' \rightarrow E''$ and $g': E' \rightarrow E''$, we can easily verify that

$$(f' \circ f) \odot (g' \circ g) = (f' \odot g') \circ (f \odot g).$$

The generalization to the symmetric tensor product $f_1 \odot \cdots \odot f_n$ of $n \geq 3$ linear maps $f_i: E \rightarrow E'$ is immediate, and left to the reader.

23.7 Bases of Symmetric Powers

The vectors $u_1 \odot \cdots \odot u_n$, where $u_1, \dots, u_n \in E$, generate $\text{Sym}^n(E)$, but they are not linearly independent. We will prove a version of Proposition 23.6 for symmetric tensor powers. For this, recall that a (finite) multiset over a set I is a function $M: I \rightarrow \mathbb{N}$, such that $M(i) \neq 0$ for finitely many $i \in I$, and that the set of all multisets over I is denoted as $\mathbb{N}^{(I)}$. We let $\text{dom}(M) = \{i \in I \mid M(i) \neq 0\}$, which is a finite set. Then, for any multiset $M \in \mathbb{N}^{(I)}$, note that the sum $\sum_{i \in I} M(i)$ makes sense, since $\sum_{i \in I} M(i) = \sum_{i \in \text{dom}(M)} M(i)$, and $\text{dom}(M)$ is finite. For every multiset $M \in \mathbb{N}^{(I)}$, for any $n \geq 2$, we define the set J_M of functions $\eta: \{1, \dots, n\} \rightarrow \text{dom}(M)$, as follows:

$$J_M = \{\eta \mid \eta: \{1, \dots, n\} \rightarrow \text{dom}(M), |\eta^{-1}(i)| = M(i), i \in \text{dom}(M), \sum_{i \in I} M(i) = n\}.$$

In other words, if $\sum_{i \in I} M(i) = n$ and $\text{dom}(M) = \{i_1, \dots, i_k\}$,¹ any function $\eta \in J_M$ specifies a sequence of length n , consisting of $M(i_1)$ occurrences of i_1 , $M(i_2)$ occurrences of i_2, \dots , $M(i_k)$ occurrences of i_k . Intuitively, any η defines a “permutation” of the sequence (of length n)

$$\underbrace{\langle i_1, \dots, i_1 \rangle}_{M(i_1)}, \underbrace{\langle i_2, \dots, i_2 \rangle}_{M(i_2)}, \dots, \underbrace{\langle i_k, \dots, i_k \rangle}_{M(i_k)}.$$

Given any $k \geq 1$, and any $u \in E$, we denote

$$\underbrace{u \odot \cdots \odot u}_k$$

as $u^{\odot k}$.

We can now prove the following Proposition.

Proposition 23.13. *Given a vector space E , if $(u_i)_{i \in I}$ is a basis for E , then the family of vectors*

$$\left(u_{i_1}^{\odot M(i_1)} \odot \cdots \odot u_{i_k}^{\odot M(i_k)} \right)_{M \in \mathbb{N}^{(I)}, \sum_{i \in I} M(i) = n, \{i_1, \dots, i_k\} = \text{dom}(M)}$$

is a basis of the symmetric n -th tensor power $\text{Sym}^n(E)$.

Proof. The proof is very similar to that of Proposition 23.6. For any nontrivial vector space F , for any family of vectors

$$(w_M)_{M \in \mathbb{N}^{(I)}, \sum_{i \in I} M(i) = n},$$

we show the existence of a symmetric multilinear map $h: \text{Sym}^n(E) \rightarrow F$, such that for every $M \in \mathbb{N}^{(I)}$ with $\sum_{i \in I} M(i) = n$, we have

$$h(u_{i_1}^{\odot M(i_1)} \odot \cdots \odot u_{i_k}^{\odot M(i_k)}) = w_M,$$

¹Note that must have $k \leq n$.

where $\{i_1, \dots, i_k\} = \text{dom}(M)$. We define the map $f: E^n \rightarrow F$ as follows:

$$f\left(\sum_{j_1 \in I} v_{j_1}^1 u_{j_1}^1, \dots, \sum_{j_n \in I} v_{j_n}^n u_{j_n}^n\right) = \sum_{\substack{M \in \mathbb{N}^{(I)} \\ \sum_{i \in I} M(i) = n}} \left(\sum_{\eta \in J_M} v_{\eta(1)}^1 \cdots v_{\eta(n)}^n \right) w_M.$$

It is not difficult to verify that f is symmetric and multilinear. By the universal mapping property of the symmetric tensor product, the linear map $f_\odot: \text{Sym}^n(E) \rightarrow F$ such that $f = f_\odot \circ \varphi$, is the desired map h . Then, by Proposition 23.3, it follows that the family

$$\left(u_{i_1}^{\odot M(i_1)} \odot \cdots \odot u_{i_k}^{\odot M(i_k)} \right)_{M \in \mathbb{N}^{(I)}, \sum_{i \in I} M(i) = n, \{i_1, \dots, i_k\} = \text{dom}(M)}$$

is linearly independent. Using the commutativity of \odot , we can also show that these vectors generate $\text{Sym}^n(E)$, and thus, they form a basis for $\text{Sym}^n(E)$. The details are left as an exercise. \square

As a consequence, when I is finite, say of size $p = \dim(E)$, the dimension of $\text{Sym}^n(E)$ is the number of finite multisets (j_1, \dots, j_p) , such that $j_1 + \cdots + j_p = n$, $j_k \geq 0$. We leave as an exercise to show that this number is $\binom{p+n-1}{n}$. Thus, if $\dim(E) = p$, then the dimension of $\text{Sym}^n(E)$ is $\binom{p+n-1}{n}$. Compare with the dimension of $E^{\otimes n}$, which is p^n . In particular, when $p = 2$, the dimension of $\text{Sym}^n(E)$ is $n + 1$. This can also be seen directly.

Remark: The number $\binom{p+n-1}{n}$ is also the number of homogeneous monomials

$$X_1^{j_1} \cdots X_p^{j_p}$$

of total degree n in p variables (we have $j_1 + \cdots + j_p = n$). This is not a coincidence! Symmetric tensor products are closely related to polynomials (for more on this, see the next remark).

Given a vector space E and a basis $(u_i)_{i \in I}$ for E , Proposition 23.13 shows that every symmetric tensor $z \in \text{Sym}^n(E)$ can be written in a unique way as

$$z = \sum_{\substack{M \in \mathbb{N}^{(I)} \\ \sum_{i \in I} M(i) = n \\ \{i_1, \dots, i_k\} = \text{dom}(M)}} \lambda_M u_{i_1}^{\odot M(i_1)} \odot \cdots \odot u_{i_k}^{\odot M(i_k)},$$

for some unique family of scalars $\lambda_M \in K$, all zero except for a finite number.

This looks like a homogeneous polynomial of total degree n , where the monomials of total degree n are the symmetric tensors

$$u_{i_1}^{\odot M(i_1)} \odot \cdots \odot u_{i_k}^{\odot M(i_k)}$$

in the “indeterminates” u_i , where $i \in I$ (recall that $M(i_1) + \cdots + M(i_k) = n$). Again, this is not a coincidence. Polynomials can be defined in terms of symmetric tensors.

23.8 Some Useful Isomorphisms for Symmetric Powers

We can show the following property of the symmetric tensor product, using the proof technique of Proposition 23.7:

$$\mathrm{Sym}^n(E \oplus F) \cong \bigoplus_{k=0}^n \mathrm{Sym}^k(E) \otimes \mathrm{Sym}^{n-k}(F).$$

23.9 Duality for Symmetric Powers

In this section, all vector spaces are assumed to have finite dimension. We define a nondegenerate pairing $\mathrm{Sym}^n(E^*) \times \mathrm{Sym}^n(E) \longrightarrow K$ as follows: Consider the multilinear map

$$(E^*)^n \times E^n \longrightarrow K$$

given by

$$(v_1^*, \dots, v_n^*, u_1, \dots, u_n) \mapsto \sum_{\sigma \in \mathfrak{S}_n} v_{\sigma(1)}^*(u_1) \cdots v_{\sigma(n)}^*(u_n).$$

Note that the expression on the right-hand side is “almost” the determinant $\det(v_j^*(u_i))$, except that the sign $\mathrm{sgn}(\sigma)$ is missing (where $\mathrm{sgn}(\sigma)$ is the signature of the permutation σ ; that is, the parity of the number of transpositions into which σ can be factored). Such an expression is called a *permanent*.

It is easily checked that this expression is symmetric w.r.t. the u_i 's and also w.r.t. the v_j^* . For any fixed $(v_1^*, \dots, v_n^*) \in (E^*)^n$, we get a symmetric multilinear map

$$l_{v_1^*, \dots, v_n^*}: (u_1, \dots, u_n) \mapsto \sum_{\sigma \in \mathfrak{S}_n} v_{\sigma(1)}^*(u_1) \cdots v_{\sigma(n)}^*(u_n)$$

from E^n to K . The map $l_{v_1^*, \dots, v_n^*}$ extends uniquely to a linear map $L_{v_1^*, \dots, v_n^*}: \mathrm{Sym}^n(E) \rightarrow K$. Now, we also have the symmetric multilinear map

$$(v_1^*, \dots, v_n^*) \mapsto L_{v_1^*, \dots, v_n^*}$$

from $(E^*)^n$ to $\mathrm{Hom}(\mathrm{Sym}^n(E), K)$, which extends to a linear map L from $\mathrm{Sym}^n(E^*)$ to $\mathrm{Hom}(\mathrm{Sym}^n(E), K)$. However, in view of the isomorphism

$$\mathrm{Hom}(U \otimes V, W) \cong \mathrm{Hom}(U, \mathrm{Hom}(V, W)),$$

we can view L as a linear map

$$L: \mathrm{Sym}^n(E^*) \otimes \mathrm{Sym}^n(E) \longrightarrow K,$$

which corresponds to a bilinear map

$$\mathrm{Sym}^n(E^*) \times \mathrm{Sym}^n(E) \longrightarrow K.$$

Now, this pairing is nondegenerate. This can be shown using bases and we leave it as an exercise to the reader (see Knapp [64], Appendix A). Therefore, we get a canonical isomorphism

$$(\mathrm{Sym}^n(E))^* \cong \mathrm{Sym}^n(E^*).$$

Since we also have an isomorphism

$$(\mathrm{Sym}^n(E))^* \cong S^n(E, K),$$

we get a canonical isomorphism

$$\mu: \mathrm{Sym}^n(E^*) \cong S^n(E, K)$$

which allows us to interpret symmetric tensors over E^* as symmetric multilinear maps.

Remark: The isomorphism $\mu: \mathrm{Sym}^n(E^*) \cong S^n(E, K)$ discussed above can be described explicitly as the linear extension of the map given by

$$\mu(v_1^* \odot \cdots \odot v_n^*)(u_1, \dots, u_n) = \sum_{\sigma \in \mathfrak{S}_n} v_{\sigma(1)}^*(u_1) \cdots v_{\sigma(n)}^*(u_n).$$

Now, the map from E^n to $\mathrm{Sym}^n(E)$ given by $(u_1, \dots, u_n) \mapsto u_1 \odot \cdots \odot u_n$ yields a surjection $\pi: E^{\otimes n} \rightarrow \mathrm{Sym}^n(E)$. Because we are dealing with vector spaces, this map has a section; that is, there is some injection $\iota: \mathrm{Sym}^n(E) \rightarrow E^{\otimes n}$ with $\pi \circ \iota = \mathrm{id}$. If our field K has characteristic 0, then there is a special section having a natural definition involving a symmetrization process defined as follows: For every permutation σ , we have the map $r_\sigma: E^n \rightarrow E^{\otimes n}$ given by

$$r_\sigma(u_1, \dots, u_n) = u_{\sigma(1)} \otimes \cdots \otimes u_{\sigma(n)}.$$

As r_σ is clearly multilinear, r_σ extends to a linear map $r_\sigma: E^{\otimes n} \rightarrow E^{\otimes n}$, and we get a map $\mathfrak{S}_n \times E^{\otimes n} \rightarrow E^{\otimes n}$, namely

$$\sigma \cdot z = r_\sigma(z).$$

It is immediately checked that this is a left action of the symmetric group \mathfrak{S}_n on $E^{\otimes n}$, and the tensors $z \in E^{\otimes n}$ such that

$$\sigma \cdot z = z, \quad \text{for all } \sigma \in \mathfrak{S}_n$$

are called *symmetrized tensors*. We define the map $\iota: E^n \rightarrow E^{\otimes n}$ by

$$\iota(u_1, \dots, u_n) = \frac{1}{n!} \sum_{\sigma \in \mathfrak{S}_n} \sigma \cdot (u_1 \otimes \cdots \otimes u_n) = \frac{1}{n!} \sum_{\sigma \in \mathfrak{S}_n} u_{\sigma(1)} \otimes \cdots \otimes u_{\sigma(n)}.$$

As the right hand side is clearly symmetric, we get a linear map $\iota: \mathrm{Sym}^n(E) \rightarrow E^{\otimes n}$. Clearly, $\iota(\mathrm{Sym}^n(E))$ is the set of symmetrized tensors in $E^{\otimes n}$. If we consider the map

$S = \iota \circ \pi: E^{\otimes n} \longrightarrow E^{\otimes n}$, it is easy to check that $S \circ S = S$. Therefore, S is a projection, and by linear algebra, we know that

$$E^{\otimes n} = S(E^{\otimes n}) \oplus \text{Ker } S = \iota(\text{Sym}^n(E)) \oplus \text{Ker } S.$$

It turns out that $\text{Ker } S = E^{\otimes n} \cap \mathfrak{I} = \text{Ker } \pi$, where \mathfrak{I} is the two-sided ideal of $T(E)$ generated by all tensors of the form $u \otimes v - v \otimes u \in E^{\otimes 2}$ (for example, see Knapp [64], Appendix A). Therefore, ι is injective,

$$E^{\otimes n} = \iota(\text{Sym}^n(E)) \oplus (E^{\otimes n} \cap \mathfrak{I}) = \iota(\text{Sym}^n(E)) \oplus \text{Ker } \pi,$$

and the symmetric tensor power $\text{Sym}^n(E)$ is naturally embedded into $E^{\otimes n}$.

23.10 Symmetric Algebras

As in the case of tensors, we can pack together all the symmetric powers $\text{Sym}^n(V)$ into an algebra

$$\text{Sym}(V) = \bigoplus_{m \geq 0} \text{Sym}^m(V),$$

called the *symmetric tensor algebra of V* . We could adapt what we did in Section 23.5 for general tensor powers to symmetric tensors but since we already have the algebra $T(V)$, we can proceed faster. If \mathfrak{I} is the two-sided ideal generated by all tensors of the form $u \otimes v - v \otimes u \in V^{\otimes 2}$, we set

$$\text{Sym}^\bullet(V) = T(V)/\mathfrak{I}.$$

Then, $\text{Sym}^\bullet(V)$ automatically inherits a multiplication operation which is commutative, and since $T(V)$ is graded, that is

$$T(V) = \bigoplus_{m \geq 0} V^{\otimes m},$$

we have

$$\text{Sym}^\bullet(V) = \bigoplus_{m \geq 0} V^{\otimes m} / (\mathfrak{I} \cap V^{\otimes m}).$$

However, it is easy to check that

$$\text{Sym}^m(V) \cong V^{\otimes m} / (\mathfrak{I} \cap V^{\otimes m}),$$

so

$$\text{Sym}^\bullet(V) \cong \text{Sym}(V).$$

When V is of finite dimension n , $T(V)$ corresponds to the algebra of polynomials with coefficients in K in n variables (this can be seen from Proposition 23.13). When V is of infinite dimension and $(u_i)_{i \in I}$ is a basis of V , the algebra $\text{Sym}(V)$ corresponds to the algebra of polynomials in infinitely many variables in I . What's nice about the symmetric tensor algebra $\text{Sym}(V)$ is that it provides an intrinsic definition of a polynomial algebra in any set I of variables.

It is also easy to see that $\text{Sym}(V)$ satisfies the following universal mapping property:

Proposition 23.14. *Given any commutative K -algebra A , for any linear map $f: V \rightarrow A$, there is a unique K -algebra homomorphism $\bar{f}: \text{Sym}(V) \rightarrow A$ so that*

$$f = \bar{f} \circ i,$$

as in the diagram below:

$$\begin{array}{ccc} V & \xrightarrow{i} & \text{Sym}(V) \\ & \searrow f & \downarrow \bar{f} \\ & & A \end{array}$$

Remark: If E is finite-dimensional, recall the isomorphism $\mu: \text{Sym}^n(E^*) \rightarrow S^n(E, K)$ defined as the linear extension of the map given by

$$\mu(v_1^* \odot \cdots \odot v_n^*)(u_1, \dots, u_n) = \sum_{\sigma \in \mathfrak{S}_n} v_{\sigma(1)}^*(u_1) \cdots v_{\sigma(n)}^*(u_n).$$

Now, we have also a multiplication operation $\text{Sym}^m(E^*) \times \text{Sym}^n(E^*) \rightarrow \text{Sym}^{m+n}(E^*)$. The following question then arises:

Can we define a multiplication $S^m(E, K) \times S^n(E, K) \rightarrow S^{m+n}(E, K)$ directly on symmetric multilinear forms, so that the following diagram commutes:

$$\begin{array}{ccc} \text{Sym}^m(E^*) \times \text{Sym}^n(E^*) & \xrightarrow{\odot} & \text{Sym}^{m+n}(E^*) \\ \downarrow \mu \times \mu & & \downarrow \mu \\ S^m(E, K) \times S^n(E, K) & \xrightarrow{\cdot} & S^{m+n}(E, K). \end{array}$$

The answer is *yes*! The solution is to define this multiplication such that, for $f \in S^m(E, K)$ and $g \in S^n(E, K)$,

$$(f \cdot g)(u_1, \dots, u_{m+n}) = \sum_{\sigma \in \text{shuffle}(m, n)} f(u_{\sigma(1)}, \dots, u_{\sigma(m)}) g(u_{\sigma(m+1)}, \dots, u_{\sigma(m+n)}),$$

where $\text{shuffle}(m, n)$ consists of all (m, n) -“shuffles,” that is, permutations σ of $\{1, \dots, m+n\}$ such that $\sigma(1) < \cdots < \sigma(m)$ and $\sigma(m+1) < \cdots < \sigma(m+n)$. We urge the reader to check this fact.

Another useful canonical isomorphism (of K -algebras) is

$$\text{Sym}(E \oplus F) \cong \text{Sym}(E) \otimes \text{Sym}(F).$$

23.11 Exterior Tensor Powers

We now consider *alternating* (also called *skew-symmetric*) multilinear maps and *exterior tensor powers* (also called *alternating tensor powers*), denoted $\bigwedge^n(E)$. In many respect, alternating multilinear maps and exterior tensor powers can be treated much like symmetric tensor powers, except that the sign $\text{sgn}(\sigma)$ needs to be inserted in front of the formulae valid for symmetric powers.

Roughly speaking, we are now in the world of determinants rather than in the world of permanents. However, there are also some fundamental differences, one of which being that the exterior tensor power $\bigwedge^n(E)$ is the trivial vector space (0) when E is finite-dimensional and when $n > \dim(E)$. As in the case of symmetric tensor powers, since we already have the tensor algebra $T(V)$, we can proceed rather quickly. But first, let us review some basic definitions and facts.

Definition 23.7. Let $f: E^n \rightarrow F$ be a multilinear map. We say that f *alternating* iff $f(u_1, \dots, u_n) = 0$ whenever $u_i = u_{i+1}$, for some i with $1 \leq i \leq n-1$, for all $u_i \in E$; that is, $f(u_1, \dots, u_n) = 0$ whenever two adjacent arguments are identical. We say that f is *skew-symmetric* (or *anti-symmetric*) iff

$$f(u_{\sigma(1)}, \dots, u_{\sigma(n)}) = \text{sgn}(\sigma) f(u_1, \dots, u_n),$$

for every permutation $\sigma \in \mathfrak{S}_n$, and all $u_i \in E$.

For $n = 1$, we agree that every linear map $f: E \rightarrow F$ is alternating. The vector space of all multilinear alternating maps $f: E^n \rightarrow F$ is denoted $\text{Alt}^n(E; F)$. Note that $\text{Alt}^1(E; F) = \text{Hom}(E, F)$. The following basic proposition shows the relationship between alternation and skew-symmetry.

Proposition 23.15. *Let $f: E^n \rightarrow F$ be a multilinear map. If f is alternating, then the following properties hold:*

(1) *For all i , with $1 \leq i \leq n-1$,*

$$f(\dots, u_i, u_{i+1}, \dots) = -f(\dots, u_{i+1}, u_i, \dots).$$

(2) *For every permutation $\sigma \in \mathfrak{S}_n$,*

$$f(u_{\sigma(1)}, \dots, u_{\sigma(n)}) = \text{sgn}(\sigma) f(u_1, \dots, u_n).$$

(3) *For all i, j , with $1 \leq i < j \leq n$,*

$$f(\dots, u_i, \dots, u_j, \dots) = 0 \quad \text{whenever } u_i = u_j.$$

Moreover, if our field K has characteristic different from 2, then every skew-symmetric multilinear map is alternating.

Proof. (i) By multilinearity applied twice, we have

$$\begin{aligned} f(\dots, u_i + u_{i+1}, u_i + u_{i+1}, \dots) &= f(\dots, u_i, u_i, \dots) + f(\dots, u_i, u_{i+1}, \dots) \\ &\quad + f(\dots, u_{i+1}, u_i, \dots) + f(\dots, u_{i+1}, u_{i+1}, \dots). \end{aligned}$$

Since f is alternating, we get

$$0 = f(\dots, u_i, u_{i+1}, \dots) + f(\dots, u_{i+1}, u_i, \dots);$$

that is, $f(\dots, u_i, u_{i+1}, \dots) = -f(\dots, u_{i+1}, u_i, \dots)$.

(ii) Clearly, the symmetric group, \mathfrak{S}_n , acts on $\text{Alt}^n(E; F)$ on the left, *via*

$$\sigma \cdot f(u_1, \dots, u_n) = f(u_{\sigma(1)}, \dots, u_{\sigma(n)}).$$

Consequently, as \mathfrak{S}_n is generated by the transpositions (permutations that swap exactly two elements), since for a transposition, (ii) is simply (i), we deduce (ii) by induction on the number of transpositions in σ .

(iii) There is a permutation σ that sends u_i and u_j respectively to u_1 and u_2 . As f is alternating,

$$f(u_{\sigma(1)}, \dots, u_{\sigma(n)}) = 0.$$

However, by (ii),

$$f(u_1, \dots, u_n) = \text{sgn}(\sigma) f(u_{\sigma(1)}, \dots, u_{\sigma(n)}) = 0.$$

Now, when f is skew-symmetric, if σ is the transposition swapping u_i and $u_{i+1} = u_i$, as $\text{sgn}(\sigma) = -1$, we get

$$f(\dots, u_i, u_i, \dots) = -f(\dots, u_i, u_i, \dots),$$

so that

$$2f(\dots, u_i, u_i, \dots) = 0,$$

and in every characteristic except 2, we conclude that $f(\dots, u_i, u_i, \dots) = 0$, namely f is alternating. \square

Proposition 23.15 shows that in every characteristic except 2, alternating and skew-symmetric multilinear maps are identical. Using Proposition 23.15 we easily deduce the following crucial fact:

Proposition 23.16. *Let $f: E^n \rightarrow F$ be an alternating multilinear map. For any families of vectors, (u_1, \dots, u_n) and (v_1, \dots, v_n) , with $u_i, v_i \in E$, if*

$$v_j = \sum_{i=1}^n a_{ij} u_i, \quad 1 \leq j \leq n,$$

then

$$f(v_1, \dots, v_n) = \left(\sum_{\sigma \in \mathfrak{S}_n} \text{sgn}(\sigma) a_{\sigma(1),1} \cdots a_{\sigma(n),n} \right) f(u_1, \dots, u_n) = \det(A) f(u_1, \dots, u_n),$$

where A is the $n \times n$ matrix, $A = (a_{ij})$.

Proof. Use property (ii) of Proposition 23.15. \square

We are now ready to define and construct exterior tensor powers.

Definition 23.8. An n -th exterior tensor power of a vector space E , where $n \geq 1$, is a vector space A together with an alternating multilinear map $\varphi: E^n \rightarrow A$, such that for every vector space F and for every alternating multilinear map $f: E^n \rightarrow F$, there is a unique linear map $f_\wedge: A \rightarrow F$ with

$$f(u_1, \dots, u_n) = f_\wedge(\varphi(u_1, \dots, u_n)),$$

for all $u_1, \dots, u_n \in E$, or for short

$$f = f_\wedge \circ \varphi.$$

Equivalently, there is a unique linear map f_\wedge such that the following diagram commutes:

$$\begin{array}{ccc} E^n & \xrightarrow{\varphi} & A \\ & \searrow f & \downarrow f_\wedge \\ & & F \end{array}$$

First, we show that any two n -th exterior tensor powers (A_1, φ_1) and (A_2, φ_2) for E are isomorphic.

Proposition 23.17. *Given any two n -th exterior tensor powers (A_1, φ_1) and (A_2, φ_2) for E , there is an isomorphism $h: A_1 \rightarrow A_2$ such that*

$$\varphi_2 = h \circ \varphi_1.$$

Proof. Replace tensor product by n exterior tensor power in the proof of Proposition 23.4. \square

We now give a construction that produces an n -th exterior tensor power of a vector space E .

Theorem 23.18. *Given a vector space E , an n -th exterior tensor power $(\bigwedge^n(E), \varphi)$ for E can be constructed ($n \geq 1$). Furthermore, denoting $\varphi(u_1, \dots, u_n)$ as $u_1 \wedge \dots \wedge u_n$, the exterior tensor power $\bigwedge^n(E)$ is generated by the vectors $u_1 \wedge \dots \wedge u_n$, where $u_1, \dots, u_n \in E$, and for every alternating multilinear map $f: E^n \rightarrow F$, the unique linear map $f_\wedge: \bigwedge^n(E) \rightarrow F$ such that $f = f_\wedge \circ \varphi$ is defined by*

$$f_\wedge(u_1 \wedge \dots \wedge u_n) = f(u_1, \dots, u_n)$$

on the generators $u_1 \wedge \dots \wedge u_n$ of $\bigwedge^n(E)$.

Proof sketch. We can give a quick proof using the tensor algebra $T(E)$. let \mathfrak{I}_a be the two-sided ideal of $T(E)$ generated by all tensors of the form $u \otimes u \in E^{\otimes 2}$. Then, let

$$\bigwedge^n(E) = E^{\otimes n} / (\mathfrak{I}_a \cap E^{\otimes n})$$

and let π be the projection $\pi: E^{\otimes n} \rightarrow \bigwedge^n(E)$. If we let $u_1 \wedge \cdots \wedge u_n = \pi(u_1 \otimes \cdots \otimes u_n)$, it is easy to check that $(\bigwedge^n(E), \wedge)$ satisfies the conditions of Theorem 23.18. \square

Remark: We can also define

$$\bigwedge(E) = T(E) / \mathfrak{I}_a = \bigoplus_{n \geq 0} \bigwedge^n(E),$$

the *exterior algebra* of E . This is the skew-symmetric counterpart of $\text{Sym}(E)$, and we will study it a little later.

For simplicity of notation, we may write $\bigwedge^n E$ for $\bigwedge^n(E)$. We also abbreviate “exterior tensor power” as “exterior power.” Clearly, $\bigwedge^1(E) \cong E$, and it is convenient to set $\bigwedge^0(E) = K$.

The fact that the map $\varphi: E^n \rightarrow \bigwedge^n(E)$ is alternating and multilinear can also be expressed as follows:

$$\begin{aligned} u_1 \wedge \cdots \wedge (u_i + v_i) \wedge \cdots \wedge u_n &= (u_1 \wedge \cdots \wedge u_i \wedge \cdots \wedge u_n) \\ &\quad + (u_1 \wedge \cdots \wedge v_i \wedge \cdots \wedge u_n), \\ u_1 \wedge \cdots \wedge (\lambda u_i) \wedge \cdots \wedge u_n &= \lambda(u_1 \wedge \cdots \wedge u_i \wedge \cdots \wedge u_n), \\ u_{\sigma(1)} \wedge \cdots \wedge u_{\sigma(n)} &= \text{sgn}(\sigma) u_1 \wedge \cdots \wedge u_n, \end{aligned}$$

for all $\sigma \in \mathfrak{S}_n$.

Theorem 23.18 yields a canonical isomorphism

$$\text{Hom}(\bigwedge^n(E), F) \cong \text{Alt}^n(E; F)$$

between the vector space of linear maps $\text{Hom}(\bigwedge^n(E), F)$ and the vector space of alternating multilinear maps $\text{Alt}^n(E; F)$, *via* the linear map $- \circ \varphi$ defined by

$$h \mapsto h \circ \varphi,$$

where $h \in \text{Hom}(\bigwedge^n(E), F)$. In particular, when $F = K$, we get a canonical isomorphism

$$\left(\bigwedge^n(E) \right)^* \cong \text{Alt}^n(E; K).$$

Tensors $\alpha \in \bigwedge^n(E)$ are called *alternating n -tensors* or *alternating tensors of degree n* and we write $\deg(\alpha) = n$. Tensors of the form $u_1 \wedge \cdots \wedge u_n$, where $u_i \in E$, are called *simple* (or *decomposable*) *alternating n -tensors*. Those alternating n -tensors that are not simple are often called *compound alternating n -tensors*. Simple tensors $u_1 \wedge \cdots \wedge u_n \in \bigwedge^n(E)$ are also called *n -vectors* and tensors in $\bigwedge^n(E^*)$ are often called (*alternating*) *n -forms*.

Given two linear maps $f: E \rightarrow E'$ and $g: E \rightarrow E'$, we can define $h: E \times E \rightarrow \bigwedge^2(E')$ by

$$h(u, v) = f(u) \wedge g(v).$$

It is immediately verified that h is alternating bilinear, and thus it induces a unique linear map

$$f \wedge g: \bigwedge^2(E) \rightarrow \bigwedge^2(E')$$

such that

$$(f \wedge g)(u \wedge v) = f(u) \wedge g(v).$$

If we also have linear maps $f': E' \rightarrow E''$ and $g': E' \rightarrow E''$, we can easily verify that

$$(f' \circ f) \wedge (g' \circ g) = (f' \wedge g') \circ (f \wedge g).$$

The generalization to the alternating product $f_1 \wedge \cdots \wedge f_n$ of $n \geq 3$ linear maps $f_i: E \rightarrow E'$ is immediate, and left to the reader.

23.12 Bases of Exterior Powers

Let E be any vector space. For any basis $(u_i)_{i \in \Sigma}$ for E , we assume that some total ordering \leq on Σ has been chosen. Call the pair $((u_i)_{i \in \Sigma}, \leq)$ an *ordered basis*. Then, for any nonempty finite subset $I \subseteq \Sigma$, let

$$u_I = u_{i_1} \wedge \cdots \wedge u_{i_m},$$

where $I = \{i_1, \dots, i_m\}$, with $i_1 < \cdots < i_m$.

Since $\bigwedge^n(E)$ is generated by the tensors of the form $v_1 \wedge \cdots \wedge v_n$, with $v_i \in E$, in view of skew-symmetry, it is clear that the tensors u_I , with $|I| = n$, generate $\bigwedge^n(E)$. Actually, they form a basis.

Proposition 23.19. *Given any vector space E , if E has finite dimension $d = \dim(E)$, then for all $n > d$, the exterior power $\bigwedge^n(E)$ is trivial; that is $\bigwedge^n(E) = (0)$. If $n \leq d$ or if E is infinite dimensional, then for every ordered basis $((u_i)_{i \in \Sigma}, \leq)$, the family (u_I) is basis of $\bigwedge^n(E)$, where I ranges over finite nonempty subsets of Σ of size $|I| = n$.*

Proof. First, assume that E has finite dimension $d = \dim(E)$ and that $n > d$. We know that $\bigwedge^n(E)$ is generated by the tensors of the form $v_1 \wedge \cdots \wedge v_n$, with $v_i \in E$. If u_1, \dots, u_d is a basis of E , as every v_i is a linear combination of the u_j , when we expand $v_1 \wedge \cdots \wedge v_n$ using multilinearity, we get a linear combination of the form

$$v_1 \wedge \cdots \wedge v_n = \sum_{(j_1, \dots, j_n)} \lambda_{(j_1, \dots, j_n)} u_{j_1} \wedge \cdots \wedge u_{j_n},$$

where each (j_1, \dots, j_n) is some sequence of integers $j_k \in \{1, \dots, d\}$. As $n > d$, each sequence (j_1, \dots, j_n) must contain two identical elements. By alternation, $u_{j_1} \wedge \cdots \wedge u_{j_n} = 0$, and so $v_1 \wedge \cdots \wedge v_n = 0$. It follows that $\bigwedge^n(E) = (0)$.

Now, assume that either $\dim(E) = d$ and $n \leq d$, or that E is infinite dimensional. The argument below shows that the u_I are nonzero and linearly independent. As usual, let $u_i^* \in E^*$ be the linear form given by

$$u_i^*(u_j) = \delta_{ij}.$$

For any nonempty subset $I = \{i_1, \dots, i_n\} \subseteq \Sigma$ with $i_1 < \cdots < i_n$, let l_I be the map given by

$$l_I(v_1, \dots, v_n) = \det(u_{i_j}^*(v_k)),$$

for all $v_k \in E$. As l_I is alternating multilinear, it induces a linear map $L_I: \bigwedge^n(E) \rightarrow K$. Observe that for any nonempty finite subset $J \subseteq \Sigma$ with $|J| = n$, we have

$$L_I(u_J) = \begin{cases} 1 & \text{if } I = J \\ 0 & \text{if } I \neq J. \end{cases}$$

Note that when $\dim(E) = d$ and $n \leq d$, or when E is infinite-dimensional, the forms $u_{i_1}^*, \dots, u_{i_n}^*$ are all distinct, so the above does hold. Since $L_I(u_I) = 1$, we conclude that $u_I \neq 0$. Now, if we have a linear combination

$$\sum_I \lambda_I u_I = 0,$$

where the above sum is finite and involves nonempty finite subset $I \subseteq \Sigma$ with $|I| = n$, for every such I , when we apply L_I we get

$$\lambda_I = 0,$$

proving linear independence. □

As a corollary, if E is finite dimensional, say $\dim(E) = d$, and if $1 \leq n \leq d$, then we have

$$\dim(\bigwedge^n(E)) = \binom{n}{d},$$

and if $n > d$, then $\dim(\bigwedge^n(E)) = 0$.

Remark: When $n = 0$, if we set $u_\emptyset = 1$, then $(u_\emptyset) = (1)$ is a basis of $\bigwedge^0(V) = K$.

It follows from Proposition 23.19 that the family $(u_I)_I$ where $I \subseteq \Sigma$ ranges over finite subsets of Σ is a basis of $\bigwedge(V) = \bigoplus_{n \geq 0} \bigwedge^n(V)$.

As a corollary of Proposition 23.19 we obtain the following useful criterion for linear independence:

Proposition 23.20. *For any vector space E , the vectors $u_1, \dots, u_n \in E$ are linearly independent iff $u_1 \wedge \dots \wedge u_n \neq 0$.*

Proof. If $u_1 \wedge \dots \wedge u_n \neq 0$, then u_1, \dots, u_n must be linearly independent. Otherwise, some u_i would be a linear combination of the other u_j 's (with $j \neq i$), and then, as in the proof of Proposition 23.19, $u_1 \wedge \dots \wedge u_n$ would be a linear combination of wedges in which two vectors are identical, and thus zero.

Conversely, assume that u_1, \dots, u_n are linearly independent. Then, we have the linear forms $u_i^* \in E^*$ such that

$$u_i^*(u_j) = \delta_{i,j} \quad 1 \leq i, j \leq n.$$

As in the proof of Proposition 23.19, we have a linear map $L_{u_1, \dots, u_n}: \bigwedge^n(E) \rightarrow K$, given by

$$L_{u_1, \dots, u_n}(v_1 \wedge \dots \wedge v_n) = \det(u_j^*(v_i)),$$

for all $v_1 \wedge \dots \wedge v_n \in \bigwedge^n(E)$. As,

$$L_{u_1, \dots, u_n}(u_1 \wedge \dots \wedge u_n) = 1,$$

we conclude that $u_1 \wedge \dots \wedge u_n \neq 0$. □

Proposition 23.20 shows that, geometrically, every nonzero wedge $u_1 \wedge \dots \wedge u_n$ corresponds to some oriented version of an n -dimensional subspace of E .

23.13 Some Useful Isomorphisms for Exterior Powers

We can show the following property of the exterior tensor product, using the proof technique of Proposition 23.7:

$$\bigwedge^n(E \oplus F) \cong \bigoplus_{k=0}^n \bigwedge^k(E) \otimes \bigwedge^{n-k}(F).$$

23.14 Duality for Exterior Powers

In this section, all vector spaces are assumed to have finite dimension. We define a nondegenerate pairing $\bigwedge^n(E^*) \times \bigwedge^n(E) \longrightarrow K$ as follows: Consider the multilinear map

$$(E^*)^n \times E^n \longrightarrow K$$

given by

$$(v_1^*, \dots, v_n^*, u_1, \dots, u_n) \mapsto \sum_{\sigma \in \mathfrak{S}_n} \text{sgn}(\sigma) v_{\sigma(1)}^*(u_1) \cdots v_{\sigma(n)}^*(u_n) = \det(v_j^*(u_i)).$$

It is easily checked that this expression is alternating w.r.t. the u_i 's and also w.r.t. the v_j^* . For any fixed $(v_1^*, \dots, v_n^*) \in (E^*)^n$, we get an alternating multilinear map

$$l_{v_1^*, \dots, v_n^*}: (u_1, \dots, u_n) \mapsto \det(v_j^*(u_i))$$

from E^n to K . By the argument used in the symmetric case, we get a bilinear map

$$\bigwedge^n(E^*) \times \bigwedge^n(E) \longrightarrow K.$$

Now, this pairing is nondegenerate. This can be shown using bases and we leave it as an exercise to the reader. Therefore, we get a canonical isomorphism

$$(\bigwedge^n(E))^* \cong \bigwedge^n(E^*).$$

Since we also have a canonical isomorphism

$$(\bigwedge^n(E))^* \cong \text{Alt}^n(E; K),$$

we get a canonical isomorphism

$$\mu: \bigwedge^n(E^*) \cong \text{Alt}^n(E; K),$$

which allows us to interpret alternating tensors over E^* as alternating multilinear maps.

The isomorphism $\mu: \bigwedge^n(E^*) \cong \text{Alt}^n(E; K)$ discussed above can be described explicitly as the linear extension of the map given by

$$\mu(v_1^* \wedge \cdots \wedge v_n^*)(u_1, \dots, u_n) = \det(v_j^*(u_i)).$$

Remark: Variants of our isomorphism μ are found in the literature. For example, there is a version μ' , where

$$\mu' = \frac{1}{n!} \mu,$$

with the factor $\frac{1}{n!}$ added in front of the determinant. Each version has its own merits and inconvenients. Morita [83] uses μ' because it is more convenient than μ when dealing with characteristic classes. On the other hand, when using μ' , some extra factor is needed in defining the wedge operation of alternating multilinear forms (see Section 23.15) and for exterior differentiation. The version μ is the one adopted by Warner [114], Knapp [64], Fulton and Harris [42], and Cartan [20, 21].

If $f: E \rightarrow F$ is any linear map, by transposition we get a linear map $f^\top: F^* \rightarrow E^*$ given by

$$f^\top(v^*) = v^* \circ f, \quad v^* \in F^*.$$

Consequently, we have

$$f^\top(v^*)(u) = v^*(f(u)), \quad \text{for all } u \in E \text{ and all } v^* \in F^*.$$

For any $p \geq 1$, the map

$$(u_1, \dots, u_p) \mapsto f(u_1) \wedge \dots \wedge f(u_p)$$

from E^n to $\bigwedge^p F$ is multilinear alternating, so it induces a linear map $\bigwedge^p f: \bigwedge^p E \rightarrow \bigwedge^p F$ defined on generators by

$$\left(\bigwedge^p f\right)(u_1 \wedge \dots \wedge u_p) = f(u_1) \wedge \dots \wedge f(u_p).$$

Combining \bigwedge^p and duality, we get a linear map $\bigwedge^p f^\top: \bigwedge^p F^* \rightarrow \bigwedge^p E^*$ defined on generators by

$$\left(\bigwedge^p f^\top\right)(v_1^* \wedge \dots \wedge v_p^*) = f^\top(v_1^*) \wedge \dots \wedge f^\top(v_p^*).$$

Proposition 23.21. *If $f: E \rightarrow F$ is any linear map between two finite-dimensional vector spaces E and F , then*

$$\mu\left(\left(\bigwedge^p f^\top\right)(\omega)\right)(u_1, \dots, u_p) = \mu(\omega)(f(u_1), \dots, f(u_p)), \quad \omega \in \bigwedge^p F^*, \quad u_1, \dots, u_p \in E.$$

Proof. It is enough to prove the formula on generators. By definition of μ , we have

$$\begin{aligned} \mu\left(\left(\bigwedge^p f^\top\right)(v_1^* \wedge \dots \wedge v_p^*)\right)(u_1, \dots, u_p) &= \mu(f^\top(v_1^*) \wedge \dots \wedge f^\top(v_p^*))(u_1, \dots, u_p) \\ &= \det(f^\top(v_j^*)(u_i)) \\ &= \det(v_j^*(f(u_i))) \\ &= \mu(v_1^* \wedge \dots \wedge v_p^*)(f(u_1), \dots, f(u_p)), \end{aligned}$$

as claimed. □

The map $\bigwedge^p f^\top$ is often denoted f^* , although this is an ambiguous notation since p is dropped. Proposition 23.21 gives us the behavior of f^* under the identification of $\bigwedge^p E^*$ and $\text{Alt}^p(E; K)$ via the isomorphism μ .

As in the case of symmetric powers, the map from E^n to $\bigwedge^n(E)$ given by $(u_1, \dots, u_n) \mapsto u_1 \wedge \dots \wedge u_n$ yields a surjection $\pi: E^{\otimes n} \rightarrow \bigwedge^n(E)$. Now, this map has some section, so there is some injection $\iota: \bigwedge^n(E) \rightarrow E^{\otimes n}$ with $\pi \circ \iota = \text{id}$. If our field K has characteristic 0, then there is a special section having a natural definition involving an antisymmetrization process.

Recall that we have a left action of the symmetric group \mathfrak{S}_n on $E^{\otimes n}$. The tensors $z \in E^{\otimes n}$ such that

$$\sigma \cdot z = \text{sgn}(\sigma) z, \quad \text{for all } \sigma \in \mathfrak{S}_n$$

are called *antisymmetrized* tensors. We define the map $\iota: E^n \rightarrow E^{\otimes n}$ by

$$\iota(u_1, \dots, u_n) = \frac{1}{n!} \sum_{\sigma \in \mathfrak{S}_n} \text{sgn}(\sigma) u_{\sigma(1)} \otimes \dots \otimes u_{\sigma(n)}.$$

As the right hand side is clearly an alternating map, we get a linear map $\iota: \bigwedge^n(E) \rightarrow E^{\otimes n}$. Clearly, $\iota(\bigwedge^n(E))$ is the set of antisymmetrized tensors in $E^{\otimes n}$. If we consider the map $A = \iota \circ \pi: E^{\otimes n} \rightarrow E^{\otimes n}$, it is easy to check that $A \circ A = A$. Therefore, A is a projection, and by linear algebra, we know that

$$E^{\otimes n} = A(E^{\otimes n}) \oplus \text{Ker } A = \iota\left(\bigwedge^n(E)\right) \oplus \text{Ker } A.$$

It turns out that $\text{Ker } A = E^{\otimes n} \cap \mathfrak{I}_a = \text{Ker } \pi$, where \mathfrak{I}_a is the two-sided ideal of $T(E)$ generated by all tensors of the form $u \otimes u \in E^{\otimes 2}$ (for example, see Knapp [64], Appendix A). Therefore, ι is injective,

$$E^{\otimes n} = \iota\left(\bigwedge^n(E)\right) \oplus (E^{\otimes n} \cap \mathfrak{I}) = \iota\left(\bigwedge^n(E)\right) \oplus \text{Ker } \pi,$$

and the exterior tensor power $\bigwedge^n(E)$ is naturally embedded into $E^{\otimes n}$.

23.15 Exterior Algebras

As in the case of symmetric tensors, we can pack together all the exterior powers $\bigwedge^n(V)$ into an algebra

$$\bigwedge(V) = \bigoplus_{m \geq 0} \bigwedge^m(V),$$

called the *exterior algebra* (or *Grassmann algebra*) of V . We mimic the procedure used for symmetric powers. If \mathfrak{I}_a is the two-sided ideal generated by all tensors of the form $u \otimes u \in V^{\otimes 2}$, we set

$$\dot{\bigwedge}(V) = T(V)/\mathfrak{I}_a.$$

Then, $\bigwedge^\bullet(V)$ automatically inherits a multiplication operation, called *wedge product*, and since $T(V)$ is graded, that is

$$T(V) = \bigoplus_{m \geq 0} V^{\otimes m},$$

we have

$$\bigwedge^\bullet(V) = \bigoplus_{m \geq 0} V^{\otimes m} / (\mathfrak{I}_a \cap V^{\otimes m}).$$

However, it is easy to check that

$$\bigwedge^m(V) \cong V^{\otimes m} / (\mathfrak{I}_a \cap V^{\otimes m}),$$

so

$$\bigwedge^\bullet(V) \cong \bigwedge(V).$$

When V has finite dimension d , we actually have a finite coproduct

$$\bigwedge(V) = \bigoplus_{m=0}^d \bigwedge^m(V),$$

and since each $\bigwedge^m(V)$ has dimension $\binom{d}{m}$, we deduce that

$$\dim(\bigwedge(V)) = 2^d = 2^{\dim(V)}.$$

The multiplication, $\wedge: \bigwedge^m(V) \times \bigwedge^n(V) \rightarrow \bigwedge^{m+n}(V)$, is skew-symmetric in the following precise sense:

Proposition 23.22. *For all $\alpha \in \bigwedge^m(V)$ and all $\beta \in \bigwedge^n(V)$, we have*

$$\beta \wedge \alpha = (-1)^{mn} \alpha \wedge \beta.$$

Proof. Since $v \wedge u = -u \wedge v$ for all $u, v \in V$, Proposition 23.22 follows by induction. \square

Since $\alpha \wedge \alpha = 0$ for every *simple* tensor $\alpha = u_1 \wedge \cdots \wedge u_n$, it seems natural to infer that $\alpha \wedge \alpha = 0$ for *every* tensor $\alpha \in \bigwedge(V)$. If we consider the case where $\dim(V) \leq 3$, we can indeed prove the above assertion. However, if $\dim(V) \geq 4$, the above fact is generally false! For example, when $\dim(V) = 4$, if u_1, u_2, u_3, u_4 are a basis for V , for $\alpha = u_1 \wedge u_2 + u_3 \wedge u_4$, we check that

$$\alpha \wedge \alpha = 2u_1 \wedge u_2 \wedge u_3 \wedge u_4,$$

which is nonzero.

The above discussion suggests that it might be useful to know when an alternating tensor is simple, that is, decomposable. It can be shown that for tensors $\alpha \in \bigwedge^2(V)$, $\alpha \wedge \alpha = 0$ iff α is simple. A general criterion for decomposability can be given in terms of some operations known as *left hook* and *right hook* (also called *interior products*); see Section 23.17.

It is easy to see that $\bigwedge(V)$ satisfies the following universal mapping property:

Proposition 23.23. *Given any K -algebra A , for any linear map $f: V \rightarrow A$, if $(f(v))^2 = 0$ for all $v \in V$, then there is a unique K -algebra homomorphism $\bar{f}: \Lambda(V) \rightarrow A$ so that*

$$f = \bar{f} \circ i,$$

as in the diagram below:

$$\begin{array}{ccc} V & \xrightarrow{i} & \Lambda(V) \\ & \searrow f & \downarrow \bar{f} \\ & & A \end{array}$$

When E is finite-dimensional, recall the isomorphism $\mu: \Lambda^n(E^*) \rightarrow \text{Alt}^n(E; K)$, defined as the linear extension of the map given by

$$\mu(v_1^* \wedge \cdots \wedge v_n^*)(u_1, \dots, u_n) = \det(v_j^*(u_i)).$$

Now, we have also a multiplication operation $\Lambda^m(E^*) \times \Lambda^n(E^*) \rightarrow \Lambda^{m+n}(E^*)$. The following question then arises:

Can we define a multiplication $\text{Alt}^m(E; K) \times \text{Alt}^n(E; K) \rightarrow \text{Alt}^{m+n}(E; K)$ directly on alternating multilinear forms, so that the following diagram commutes:

$$\begin{array}{ccc} \Lambda^m(E^*) \times \Lambda^n(E^*) & \xrightarrow{\wedge} & \Lambda^{m+n}(E^*) \\ \downarrow \mu \times \mu & & \downarrow \mu \\ \text{Alt}^m(E; K) \times \text{Alt}^n(E; K) & \xrightarrow{\wedge} & \text{Alt}^{m+n}(E; K). \end{array}$$

As in the symmetric case, the answer is *yes!* The solution is to define this multiplication such that, for $f \in \text{Alt}^m(E; K)$ and $g \in \text{Alt}^n(E; K)$,

$$(f \wedge g)(u_1, \dots, u_{m+n}) = \sum_{\sigma \in \text{shuffle}(m, n)} \text{sgn}(\sigma) f(u_{\sigma(1)}, \dots, u_{\sigma(m)}) g(u_{\sigma(m+1)}, \dots, u_{\sigma(m+n)}),$$

where $\text{shuffle}(m, n)$ consists of all (m, n) -“shuffles;” that is, permutations σ of $\{1, \dots, m+n\}$ such that $\sigma(1) < \cdots < \sigma(m)$ and $\sigma(m+1) < \cdots < \sigma(m+n)$. For example, when $m = n = 1$, we have

$$(f \wedge g)(u, v) = f(u)g(v) - g(u)f(v).$$

When $m = 1$ and $n \geq 2$, check that

$$(f \wedge g)(u_1, \dots, u_{m+1}) = \sum_{i=1}^{m+1} (-1)^{i-1} f(u_i) g(u_1, \dots, \hat{u}_i, \dots, u_{m+1}),$$

where the hat over the argument u_i means that it should be omitted.

As a result of all this, the coproduct

$$\text{Alt}(E) = \bigoplus_{n \geq 0} \text{Alt}^n(E; K)$$

is an algebra under the above multiplication, and this algebra is isomorphic to $\bigwedge(E^*)$. For the record, we state

Proposition 23.24. *When E is finite dimensional, the maps $\mu: \bigwedge^n(E^*) \rightarrow \text{Alt}^n(E; K)$ induced by the linear extensions of the maps given by*

$$\mu(v_1^* \wedge \cdots \wedge v_n^*)(u_1, \dots, u_n) = \det(v_j^*(u_i))$$

yield a canonical isomorphism of algebras $\mu: \bigwedge(E^) \rightarrow \text{Alt}(E)$, where the multiplication in $\text{Alt}(E)$ is defined by the maps $\wedge: \text{Alt}^m(E; K) \times \text{Alt}^n(E; K) \rightarrow \text{Alt}^{m+n}(E; K)$, with*

$$(f \wedge g)(u_1, \dots, u_{m+n}) = \sum_{\sigma \in \text{shuffle}(m, n)} \text{sgn}(\sigma) f(u_{\sigma(1)}, \dots, u_{\sigma(m)}) g(u_{\sigma(m+1)}, \dots, u_{\sigma(m+n)}),$$

where $\text{shuffle}(m, n)$ consists of all (m, n) -“shuffles,” that is, permutations σ of $\{1, \dots, m+n\}$ such that $\sigma(1) < \cdots < \sigma(m)$ and $\sigma(m+1) < \cdots < \sigma(m+n)$.

Remark: The algebra $\bigwedge(E)$ is a graded algebra. Given two graded algebras E and F , we can make a new tensor product $E \hat{\otimes} F$, where $E \hat{\otimes} F$ is equal to $E \otimes F$ as a vector space, but with a skew-commutative multiplication given by

$$(a \otimes b) \wedge (c \otimes d) = (-1)^{\deg(b)\deg(c)} (ac) \otimes (bd),$$

where $a \in E^m, b \in F^p, c \in E^n, d \in F^q$. Then, it can be shown that

$$\bigwedge(E \oplus F) \cong \bigwedge(E) \hat{\otimes} \bigwedge(F).$$

23.16 The Hodge *-Operator

In order to define a generalization of the Laplacian that applies to differential forms on a Riemannian manifold, we need to define isomorphisms

$$\bigwedge^k V \rightarrow \bigwedge^{n-k} V,$$

for any Euclidean vector space V of dimension n and any k , with $0 \leq k \leq n$. If $\langle -, - \rangle$ denotes the inner product on V , we define an inner product on $\bigwedge^k V$, also denoted $\langle -, - \rangle$, by setting

$$\langle u_1 \wedge \cdots \wedge u_k, v_1 \wedge \cdots \wedge v_k \rangle = \det(\langle u_i, v_j \rangle),$$

for all $u_i, v_i \in V$, and extending $\langle -, - \rangle$ by bilinearity.

It is easy to show that if (e_1, \dots, e_n) is an orthonormal basis of V , then the basis of $\bigwedge^k V$ consisting of the e_I (where $I = \{i_1, \dots, i_k\}$, with $1 \leq i_1 < \dots < i_k \leq n$) is an orthonormal basis of $\bigwedge^k V$. Since the inner product on V induces an inner product on V^* (recall that $\langle \omega_1, \omega_2 \rangle = \langle \omega_1^\sharp, \omega_2^\sharp \rangle$, for all $\omega_1, \omega_2 \in V^*$), we also get an inner product on $\bigwedge^k V^*$.

Recall that an *orientation* of a vector space V of dimension n is given by the choice of some basis (e_1, \dots, e_n) . We say that a basis (u_1, \dots, u_n) of V is *positively oriented* iff $\det(u_1, \dots, u_n) > 0$ (where $\det(u_1, \dots, u_n)$ denotes the determinant of the matrix whose j th column consists of the coordinates of u_j over the basis (e_1, \dots, e_n)), otherwise it is *negatively oriented*. An *oriented vector space* is a vector space V together with an orientation of V . If V is oriented by the basis (e_1, \dots, e_n) , then V^* is oriented by the dual basis (e_1^*, \dots, e_n^*) .

If V is an oriented vector space of dimension n , then we can define a linear map

$$*: \bigwedge^k V \rightarrow \bigwedge^{n-k} V,$$

called the *Hodge *-operator*, as follows: For any choice of a positively oriented orthonormal basis (e_1, \dots, e_n) of V , set

$$*(e_1 \wedge \dots \wedge e_k) = e_{k+1} \wedge \dots \wedge e_n.$$

In particular, for $k = 0$ and $k = n$, we have

$$\begin{aligned} *(1) &= e_1 \wedge \dots \wedge e_n \\ *(e_1 \wedge \dots \wedge e_n) &= 1. \end{aligned}$$

It is easy to see that the definition of $*$ does not depend on the choice of positively oriented orthonormal basis.

The Hodge *-operators $*: \bigwedge^k V \rightarrow \bigwedge^{n-k} V$ induces a linear bijection $*: \bigwedge(V) \rightarrow \bigwedge(V)$. We also have Hodge *-operators $*: \bigwedge^k V^* \rightarrow \bigwedge^{n-k} V^*$.

The following proposition is easy to show:

Proposition 23.25. *If V is any oriented vector space of dimension n , for every k with $0 \leq k \leq n$, we have*

$$(i) \quad ** = (-\text{id})^{k(n-k)}.$$

$$(ii) \quad \langle x, y \rangle = *(x \wedge *y) = *(y \wedge *x), \text{ for all } x, y \in \bigwedge^k V.$$

If (e_1, \dots, e_n) is an orthonormal basis of V and (v_1, \dots, v_n) is any other basis of V , it is easy to see that

$$v_1 \wedge \dots \wedge v_n = \sqrt{\det(\langle v_i, v_j \rangle)} e_1 \wedge \dots \wedge e_n,$$

from which it follows that

$$*(1) = \frac{1}{\sqrt{\det(\langle v_i, v_j \rangle)}} v_1 \wedge \dots \wedge v_n$$

(see Jost [61], Chapter 2, Lemma 2.1.3).

23.17 Testing Decomposability; Left and Right Hooks

In this section, all vector spaces are assumed to have finite dimension. Say $\dim(E) = n$. Using our nonsingular pairing

$$\langle -, - \rangle: \bigwedge^p E^* \times \bigwedge^p E \longrightarrow K \quad (1 \leq p \leq n)$$

defined on generators by

$$\langle u_1^* \wedge \cdots \wedge u_p^*, v_1 \wedge \cdots \wedge v_p \rangle = \det(u_i^*(v_j)),$$

we define various contraction operations

$$\lrcorner : \bigwedge^p E \times \bigwedge^{p+q} E^* \longrightarrow \bigwedge^q E^* \quad (\text{left hook})$$

and

$$\lrcorner : \bigwedge^{p+q} E^* \times \bigwedge^p E \longrightarrow \bigwedge^q E^* \quad (\text{right hook}),$$

as well as the versions obtained by replacing E by E^* and E^{**} by E . We begin with the *left interior product or left hook*, \lrcorner .

Let $u \in \bigwedge^p E$. For any q such that $p + q \leq n$, multiplication on the right by u is a linear map

$$\wedge_R(u): \bigwedge^q E \longrightarrow \bigwedge^{p+q} E$$

given by

$$v \mapsto v \wedge u$$

where $v \in \bigwedge^q E$. The transpose of $\wedge_R(u)$ yields a linear map

$$(\wedge_R(u))^t: (\bigwedge^{p+q} E)^* \longrightarrow (\bigwedge^q E)^*,$$

which, using the isomorphisms $(\bigwedge^{p+q} E)^* \cong \bigwedge^{p+q} E^*$ and $(\bigwedge^q E)^* \cong \bigwedge^q E^*$, can be viewed as a map

$$(\wedge_R(u))^t: \bigwedge^{p+q} E^* \longrightarrow \bigwedge^q E^*$$

given by

$$z^* \mapsto z^* \circ \wedge_R(u),$$

where $z^* \in \bigwedge^{p+q} E^*$.

We denote $z^* \circ \wedge_R(u)$ by

$$u \lrcorner z^*.$$

In terms of our pairing, the q -vector $u \lrcorner z^*$ is uniquely defined by

$$\langle u \lrcorner z^*, v \rangle = \langle z^*, v \wedge u \rangle, \quad \text{for all } u \in \bigwedge^p E, v \in \bigwedge^q E \text{ and } z^* \in \bigwedge^{p+q} E^*.$$

It is immediately verified that

$$(u \wedge v) \lrcorner z^* = u \lrcorner (v \lrcorner z^*),$$

so \lrcorner defines a left action

$$\lrcorner : \bigwedge^p E \times \bigwedge^{p+q} E^* \longrightarrow \bigwedge^q E^*.$$

By interchanging E and E^* and using the isomorphism

$$\left(\bigwedge^k F\right)^* \cong \bigwedge^k F^*,$$

we can also define a left action

$$\lrcorner : \bigwedge^p E^* \times \bigwedge^{p+q} E \longrightarrow \bigwedge^q E.$$

In terms of our pairing, $u^* \lrcorner z$ is uniquely defined by

$$\langle v^*, u^* \lrcorner z \rangle = \langle v^* \wedge u^*, z \rangle, \quad \text{for all } u^* \in \bigwedge^p E^*, v^* \in \bigwedge^q E^* \text{ and } z \in \bigwedge^{p+q} E.$$

In order to proceed any further, we need some combinatorial properties of the basis of $\bigwedge^p E$ constructed from a basis (e_1, \dots, e_n) of E . Recall that for any (nonempty) subset $I \subseteq \{1, \dots, n\}$, we let

$$e_I = e_{i_1} \wedge \cdots \wedge e_{i_p},$$

where $I = \{i_1, \dots, i_p\}$ with $i_1 < \cdots < i_p$. We also let $e_\emptyset = 1$.

Given any two subsets $H, L \subseteq \{1, \dots, n\}$, let

$$\rho_{H,L} = \begin{cases} 0 & \text{if } H \cap L \neq \emptyset, \\ (-1)^\nu & \text{if } H \cap L = \emptyset, \end{cases}$$

where

$$\nu = |\{(h, l) \mid (h, l) \in H \times L, h > l\}|.$$

Proposition 23.26. *For any basis (e_1, \dots, e_n) of E the following properties hold:*

(1) *If $H \cap L = \emptyset$, $|H| = h$, and $|L| = l$, then*

$$\rho_{H,L} \rho_{L,H} = (-1)^{hl}.$$

(2) For $H, L \subseteq \{1, \dots, m\}$, we have

$$e_H \wedge e_L = \rho_{H,L} e_{H \cup L}.$$

(3) For the left hook

$$\lrcorner : \bigwedge^p E \times \bigwedge^{p+q} E^* \longrightarrow \bigwedge^q E^*,$$

we have

$$\begin{aligned} e_H \lrcorner e_L^* &= 0 \quad \text{if } H \not\subseteq L \\ e_H \lrcorner e_L^* &= \rho_{L-H, H} e_{L-H}^* \quad \text{if } H \subseteq L. \end{aligned}$$

Similar formulae hold for $\lrcorner : \bigwedge^p E^* \times \bigwedge^{p+q} E \longrightarrow \bigwedge^q E$. Using Proposition 23.26, we have the

Proposition 23.27. *For the left hook*

$$\lrcorner : \bigwedge^p E \times \bigwedge^{p+q} E^* \longrightarrow \bigwedge^q E^*,$$

for every $u \in E$, we have

$$u \lrcorner (x^* \wedge y^*) = (-1)^s (u \lrcorner x^*) \wedge y^* + x^* \wedge (u \lrcorner y^*),$$

where $y \in \bigwedge^s E^*$.

Proof. We can prove the above identity assuming that x^* and y^* are of the form e_I^* and e_J^* using Proposition 23.26, but this is rather tedious. There is also a proof involving determinants; see Warner [114], Chapter 2. \square

Thus, \lrcorner is almost an anti-derivation, except that the sign $(-1)^s$ is applied to the wrong factor.

It is also possible to define a *right interior product or right hook* \lrcorner , using multiplication on the left rather than multiplication on the right. Then, \lrcorner defines a right action

$$\lrcorner : \bigwedge^{p+q} E^* \times \bigwedge^p E \longrightarrow \bigwedge^q E^*$$

such that

$$\langle z^*, u \wedge v \rangle = \langle z^* \lrcorner u, v \rangle, \quad \text{for all } u \in \bigwedge^p E, v \in \bigwedge^q E, \text{ and } z^* \in \bigwedge^{p+q} E^*.$$

Similarly, we have the right action

$$\lrcorner : \bigwedge^{p+q} E \times \bigwedge^p E^* \longrightarrow \bigwedge^q E,$$

such that

$$\langle u^* \wedge v^*, z \rangle = \langle v^*, z \lrcorner u^* \rangle, \quad \text{for all } u^* \in \bigwedge^p E^*, v^* \in \bigwedge^q E^*, \text{ and } z \in \bigwedge^{p+q} E.$$

Since the left hook $\lrcorner : \bigwedge^p E \times \bigwedge^{p+q} E^* \longrightarrow \bigwedge^q E^*$ is defined by

$$\langle u \lrcorner z^*, v \rangle = \langle z^*, v \wedge u \rangle, \quad \text{for all } u \in \bigwedge^p E, v \in \bigwedge^q E \text{ and } z^* \in \bigwedge^{p+q} E^*,$$

the right hook

$$\lrcorner : \bigwedge^{p+q} E^* \times \bigwedge^p E \longrightarrow \bigwedge^q E^*$$

by

$$\langle z^* \lrcorner u, v \rangle = \langle z^*, u \wedge v \rangle, \quad \text{for all } u \in \bigwedge^p E, v \in \bigwedge^q E, \text{ and } z^* \in \bigwedge^{p+q} E^*,$$

and $v \wedge u = (-1)^{pq} u \wedge v$, we conclude that

$$u \lrcorner z^* = (-1)^{pq} z^* \lrcorner u,$$

where $u \in \bigwedge^p E$ and $z \in \bigwedge^{p+q} E^*$.

Using the above property and Proposition 23.27, we get the following version of Proposition 23.27 for the right hook:

Proposition 23.28. *For the right hook*

$$\lrcorner : \bigwedge^{p+q} E^* \times \bigwedge^p E \longrightarrow \bigwedge^q E^*,$$

for every $u \in E$, we have

$$(x^* \wedge y^*) \lrcorner u = (x^* \lrcorner u) \wedge y^* + (-1)^r x^* \wedge (y^* \lrcorner u),$$

where $x^* \in \bigwedge^r E^*$.

Thus, \lrcorner is an anti-derivation.

For $u \in E$, the right hook $z^* \lrcorner u$ is also denoted $i(u)z^*$, and called *insertion operator* or *interior product*. This operator plays an important role in differential geometry. If we view $z^* \in \bigwedge^{n+1}(E^*)$ as an alternating multilinear map in $\text{Alt}^{n+1}(E; K)$, then $i(u)z^* \in \text{Alt}^n(E; K)$ is given by

$$(i(u)z^*)(v_1, \dots, v_n) = z^*(u, v_1, \dots, v_n).$$



Note that certain authors, such as Shafarevitch [98], denote our right hook $z^* \lrcorner u$ (which is also the right hook in Bourbaki [14] and Fulton and Harris [42]) by $u \lrcorner z^*$.

Using the two versions of \lrcorner , we can define linear maps $\gamma: \bigwedge^p E \rightarrow \bigwedge^{n-p} E^*$ and $\delta: \bigwedge^p E^* \rightarrow \bigwedge^{n-p} E$. For any basis (e_1, \dots, e_n) of E , if we let $M = \{1, \dots, n\}$, $e = e_1 \wedge \dots \wedge e_n$, and $e^* = e_1^* \wedge \dots \wedge e_n^*$, then

$$\gamma(u) = u \lrcorner e^* \quad \text{and} \quad \delta(v) = v^* \lrcorner e,$$

for all $u \in \bigwedge^p E$ and all $v^* \in \bigwedge^p E^*$. The following proposition is easily shown.

Proposition 23.29. *The linear maps $\gamma: \bigwedge^p E \rightarrow \bigwedge^{n-p} E^*$ and $\delta: \bigwedge^p E^* \rightarrow \bigwedge^{n-p} E$ are isomorphisms. The isomorphisms γ and δ map decomposable vectors to decomposable vectors. Furthermore, if $z \in \bigwedge^p E$ is decomposable, then $\langle \gamma(z), z \rangle = 0$, and similarly for $z \in \bigwedge^p E^*$. If (e'_1, \dots, e'_n) is any other basis of E and $\gamma': \bigwedge^p E \rightarrow \bigwedge^{n-p} E^*$ and $\delta': \bigwedge^p E^* \rightarrow \bigwedge^{n-p} E$ are the corresponding isomorphisms, then $\gamma' = \lambda \gamma$ and $\delta' = \lambda^{-1} \delta$ for some nonzero $\lambda \in \Omega$.*

Proof. Using Proposition 23.26, for any subset $J \subseteq \{1, \dots, n\} = M$ such that $|J| = p$, we have

$$\gamma(e_J) = e_J \lrcorner e^* = \rho_{M-J, J} e_{M-J}^* \quad \text{and} \quad \delta(e_J^*) = e_J^* \lrcorner e = \rho_{M-J, J} e_{M-J}.$$

Thus,

$$\delta \circ \gamma(e_J) = \rho_{M-J, J} \rho_{J, M-J} e_J = (-1)^{p(n-p)} e_J.$$

A similar result holds for $\gamma \circ \delta$. This implies that

$$\delta \circ \gamma = (-1)^{p(n-p)} \text{id} \quad \text{and} \quad \gamma \circ \delta = (-1)^{p(n-p)} \text{id}.$$

Thus, γ and δ are isomorphisms. If $z \in \bigwedge^p E$ is decomposable, then $z = u_1 \wedge \dots \wedge u_p$ where u_1, \dots, u_p are linearly independent since $z \neq 0$, and we can pick a basis of E of the form (u_1, \dots, u_n) . Then, the above formulae show that

$$\gamma(z) = \pm u_{p+1}^* \wedge \dots \wedge u_n^*.$$

Clearly

$$\langle \gamma(z), z \rangle = 0.$$

If (e'_1, \dots, e'_n) is any other basis of E , because $\bigwedge^m E$ has dimension 1, we have

$$e'_1 \wedge \dots \wedge e'_n = \lambda e_1 \wedge \dots \wedge e_n$$

for some nonnull $\lambda \in \Omega$, and the rest is trivial. \square

We are now ready to tackle the problem of finding criteria for decomposability. We need a few preliminary results.

Proposition 23.30. *Given $z \in \bigwedge^p E$ with $z \neq 0$, the smallest vector space $W \subseteq E$ such that $z \in \bigwedge^p W$ is generated by the vectors of the form*

$$u^* \lrcorner z, \quad \text{with } u^* \in \bigwedge^{p-1} E^*.$$

Proof. First, let W be any subspace such that $z \in \bigwedge^p(E)$ and let $(e_1, \dots, e_r, e_{r+1}, \dots, e_n)$ be a basis of E such that (e_1, \dots, e_r) is a basis of W . Then, $u^* = \sum_I e_I^*$, where $I \subseteq \{1, \dots, n\}$ and $|I| = p-1$, and $z = \sum_J e_J$, where $J \subseteq \{1, \dots, r\}$ and $|J| = p \leq r$. It follows immediately from the formula of Proposition 23.26 (3) that $u^* \lrcorner z \in W$.

Next, we prove that if W is the smallest subspace of E such that $z \in \bigwedge^p(W)$, then W is generated by the vectors of the form $u^* \lrcorner z$, where $u^* \in \bigwedge^{p-1} E^*$. Suppose not, then the vectors $u^* \lrcorner z$ with $u^* \in \bigwedge^{p-1} E^*$ span a proper subspace U of W . We prove that for every subspace W' of W with $\dim(W') = \dim(W) - 1 = r - 1$, it is not possible that $u^* \lrcorner z \in W'$ for all $u^* \in \bigwedge^{p-1} E^*$. But then, as U is a proper subspace of W , it is contained in some subspace W' with $\dim(W') = r - 1$, and we have a contradiction.

Let $w \in W - W'$ and pick a basis of W formed by a basis (e_1, \dots, e_{r-1}) of W' and w . We can write $z = z' + w \wedge z''$, where $z' \in \bigwedge^p W'$ and $z'' \in \bigwedge^{p-1} W'$, and since W is the smallest subspace containing z , we have $z'' \neq 0$. Consequently, if we write $z'' = \sum_I e_I$ in terms of the basis (e_1, \dots, e_{r-1}) of W' , there is some e_I , with $I \subseteq \{1, \dots, r-1\}$ and $|I| = p-1$, so that the coefficient λ_I is nonzero. Now, using any basis of E containing (e_1, \dots, e_{r-1}, w) , by Proposition 23.26 (3), we see that

$$e_I^* \lrcorner (w \wedge e_I) = \lambda w, \quad \lambda = \pm 1.$$

It follows that

$$e_I^* \lrcorner z = e_I^* \lrcorner (z' + w \wedge z'') = e_I^* \lrcorner z' + e_I^* \lrcorner (w \wedge z'') = e_I^* \lrcorner z' + \lambda w,$$

with $e_I^* \lrcorner z' \in W'$, which shows that $e_I^* \lrcorner z \notin W'$. Therefore, W is indeed generated by the vectors of the form $u^* \lrcorner z$, where $u^* \in \bigwedge^{p-1} E^*$. \square

Proposition 23.31. *Any nonzero $z \in \bigwedge^p E$ is decomposable iff*

$$(u^* \lrcorner z) \wedge z = 0, \quad \text{for all } u^* \in \bigwedge^{p-1} E^*.$$

Proof. Clearly, $z \in \bigwedge^p E$ is decomposable iff the smallest vector space W such that $z \in \bigwedge^p W$ has dimension p . If $\dim(W) = p$, we have $z = e_1 \wedge \dots \wedge e_p$ where e_1, \dots, e_p form a basis of W . By Proposition 23.30, for every $u^* \in \bigwedge^{p-1} E^*$, we have $u^* \lrcorner z \in W$, so each $u^* \lrcorner z$ is a linear combination of the e_i 's and $(u^* \lrcorner z) \wedge z = (u^* \lrcorner z) \wedge e_1 \wedge \dots \wedge e_p = 0$.

Now, assume that $(u^* \lrcorner z) \wedge z = 0$ for all $u^* \in \bigwedge^{p-1} E^*$, and that $\dim(W) = n > p$. If e_1, \dots, e_n is a basis of W , then we have $z = \sum_I \lambda_I e_I$, where $I \subseteq \{1, \dots, n\}$ and $|I| = p$. Recall that $z \neq 0$, and so, some λ_I is nonzero. By Proposition 23.30, each e_i can be written as $u^* \lrcorner z$ for some $u^* \in \bigwedge^{p-1} E^*$, and since $(u^* \lrcorner z) \wedge z = 0$ for all $u^* \in \bigwedge^{p-1} E^*$, we get

$$e_j \wedge z = 0 \quad \text{for } j = 1, \dots, n.$$

By wedging $z = \sum_I \lambda_I e_I$ with each e_j , as $n > p$, we deduce $\lambda_I = 0$ for all I , so $z = 0$, a contradiction. Therefore, $n = p$ and z is decomposable. \square

In Proposition 23.31, we can let u^* range over a basis of $\bigwedge^{p-1} E^*$, and then the conditions are

$$(e_H^* \lrcorner z) \wedge z = 0$$

for all $H \subseteq \{1, \dots, n\}$, with $|H| = p - 1$. Since $(e_H^* \lrcorner z) \wedge z \in \bigwedge^{p+1} E$, this is equivalent to

$$e_J^*((e_H^* \lrcorner z) \wedge z) = 0$$

for all $H, J \subseteq \{1, \dots, n\}$, with $|H| = p - 1$ and $|J| = p + 1$. Then, for all $I, I' \subseteq \{1, \dots, n\}$ with $|I| = |I'| = p$, we can show that

$$e_J^*((e_H^* \lrcorner e_I) \wedge e_{I'}) = 0,$$

unless there is some $i \in \{1, \dots, n\}$ such that

$$I - H = \{i\}, \quad J - I' = \{i\}.$$

In this case,

$$e_J^*((e_H^* \lrcorner e_{H \cup \{i\}}) \wedge e_{J - \{i\}}) = \rho_{\{i\}, H} \rho_{\{i\}, J - \{i\}}.$$

If we let

$$\epsilon_{i, J, H} = \rho_{\{i\}, H} \rho_{\{i\}, J - \{i\}},$$

we have $\epsilon_{i, J, H} = +1$ if the parity of the number of $j \in J$ such that $j < i$ is the same as the parity of the number of $h \in H$ such that $h < i$, and $\epsilon_{i, J, H} = -1$ otherwise.

Finally, we obtain the following criterion in terms of quadratic equations (*Plücker's equations*) for the decomposability of an alternating tensor:

Proposition 23.32. (*Grassmann-Plücker's Equations*) For $z = \sum_I \lambda_I e_I \in \bigwedge^p E$, the conditions for $z \neq 0$ to be decomposable are

$$\sum_{i \in J - H} \epsilon_{i, J, H} \lambda_{H \cup \{i\}} \lambda_{J - \{i\}} = 0,$$

for all $H, J \subseteq \{1, \dots, n\}$ such that $|H| = p - 1$ and $|J| = p + 1$.

Using these criteria, it is a good exercise to prove that if $\dim(E) = n$, then every tensor in $\bigwedge^{n-1}(E)$ is decomposable. This can also be shown directly.

It should be noted that the equations given by Proposition 23.32 are not independent. For example, when $\dim(E) = n = 4$ and $p = 2$, these equations reduce to the single equation

$$\lambda_{12}\lambda_{34} - \lambda_{13}\lambda_{24} + \lambda_{14}\lambda_{23} = 0.$$

When the field K is the field of complex numbers, this is the homogeneous equation of a quadric in \mathbb{CP}^5 known as the *Klein quadric*. The points on this quadric are in one-to-one correspondence with the lines in \mathbb{CP}^3 .

23.18 Vector-Valued Alternating Forms

In this section, the vector space E is assumed to have finite dimension. We know that there is a canonical isomorphism $\bigwedge^n(E^*) \cong \text{Alt}^n(E; K)$ between alternating n -forms and alternating multilinear maps. As in the case of general tensors, the isomorphisms

$$\begin{aligned}\text{Alt}^n(E; F) &\cong \text{Hom}(\bigwedge^n(E), F) \\ \text{Hom}(\bigwedge^n(E), F) &\cong (\bigwedge^n(E))^* \otimes F \\ (\bigwedge^n(E))^* &\cong \bigwedge^n(E^*)\end{aligned}$$

yield a canonical isomorphism

$$\text{Alt}^n(E; F) \cong \left(\bigwedge^n(E^*) \right) \otimes F.$$

Note that F may have infinite dimension. This isomorphism allows us to view the tensors in $\bigwedge^n(E^*) \otimes F$ as *vector valued alternating forms*, a point of view that is useful in differential geometry. If (f_1, \dots, f_r) is a basis of F , every tensor $\omega \in \bigwedge^n(E^*) \otimes F$ can be written as some linear combination

$$\omega = \sum_{i=1}^r \alpha_i \otimes f_i,$$

with $\alpha_i \in \bigwedge^n(E^*)$. We also let

$$\bigwedge(E; F) = \bigoplus_{n=0} \left(\bigwedge^n(E^*) \right) \otimes F = \left(\bigwedge(E) \right) \otimes F.$$

Given three vector spaces, F, G, H , if we have some bilinear map $\Phi: F \otimes G \rightarrow H$, then we can define a multiplication operation

$$\wedge_\Phi: \bigwedge(E; F) \times \bigwedge(E; G) \rightarrow \bigwedge(E; H)$$

as follows: For every pair (m, n) , we define the multiplication

$$\wedge_\Phi: \left(\left(\bigwedge^m(E^*) \right) \otimes F \right) \times \left(\left(\bigwedge^n(E^*) \right) \otimes G \right) \longrightarrow \left(\bigwedge^{m+n}(E^*) \right) \otimes H$$

by

$$(\alpha \otimes f) \wedge_\Phi (\beta \otimes g) = (\alpha \wedge \beta) \otimes \Phi(f, g).$$

As in Section 23.15 (following H. Cartan [21]) we can also define a multiplication

$$\wedge_\Phi: \text{Alt}^m(E; F) \times \text{Alt}^n(E; G) \longrightarrow \text{Alt}^{m+n}(E; H)$$

directly on alternating multilinear maps as follows: For $f \in \text{Alt}^m(E; F)$ and $g \in \text{Alt}^n(E; G)$,

$$(f \wedge_{\Phi} g)(u_1, \dots, u_{m+n}) = \sum_{\sigma \in \text{shuffle}(m, n)} \text{sgn}(\sigma) \Phi(f(u_{\sigma(1)}, \dots, u_{\sigma(m)}), g(u_{\sigma(m+1)}, \dots, u_{\sigma(m+n)})),$$

where $\text{shuffle}(m, n)$ consists of all (m, n) -“shuffles,” that is, permutations σ of $\{1, \dots, m+n\}$ such that $\sigma(1) < \dots < \sigma(m)$ and $\sigma(m+1) < \dots < \sigma(m+n)$.

In general, not much can be said about \wedge_{Φ} , unless Φ has some additional properties. In particular, \wedge_{Φ} is generally not associative. We also have the map

$$\mu: \left(\bigwedge^n (E^*) \right) \otimes F \longrightarrow \text{Alt}^n(E; F)$$

defined on generators by

$$\mu((v_1^* \wedge \dots \wedge v_n^*) \otimes a)(u_1, \dots, u_n) = (\det(v_j^*(u_i)))a.$$

Proposition 23.33. *The map*

$$\mu: \left(\bigwedge^n (E^*) \right) \otimes F \longrightarrow \text{Alt}^n(E; F)$$

defined as above is a canonical isomorphism for every $n \geq 0$. Furthermore, given any three vector spaces, F, G, H , and any bilinear map $\Phi: F \times G \rightarrow H$, for all $\omega \in (\bigwedge^n (E^)) \otimes F$ and all $\eta \in (\bigwedge^n (E^*)) \otimes G$,*

$$\mu(\alpha \wedge_{\Phi} \beta) = \mu(\alpha) \wedge_{\Phi} \mu(\beta).$$

Proof. Since we already know that $(\bigwedge^n (E^*)) \otimes F$ and $\text{Alt}^n(E; F)$ are isomorphic, it is enough to show that μ maps some basis of $(\bigwedge^n (E^*)) \otimes F$ to linearly independent elements. Pick some bases (e_1, \dots, e_p) in E and $(f_j)_{j \in J}$ in F . Then, we know that the vectors $e_I^* \otimes f_j$, where $I \subseteq \{1, \dots, p\}$ and $|I| = n$, form a basis of $(\bigwedge^n (E^*)) \otimes F$. If we have a linear dependence

$$\sum_{I, j} \lambda_{I, j} \mu(e_I^* \otimes f_j) = 0,$$

applying the above combination to each $(e_{i_1}, \dots, e_{i_n})$ ($I = \{i_1, \dots, i_n\}$, $i_1 < \dots < i_n$), we get the linear combination

$$\sum_j \lambda_{I, j} f_j = 0,$$

and by linear independence of the f_j 's, we get $\lambda_{I, j} = 0$ for all I and all j . Therefore, the $\mu(e_I^* \otimes f_j)$ are linearly independent, and we are done. The second part of the proposition is easily checked (a simple computation). \square

A special case of interest is the case where $F = G = H$ is a Lie algebra and $\Phi(a, b) = [a, b]$ is the Lie bracket of F . In this case, using a basis (f_1, \dots, f_r) of F , if we write $\omega = \sum_i \alpha_i \otimes f_i$ and $\eta = \sum_j \beta_j \otimes f_j$, we have

$$[\omega, \eta] = \sum_{i,j} \alpha_i \wedge \beta_j \otimes [f_i, f_j].$$

Consequently,

$$[\eta, \omega] = (-1)^{mn+1} [\omega, \eta].$$

The following proposition will be useful in dealing with vector-valued differential forms:

Proposition 23.34. *If (e_1, \dots, e_p) is any basis of E , then every element $\omega \in (\bigwedge^n(E^*)) \otimes F$ can be written in a unique way as*

$$\omega = \sum_I e_I^* \otimes f_I, \quad f_I \in F,$$

where the e_I^* are defined as in Section 23.12.

Proof. Since, by Proposition 23.19, the e_I^* form a basis of $\bigwedge^n(E^*)$, elements of the form $e_I^* \otimes f$ span $(\bigwedge^n(E^*)) \otimes F$. Now, if we apply $\mu(\omega)$ to $(e_{i_1}, \dots, e_{i_n})$, where $I = \{i_1, \dots, i_n\} \subseteq \{1, \dots, p\}$, we get

$$\mu(\omega)(e_{i_1}, \dots, e_{i_n}) = \mu(e_I^* \otimes f_I)(e_{i_1}, \dots, e_{i_n}) = f_I.$$

Therefore, the f_I are uniquely determined by ω . □

Proposition can also be formulated in terms of alternating multilinear maps, a fact that will be useful to deal with differential forms.

Define the product, $\cdot: \text{Alt}^n(E; \mathbb{R}) \times F \rightarrow \text{Alt}^n(E; F)$, as follows: For all $\omega \in \text{Alt}^n(E; \mathbb{R})$ and all $f \in F$,

$$(\omega \cdot f)(u_1, \dots, u_n) = \omega(u_1, \dots, u_n)f,$$

for all $u_1, \dots, u_n \in E$. Then, it is immediately verified that for every $\omega \in (\bigwedge^n(E^*)) \otimes F$ of the form

$$\omega = u_1^* \wedge \dots \wedge u_n^* \otimes f,$$

we have

$$\mu(u_1^* \wedge \dots \wedge u_n^* \otimes f) = \mu(u_1^* \wedge \dots \wedge u_n^*) \cdot f.$$

Then, Proposition 23.34 yields

Proposition 23.35. *If (e_1, \dots, e_p) is any basis of E , then every element $\omega \in \text{Alt}^n(E; F)$ can be written in a unique way as*

$$\omega = \sum_I e_I^* \cdot f_I, \quad f_I \in F,$$

where the e_I^* are defined as in Section 23.12.

23.19 The Pfaffian Polynomial

Let $\mathfrak{so}(2n)$ denote the vector space (actually, Lie algebra) of $2n \times 2n$ real skew-symmetric matrices. It is well-known that every matrix $A \in \mathfrak{so}(2n)$ can be written as

$$A = PDP^\top,$$

where P is an orthogonal matrix and where D is a block diagonal matrix

$$D = \begin{pmatrix} D_1 & & & \\ & D_2 & & \\ & & \ddots & \\ & & & D_n \end{pmatrix}$$

consisting of 2×2 blocks of the form

$$D_i = \begin{pmatrix} 0 & -a_i \\ a_i & 0 \end{pmatrix}.$$

For a proof, see Horn and Johnson [57] (Corollary 2.5.14), Gantmacher [46] (Chapter IX), or Gallier [44] (Chapter 11).

Since $\det(D_i) = a_i^2$ and $\det(A) = \det(PDP^\top) = \det(D) = \det(D_1) \cdots \det(D_n)$, we get

$$\det(A) = (a_1 \cdots a_n)^2.$$

The Pfaffian is a polynomial function $\text{Pf}(A)$ in skew-symmetric $2n \times 2n$ matrices A (a polynomial in $(2n-1)n$ variables) such that

$$\text{Pf}(A)^2 = \det(A),$$

and for every arbitrary matrix B ,

$$\text{Pf}(BAB^\top) = \text{Pf}(A) \det(B).$$

The Pfaffian shows up in the definition of the Euler class of a vector bundle. There is a simple way to define the Pfaffian using some exterior algebra. Let (e_1, \dots, e_{2n}) be any basis of \mathbb{R}^{2n} . For any matrix $A \in \mathfrak{so}(2n)$, let

$$\omega(A) = \sum_{i < j} a_{ij} e_i \wedge e_j,$$

where $A = (a_{ij})$. Then, $\bigwedge^n \omega(A)$ is of the form $C e_1 \wedge e_2 \wedge \cdots \wedge e_{2n}$ for some constant $C \in \mathbb{R}$.

Definition 23.9. For every skew symmetric matrix $A \in \mathfrak{so}(2n)$, the *Pfaffian polynomial* or *Pfaffian*, is the degree n polynomial $\text{Pf}(A)$ defined by

$$\bigwedge^n \omega(A) = n! \text{Pf}(A) e_1 \wedge e_2 \wedge \cdots \wedge e_{2n}.$$

Clearly, $\text{Pf}(A)$ is independent of the basis chosen. If A is the block diagonal matrix D , a simple calculation shows that

$$\omega(D) = -(a_1 e_1 \wedge e_2 + a_2 e_3 \wedge e_4 + \cdots + a_n e_{2n-1} \wedge e_{2n})$$

and that

$$\bigwedge^n \omega(D) = (-1)^n n! a_1 \cdots a_n e_1 \wedge e_2 \wedge \cdots \wedge e_{2n},$$

and so

$$\text{Pf}(D) = (-1)^n a_1 \cdots a_n.$$

Since $\text{Pf}(D)^2 = (a_1 \cdots a_n)^2 = \det(A)$, we seem to be on the right track.

Proposition 23.36. *For every skew symmetric matrix $A \in \mathfrak{so}(2n)$ and every arbitrary matrix B , we have:*

$$(i) \text{Pf}(A)^2 = \det(A)$$

$$(ii) \text{Pf}(BAB^\top) = \text{Pf}(A) \det(B).$$

Proof. If we assume that (ii) is proved then, since we can write $A = PDP^\top$ for some orthogonal matrix P and some block diagonal matrix D as above, as $\det(P) = \pm 1$ and $\text{Pf}(D)^2 = \det(A)$, we get

$$\text{Pf}(A)^2 = \text{Pf}(PDP^\top)^2 = \text{Pf}(D)^2 \det(P)^2 = \det(A),$$

which is (i). Therefore, it remains to prove (ii).

Let $f_i = B e_i$ for $i = 1, \dots, 2n$, where (e_1, \dots, e_{2n}) is any basis of \mathbb{R}^{2n} . Since $f_i = \sum_k b_{ki} e_k$, we have

$$\tau = \sum_{i,j} a_{ij} f_i \wedge f_j = \sum_{i,j} \sum_{k,l} b_{ki} a_{ij} b_{lj} e_k \wedge e_l = \sum_{k,l} (BAB^\top)_{kl} e_k \wedge e_l,$$

and so, as BAB^\top is skew symmetric and $e_k \wedge e_l = -e_l \wedge e_k$, we get

$$\tau = 2\omega(BAB^\top).$$

Consequently,

$$\bigwedge^n \tau = 2^n \bigwedge^n \omega(BAB^\top) = 2^n n! \text{Pf}(BAB^\top) e_1 \wedge e_2 \wedge \cdots \wedge e_{2n}.$$

Now,

$$\bigwedge^n \tau = C f_1 \wedge f_2 \wedge \cdots \wedge f_{2n},$$

for some $C \in \mathbb{R}$. If B is singular, then the f_i are linearly dependent, which implies that $f_1 \wedge f_2 \wedge \cdots \wedge f_{2n} = 0$, in which case

$$\text{Pf}(BAB^\top) = 0,$$

as $e_1 \wedge e_2 \wedge \cdots \wedge e_{2n} \neq 0$. Therefore, if B is singular, $\det(B) = 0$ and

$$\text{Pf}(BAB^\top) = 0 = \text{Pf}(A) \det(B).$$

If B is invertible, as $\tau = \sum_{i,j} a_{ij} f_i \wedge f_j = 2 \sum_{i < j} a_{ij} f_i \wedge f_j$, we have

$$\bigwedge^n \tau = 2^n n! \text{Pf}(A) f_1 \wedge f_2 \wedge \cdots \wedge f_{2n}.$$

However, as $f_i = B e_i$, we have

$$f_1 \wedge f_2 \wedge \cdots \wedge f_{2n} = \det(B) e_1 \wedge e_2 \wedge \cdots \wedge e_{2n},$$

so

$$\bigwedge^n \tau = 2^n n! \text{Pf}(A) \det(B) e_1 \wedge e_2 \wedge \cdots \wedge e_{2n}$$

and as

$$\bigwedge^n \tau = 2^n n! \text{Pf}(BAB^\top) e_1 \wedge e_2 \wedge \cdots \wedge e_{2n},$$

we get

$$\text{Pf}(BAB^\top) = \text{Pf}(A) \det(B),$$

as claimed. □

Remark: It can be shown that the polynomial $\text{Pf}(A)$ is the unique polynomial with integer coefficients such that $\text{Pf}(A)^2 = \det(A)$ and $\text{Pf}(\text{diag}(S, \dots, S)) = +1$, where

$$S = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix};$$

see Milnor and Stasheff [82] (Appendix C, Lemma 9). There is also an explicit formula for $\text{Pf}(A)$, namely:

$$\text{Pf}(A) = \frac{1}{2^n n!} \sum_{\sigma \in \mathfrak{S}_{2n}} \text{sgn}(\sigma) \prod_{i=1}^n a_{\sigma(2i-1) \sigma(2i)}.$$



Beware, some authors use a different sign convention and require the Pfaffian to have the value $+1$ on the matrix $\text{diag}(S', \dots, S')$, where

$$S' = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}.$$

For example, if \mathbb{R}^{2n} is equipped with an inner product $\langle -, - \rangle$, then some authors define $\omega(A)$ as

$$\omega(A) = \sum_{i < j} \langle A e_i, e_j \rangle e_i \wedge e_j,$$

where $A = (a_{ij})$. But then, $\langle A e_i, e_j \rangle = a_{ji}$ and **not** a_{ij} , and this Pfaffian takes the value $+1$ on the matrix $\text{diag}(S', \dots, S')$. This version of the Pfaffian differs from our version by the factor $(-1)^n$. In this respect, Madsen and Tornehave [74] seem to have an incorrect sign in Proposition B6 of Appendix C.

We will also need another property of Pfaffians. Recall that the ring $M_n(\mathbb{C})$ of $n \times n$ matrices over \mathbb{C} is embedded in the ring $M_{2n}(\mathbb{R})$ of $2n \times 2n$ matrices with real coefficients, using the injective homomorphism that maps every entry $z = a + ib \in \mathbb{C}$ to the 2×2 matrix

$$\begin{pmatrix} a & -b \\ b & a \end{pmatrix}.$$

If $A \in M_n(\mathbb{C})$, let $A_{\mathbb{R}} \in M_{2n}(\mathbb{R})$ denote the real matrix obtained by the above process. Observe that every skew Hermitian matrix $A \in \mathfrak{u}(n)$ (i.e., with $A^* = \overline{A}^{\top} = -A$) yields a matrix $A_{\mathbb{R}} \in \mathfrak{so}(2n)$.

Proposition 23.37. *For every skew Hermitian matrix $A \in \mathfrak{u}(n)$, we have*

$$\text{Pf}(A_{\mathbb{R}}) = i^n \det(A).$$

Proof. It is well-known that a skew Hermitian matrix can be diagonalized with respect to a unitary matrix U and that the eigenvalues are pure imaginary or zero, so we can write

$$A = U \text{diag}(ia_1, \dots, ia_n) U^*,$$

for some reals $a_j \in \mathbb{R}$. Consequently, we get

$$A_{\mathbb{R}} = U_{\mathbb{R}} \text{diag}(D_1, \dots, D_n) U_{\mathbb{R}}^{\top},$$

where

$$D_j = \begin{pmatrix} 0 & -a_j \\ a_j & 0 \end{pmatrix}$$

and

$$\text{Pf}(A_{\mathbb{R}}) = \text{Pf}(\text{diag}(D_1, \dots, D_n)) = (-1)^n a_1 \cdots a_n,$$

as we saw before. On the other hand,

$$\det(A) = \det(\text{diag}(ia_1, \dots, ia_n)) = i^n a_1 \cdots a_n,$$

and as $(-1)^n = i^n i^n$, we get

$$\text{Pf}(A_{\mathbb{R}}) = i^n \det(A),$$

as claimed. \square

\square



Madsen and Tornehave [74] state Proposition 23.37 using the factor $(-i)^n$, which is wrong.

Chapter 24

Introduction to Modules; Modules over a PID

24.1 Modules over a Commutative Ring

In this chapter, we introduce modules over a commutative ring (with unity). After a quick overview of fundamental concepts such as free modules, torsion modules, and some basic results about them, we focus on finitely generated modules over a PID and we prove the structure theorems for this class of modules (invariant factors and elementary divisors). Our main goal is not to give a comprehensive exposition of modules, but instead to apply the structure theorem to the $K[X]$ -module E_f defined by a linear map f acting on a finite-dimensional vector space E , and to obtain several normal forms for f , including the rational canonical form.

A module is the generalization of a vector space E over a field K obtained replacing the field K by a commutative ring A (with unity 1). Although formally, the definition is the same, the fact that some nonzero elements of A are not invertible has some serious consequences. For example, it is possible that $\lambda \cdot u = 0$ for some nonzero $\lambda \in A$ and some nonzero $u \in E$, and a module may no longer have a basis.

For the sake of completeness, we give the definition of a module, although it is the same as Definition 2.9 with the field K replaced by a ring A . In this chapter, *all rings under consideration are assumed to be commutative and to have an identity element 1.*

Definition 24.1. Given a ring A , a (*left*) *module over A* (or *A -module*) is a set M (of vectors) together with two operations $+: M \times M \rightarrow M$ (called *vector addition*),¹ and $\cdot: A \times M \rightarrow M$ (called *scalar multiplication*) satisfying the following conditions for all $\alpha, \beta \in A$ and all $u, v \in M$;

(M0) M is an abelian group w.r.t. $+$, with identity element 0;

¹The symbol $+$ is overloaded, since it denotes both addition in the ring A and addition of vectors in M . It is usually clear from the context which $+$ is intended.

$$(M1) \quad \alpha \cdot (u + v) = (\alpha \cdot u) + (\alpha \cdot v);$$

$$(M2) \quad (\alpha + \beta) \cdot u = (\alpha \cdot u) + (\beta \cdot u);$$

$$(M3) \quad (\alpha * \beta) \cdot u = \alpha \cdot (\beta \cdot u);$$

$$(M4) \quad 1 \cdot u = u.$$

Given $\alpha \in A$ and $v \in M$, the element $\alpha \cdot v$ is also denoted by αv . The ring A is often called the ring of scalars.

Unless specified otherwise or unless we are dealing with several different rings, in the rest of this chapter, we assume that all A -modules are defined with respect to a fixed ring A . Thus, we will refer to a A -module simply as a module.

From (M0), a module always contains the null vector 0, and thus is nonempty. From (M1), we get $\alpha \cdot 0 = 0$, and $\alpha \cdot (-v) = -(\alpha \cdot v)$. From (M2), we get $0 \cdot v = 0$, and $(-\alpha) \cdot v = -(\alpha \cdot v)$. The ring A itself can be viewed as a module over itself, addition of vectors being addition in the ring, and multiplication by a scalar being multiplication in the ring.

When the ring A is a field, an A -module is a vector space. When $A = \mathbb{Z}$, a \mathbb{Z} -module is just an abelian group, with the action given by

$$\begin{aligned} 0 \cdot u &= 0, \\ n \cdot u &= \underbrace{u + \cdots + u}_n, & n > 0 \\ n \cdot u &= -(-n) \cdot u, & n < 0. \end{aligned}$$

All definitions from Section 2.3, linear combinations, linear independence and linear dependence, subspaces renamed as *submodules*, apply unchanged to modules. Proposition 2.7 also holds for the module spanned by a set of vectors. The definition of a basis (Definition 2.12) also applies to modules, but the only result from Section 2.4 that holds for modules is Proposition 2.14. Unfortunately, it is longer true that every module has a basis. For example, for any nonzero integer $m \in \mathbb{Z}$, the \mathbb{Z} -module $\mathbb{Z}/m\mathbb{Z}$ has no basis. Similarly, \mathbb{Q} , as a \mathbb{Z} -module, has no basis. In fact, any two distinct nonzero elements p_1/q_1 and p_2/q_2 are linearly dependent, since

$$(p_2 q_1) \left(\frac{p_1}{q_1} \right) - (p_1 q_2) \left(\frac{p_2}{q_2} \right) = 0.$$

Definition 2.13 can be generalized to rings and yields free modules.

Definition 24.2. Given a commutative ring A and any (nonempty) set I , let $A^{(I)}$ be the subset of the cartesian product A^I consisting of all families $(\lambda_i)_{i \in I}$ with finite support of scalars in A .² We define addition and multiplication by a scalar as follows:

$$(\lambda_i)_{i \in I} + (\mu_i)_{i \in I} = (\lambda_i + \mu_i)_{i \in I},$$

²Where A^I denotes the set of all functions from I to A .

and

$$\lambda \cdot (\mu_i)_{i \in I} = (\lambda \mu_i)_{i \in I}.$$

It is immediately verified that addition and multiplication by a scalar are well defined. Thus, $A^{(I)}$ is a module. Furthermore, because families with finite support are considered, the family $(e_i)_{i \in I}$ of vectors e_i , defined such that $(e_i)_j = 0$ if $j \neq i$ and $(e_i)_i = 1$, is clearly a basis of the module $A^{(I)}$. When $I = \{1, \dots, n\}$, we denote $A^{(I)}$ by A^n . The function $\iota: I \rightarrow A^{(I)}$, such that $\iota(i) = e_i$ for every $i \in I$, is clearly an injection.

Definition 24.3. An A -module M is *free* iff it has a basis.

The module $A^{(I)}$ is a free module.

All definitions from Section 2.5 apply to modules, linear maps, kernel, image, except the definition of rank, which has to be defined differently. Propositions 2.15, 2.16, 2.17, and 2.18 hold for modules. However, the other propositions do not generalize to modules. The definition of an isomorphism generalizes to modules. As a consequence, a module is free iff it is isomorphic to a module of the form $A^{(I)}$.

Section 2.6 generalizes to modules. Given a submodule N of a module M , we can define the quotient module M/N .

If \mathfrak{a} is an ideal in A and if M is an A -module, we define $\mathfrak{a}M$ as the set of finite sums of the form

$$a_1 m_1 + \dots + a_k m_k, \quad a_i \in \mathfrak{a}, m_i \in M.$$

It is immediately verified that $\mathfrak{a}M$ is a submodule of M .

Interestingly, the part of Theorem 2.13 that asserts that any two bases of a vector space have the same cardinality holds for modules. One way to prove this fact is to “pass” to a vector space by a quotient process.

Theorem 24.1. *For any free module M , any two bases of M have the same cardinality.*

Proof sketch. We give the argument for finite bases, but it also holds for infinite bases. The trick is to pick any maximal ideal \mathfrak{m} in A (whose existence is guaranteed by Theorem 31.3). Then, A/\mathfrak{m} is a field, and $M/\mathfrak{m}M$ can be made into a vector space over A/\mathfrak{m} ; we leave the details as an exercise. If (u_1, \dots, u_n) is a basis of M , then it is easy to see that the image of this basis is a basis of the vector space $M/\mathfrak{m}M$. By Theorem 2.13, the number n of elements in any basis of $M/\mathfrak{m}M$ is an invariant, so any two bases of M must have the same number of elements. \square

The common number of elements in any basis of a free module is called the *dimension* (or *rank*) of the free module.

One should realize that the notion of linear independence in a module is a little tricky. According to the definition, the one-element sequence (u) consisting of a single nonzero

vector is linearly independent if for all $\lambda \in A$, if $\lambda u = 0$ then $\lambda = 0$. However, there are free modules that contain nonzero vectors that are not linearly independent! For example, the ring $A = \mathbb{Z}/6\mathbb{Z}$ viewed as a module over itself has the basis (1) , but the zero-divisors, such as 2 or 4, are not linearly independent. Using language introduced in Definition 24.4, a free module may have torsion elements. There are also nonfree modules such that every nonzero vector is linearly independent, such as \mathbb{Q} over \mathbb{Z} .

All definitions from Section 3.1 about matrices apply to free modules, and so do all the proposition. Similarly, all definitions from Section 4.1 about direct sums and direct products apply to modules. All propositions that do not involve extending bases still hold. The important proposition 4.10 survives in the following form.

Proposition 24.2. *Let $f: E \rightarrow F$ be a surjective linear between two A -modules with F a free module. Given any basis (v_1, \dots, v_r) of F , for any r vectors $u_1, \dots, u_r \in E$ such that $f(u_i) = v_i$ for $i = 1, \dots, r$, the vectors (u_1, \dots, u_r) are linearly independent and the module E is the direct sum*

$$E = \text{Ker}(f) \oplus U,$$

where U is the free submodule of E spanned by the basis (u_1, \dots, u_r) .

Proof. Pick any $w \in E$, write $f(w)$ over the basis (v_1, \dots, v_r) as $f(w) = a_1v_1 + \dots + a_rv_r$, and let $u = a_1u_1 + \dots + a_ru_r$. Observe that

$$\begin{aligned} f(w - u) &= f(w) - f(u) \\ &= a_1v_1 + \dots + a_rv_r - (a_1f(u_1) + \dots + a_rf(u_r)) \\ &= a_1v_1 + \dots + a_rv_r - (a_1v_1 + \dots + a_rv_r) \\ &= 0. \end{aligned}$$

Therefore, $h = w - u \in \text{Ker}(f)$, and since $w = h + u$ with $h \in \text{Ker}(f)$ and $u \in U$, we have $E = \text{Ker}(f) + U$.

If $u = a_1u_1 + \dots + a_ru_r \in U$ also belongs to $\text{Ker}(f)$, then

$$0 = f(u) = f(a_1u_1 + \dots + a_ru_r) = a_1v_1 + \dots + a_rv_r,$$

and since (v_1, \dots, v_r) is a basis, $a_i = 0$ for $i = 1, \dots, r$, which shows that $\text{Ker}(f) \cap U = (0)$. Therefore, we have a direct sum

$$E = \text{Ker}(f) \oplus U.$$

Finally, if

$$a_1u_1 + \dots + a_ru_r = 0,$$

the above reasoning shows that $a_i = 0$ for $i = 1, \dots, r$, so (u_1, \dots, u_r) are linearly independent. Therefore, the module U is a free module. \square

One should be aware that if we have a direct sum of modules

$$U = U_1 \oplus \cdots \oplus U_m,$$

every vector $u \in U$ can be written in a unique way as

$$u = u_1 + \cdots + u_m,$$

with $u_i \in U_i$ but, unlike the case of vector spaces, this does not imply that any m nonzero vectors (u_1, \dots, u_m) are linearly independent. For example,

$$\mathbb{Z} = \mathbb{Z}/2\mathbb{Z} \oplus \mathbb{Z}/2\mathbb{Z}$$

where \mathbb{Z} and $\mathbb{Z}/2\mathbb{Z}$ are viewed as \mathbb{Z} -modules, but $(1, 0)$ and $(0, 1)$ are not linearly independent, since

$$2(1, 0) + 2(0, 1) = (0, 0).$$

A useful fact is that every module is a quotient of some free module. Indeed, if M is an A -module, pick any spanning set I for M (such a set exists, for example, $I = M$), and consider the unique homomorphism $\varphi: A^{(I)} \rightarrow M$ extending the identity function from I to itself. Then we have an isomorphism $A^{(I)}/\text{Ker}(\varphi) \approx M$.

In particular, if M is finitely generated, we can pick I to be a finite set of generators, in which case we get an isomorphism $A^n/\text{Ker}(\varphi) \approx M$, for some natural number n . A finitely generated module is sometimes called a module of *finite type*.

The case $n = 1$ is of particular interest. A module M is said to be *cyclic* if it is generated by a single element. In this case $M = Ax$, for some $x \in M$. We have the linear map $m_x: A \rightarrow M$ given by $a \mapsto ax$ for every $a \in A$, and it is obviously surjective since $M = Ax$. Since the kernel $\mathfrak{a} = \text{Ker}(m_x)$ of m_x is an ideal in A , we get an isomorphism $A/\mathfrak{a} \approx Ax$. Conversely, for any ideal \mathfrak{a} of A , if $M = A/\mathfrak{a}$, we see that M is generated by the image x of 1 in M , so M is a cyclic module.

The ideal $\mathfrak{a} = \text{Ker}(m_x)$ is the set of all $a \in A$ such that $ax = 0$. This is called the *annihilator* of x , and it is the special case of the following more general situation.

Definition 24.4. If M is any A -module, for any subset S of M , the set of all $a \in A$ such that $ax = 0$ for all $x \in S$ is called the *annihilator* of S , and it is denoted by $\text{Ann}(S)$. If $S = \{x\}$, we write $\text{Ann}(x)$ instead of $\text{Ann}(\{x\})$. A nonzero element $x \in M$ is called a *torsion element* iff $\text{Ann}(x) \neq (0)$. The set consisting of all torsion elements in M and 0 is denoted by M_{tor} .

It is immediately verified that $\text{Ann}(S)$ is an ideal of A , and by definition,

$$M_{\text{tor}} = \{x \in M \mid (\exists a \in A, a \neq 0)(ax = 0)\}.$$

If a ring has zero divisors, then the set of all torsion elements in an A -module M may not be a submodule of M . For example, if $M = A = \mathbb{Z}/6\mathbb{Z}$, then $M_{\text{tor}} = \{2, 3, 4\}$, but $3 + 4 = 1$ is not a torsion element. Also, a free module may not be torsion-free because there may be torsion elements, as the example of $\mathbb{Z}/6\mathbb{Z}$ as a free module over itself shows.

However, if A is an integral domain, then a free module is torsion-free and M_{tor} is a submodule of M . (Recall that an integral domain is commutative).

Proposition 24.3. *If A is an integral domain, then for any A -module M , the set M_{tor} of torsion elements in M is a submodule of M .*

Proof. If $x, y \in M$ are torsion elements ($x, y \neq 0$), then there exist some nonzero elements $a, b \in A$ such that $ax = 0$ and $by = 0$. Since A is an integral domain, $ab \neq 0$, and then for all $\lambda, \mu \in A$, we have

$$ab(\lambda x + \mu y) = b\lambda ax + a\mu by = 0.$$

Therefore, M_{tor} is a submodule of M . □

The module M_{tor} is called the *torsion submodule* of M . If $M_{\text{tor}} = (0)$, then we say that M is *torsion-free*, and if $M = M_{\text{tor}}$, then we say that M is a *torsion module*.

If M is not finitely generated, then it is possible that $M_{\text{tor}} \neq 0$, yet the annihilator of M_{tor} is reduced to 0 (find an example). However, if M is finitely generated, this cannot happen, since if x_1, \dots, x_n generate M and if a_1, \dots, a_n annihilate x_1, \dots, x_n , then $a_1 \cdots a_n$ annihilates every element of M .

Proposition 24.4. *If A is an integral domain, then for any A -module M , the quotient module M/M_{tor} is torsion free.*

Proof. Let \bar{x} be an element of M/M_{tor} and assume that $a\bar{x} = 0$ for some $a \neq 0$ in A . This means that $ax \in M_{\text{tor}}$, so there is some $b \neq 0$ in A such that $bax = 0$. Since $a, b \neq 0$ and A is an integral domain, $ba \neq 0$, so $x \in M_{\text{tor}}$, which means that $\bar{x} = 0$. □

If A is an integral domain and if F is a free A -module with basis (u_1, \dots, u_n) , then F can be embedded in a K -vector space F_K isomorphic to K^n , where $K = \text{Frac}(A)$ is the fraction field of A . Similarly, any submodule M of F is embedded into a subspace M_K of F_K . Note that any linearly independent vectors (u_1, \dots, u_m) in the A -module M remain linearly independent in the vector space M_K , because any linear dependence over K is of the form

$$\frac{a_1}{b_1} u_1 + \cdots + \frac{a_m}{b_m} u_m = 0$$

for some $a_i, b_i \in A$, with $b_1 \cdots b_m \neq 0$, so if we multiply by $b_1 \cdots b_m \neq 0$, we get a linear dependence in the A -module M . Then, we see that the maximum number of linearly independent vectors in the A -module M is at most n . The maximum number of linearly independent vectors in a finitely generated submodule of a free module (over an integral domain) is called the *rank* of the module M . If (u_1, \dots, u_m) are linearly independent where

m is the rank of m , then for every nonzero $v \in M$, there are some $a, a_1, \dots, a_m \in A$, not all zero, such that

$$av = a_1u_1 + \dots + a_mu_m.$$

We must have $a \neq 0$, since otherwise, linear independence of the u_i would imply that $a_1 = \dots = a_m = 0$, contradicting the fact that $a, a_1, \dots, a_m \in A$ are not all zero.

Unfortunately, in general, a torsion-free module is not free. For example, \mathbb{Q} as a \mathbb{Z} -module is torsion-free but not free. If we restrict ourselves to finitely generated modules over PID's, then such modules split as the direct sum of their torsion module with a free module, and a torsion module has a nice decomposition in terms of cyclic modules.

The following proposition shows that over a PID, submodules of a free module are free. There are various ways of proving this result. We give a proof due to Lang [67] (see Chapter III, Section 7).

Proposition 24.5. *If A is a PID and if F is a free A -module of dimension n , then every submodule M of F is a free module of dimension at most n .*

Proof. Let (u_1, \dots, u_n) be a basis of F , and let $M_r = M \cap (Au_1 \oplus \dots \oplus Au_r)$, the intersection of M with the free module generated by (u_1, \dots, u_r) , for $r = 1, \dots, n$. We prove by induction on r that each M_r is free and of dimension at most r . Since $M = M_r$ for some r , this will prove our result.

Consider $M_1 = M \cap Au_1$. If $M_1 = (0)$, we are done. Otherwise let

$$\mathfrak{a} = \{a \in A \mid au_1 \in M\}.$$

It is immediately verified that \mathfrak{a} is an ideal, and since A is a PID, $\mathfrak{a} = a_1A$, for some $a_1 \in A$. Since we are assuming that $M_1 \neq (0)$, we have $a_1 \neq 0$, and $a_1u_1 \in M$. If $x \in M_1$, then $x = au_1$ for some $a \in A$, so $a \in a_1A$, and thus $a = ba_1$ for some $b \in A$. It follows that $M_1 = Aa_1u_1$, which is free.

Assume inductively that M_r is free of dimension at most $r < n$, and let

$$\mathfrak{a} = \{a \in A \mid (\exists b_1 \in A) \dots (\exists b_r \in A)(b_1u_1 + \dots + b_ru_r + au_{r+1} \in M)\}.$$

It is immediately verified that \mathfrak{a} is an ideal, and since A is a PID, $\mathfrak{a} = a_{r+1}A$, for some $a_{r+1} \in A$. If $a_{r+1} = 0$, then $M_{r+1} = M_r$, and we are done.

If $a_{r+1} \neq 0$, then there is some $v_1 \in Au_1 \oplus \dots \oplus Au_r$ such that

$$w = v_1 + a_{r+1}u_{r+1} \in M.$$

For any $x \in M_{r+1}$, there is some $v \in Au_1 \oplus \dots \oplus Au_r$ and some $a \in A$ such that $x = v + au_{r+1}$. Then, $a \in a_{r+1}A$, so there is some $b \in A$ such that $a = ba_{r+1}$. As a consequence

$$x - bw = v - bv_1 \in M_r,$$

and so $x = x - bw + bw$ with $x - bw \in M_r$, which shows that

$$M_{r+1} = M_r + Aw.$$

On the other hand, if $u \in M_r \cap Aw$, then since $w = v_1 + a_{r+1}u_{r+1}$ we have

$$u = bv_1 + ba_{r+1}u_{r+1},$$

for some $b \in A$, with $u, v_1 \in Au_1 \oplus \cdots \oplus Au_r$, and if $b \neq 0$, this yields the nontrivial linear combination

$$bv_1 - u + ba_{r+1}u_{r+1} = 0,$$

contradicting the fact that (u_1, \dots, u_{r+1}) are linearly independent. Therefore,

$$M_{r+1} = M_r \oplus Aw,$$

which shows that M_{r+1} is free of dimension at most $r + 1$. \square

Proposition 24.5 implies that if M is a finitely generated module over a PID, then any submodule N of M is also finitely generated.

Indeed, if (u_1, \dots, u_n) generate M , then we have a surjection $\varphi: A^n \rightarrow M$ from the free module A^n onto M . The inverse image $\varphi^{-1}(N)$ of N is a submodule of the free module A^n , therefore by Proposition 24.5, $\varphi^{-1}(N)$ is free and finitely generated. This implies that N is finitely generated (and that it has a number of generators $\leq n$).

We can also prove that a finitely generated torsion-free module over a PID is actually free. We will give another proof of this fact later, but the following proof is instructive.

Proposition 24.6. *If A is a PID and if M is a finitely generated module which is torsion-free, then M is free.*

Proof. Let (y_1, \dots, y_n) be some generators for M , and let (u_1, \dots, u_m) be a maximal subsequence of (y_1, \dots, y_n) which is linearly independent. If $m = n$, we are done. Otherwise, due to the maximality of m , for $i = 1, \dots, n$, there is some $a_i \neq 0$ such that $a_i y_i$ can be expressed as a linear combination of (u_1, \dots, u_m) . If we let $a = a_1 \dots a_n$, then $a_1 \dots a_n y_i \in Au_1 \oplus \cdots \oplus Au_m$ for $i = 1, \dots, n$, which shows that

$$aM \subseteq Au_1 \oplus \cdots \oplus Au_m.$$

Now, A is an integral domain, and since $a_i \neq 0$ for $i = 1, \dots, n$, we have $a = a_1 \dots a_n \neq 0$, and because M is torsion-free, the map $x \mapsto ax$ is injective. It follows that M is isomorphic to a submodule of the free module $Au_1 \oplus \cdots \oplus Au_m$. By Proposition 24.5, this submodule is free, and thus, M is free. \square

Although we will obtain this result as a corollary of the structure theorem for finitely generated modules over a PID, we are in the position to give a quick proof of the following theorem.

Theorem 24.7. *Let M be a finitely generated module over a PID. Then M/M_{tor} is free, and there exist a free submodule F of M such that M is the direct sum*

$$M = M_{\text{tor}} \oplus F.$$

The dimension of F is uniquely determined.

Proof. By Proposition 24.4 M/M_{tor} is torsion-free, and since M is finitely generated, it is also finitely generated. By Proposition 24.6, M/M_{tor} is free. We have the quotient linear map $\pi: M \rightarrow M/M_{\text{tor}}$, which is surjective, and M/M_{tor} is free, so by Proposition 24.2, there is a free module F isomorphic to M/M_{tor} such that

$$M = \text{Ker}(\pi) \oplus F = M_{\text{tor}} \oplus F.$$

Since F is isomorphic to M/M_{tor} , the dimension of F is uniquely determined. \square

Theorem 24.7 reduces the study of finitely generated module over a PID to the study of finitely generated torsion modules. This is the path followed by Lang [67] (Chapter III, section 7).

24.2 Finite Presentations of Modules

Since modules are generally not free, it is natural to look for techniques for dealing with nonfree modules. The hint is that if M is an A -module and if $(u_i)_{i \in I}$ is any set of generators for M , then we know that there is a surjective homomorphism $\varphi: A^{(I)} \rightarrow M$ from the free module $A^{(I)}$ generated by I onto M . Furthermore M is isomorphic to $A^{(I)}/\text{Ker}(\varphi)$. Then, we can pick a set of generators $(v_j)_{j \in J}$ for $\text{Ker}(\varphi)$, and again there is a surjective map $\psi: A^{(J)} \rightarrow \text{Ker}(\varphi)$ from the free module $A^{(J)}$ generated by J onto $\text{Ker}(\varphi)$. The map ψ can be viewed a linear map from $A^{(J)}$ to $A^{(I)}$, we have

$$\text{Im}(\psi) = \text{Ker}(\varphi),$$

and φ is surjective. Note that M is isomorphic to $A^{(I)}/\text{Im}(\psi)$. In such a situation we say that we have an *exact sequence* and this is denoted by the diagram

$$A^{(J)} \xrightarrow{\psi} A^{(I)} \xrightarrow{\varphi} M \longrightarrow 0.$$

Definition 24.5. Given an A -module M , a *presentation* of M is an exact sequence

$$A^{(J)} \xrightarrow{\psi} A^{(I)} \xrightarrow{\varphi} M \longrightarrow 0$$

which means that

1. $\text{Im}(\psi) = \text{Ker}(\varphi)$.

2. φ is surjective.

Consequently, M is isomorphic to $A^{(I)}/\text{Im}(\psi)$. If I and J are both finite, we say that this is a *finite presentation* of M .

Observe that in the case of a finite presentation, I and J are finite, and if $|J| = n$ and $|I| = m$, then ψ is a linear map $\psi: A^n \rightarrow A^m$, so it is given by some $m \times n$ matrix R with coefficients in A called the *presentation matrix* of M . Every column R^j of R may be thought of as a relation

$$a_{j1}e_1 + \cdots + a_{jm}e_m = 0$$

among the generators e_1, \dots, e_m of A^m , so we have n relations among these generators. Also the images of e_1, \dots, e_m in M are generators of M , so we can think of the above relations as relations among the generators of M . The submodule of A^m spanned by the columns of R is the *set of relations* of M , and the columns of R are called a *complete set of relations* for M . The vectors e_1, \dots, e_m are called a set of *generators* for M . We may also say that the generators e_1, \dots, e_m and the relations R^1, \dots, R^n (the columns of R) are a (finite) presentation of the module M .

For example, the \mathbb{Z} -module presented by the 1×1 matrix $R = (5)$ is the quotient, $\mathbb{Z}/5\mathbb{Z}$, of \mathbb{Z} by the submodule $5\mathbb{Z}$ corresponding to the single relation

$$5e_1 = 0.$$

But $\mathbb{Z}/5\mathbb{Z}$ has other presentations. For example, if we consider the matrix of relations

$$R = \begin{pmatrix} 2 & -1 \\ 1 & 2 \end{pmatrix},$$

presenting the module M , then we have the relations

$$\begin{aligned} 2e_1 + e_2 &= 0 \\ -e_1 + 2e_2 &= 0. \end{aligned}$$

From the first equation, we get $e_2 = -2e_1$, and substituting into the second equation we get

$$-5e_1 = 0.$$

It follows that the generator e_2 can be eliminated and M is generated by the single generator e_1 satisfying the relation

$$5e_1 = 0,$$

which shows that $M \approx \mathbb{Z}/5\mathbb{Z}$.

The above example shows that many different matrices can present the same module. Here are some useful rules for manipulating a relation matrix without changing the isomorphism class of the module M it presents.

Proposition 24.8. *If R is an $m \times n$ matrix presenting an A -module M , then the matrices S of the form listed below present the same module (a module isomorphic to M):*

- (1) $S = QRP^{-1}$, where Q is a $m \times m$ invertible matrix and P a $n \times n$ invertible matrix (both over A).
- (2) S is obtained from R by deleting a column of zeros.
- (3) The j th column of R is e_i , and S is obtained from R by deleting the i th row and the j th column.

Proof. (1) By definition, we have an isomorphism $M \approx A^m/RA^n$, where we denote by RA^n the image of A^n by the linear map defined by R . Going from R to QRP^{-1} corresponds to making a change of basis in A^m and a change of basis in A^n , and this yields a quotient module isomorphic to M .

(2) A zero column does not contribute to the span of the columns of R , so it can be eliminated.

(3) If the j th column of R is e_i , then when taking the quotient A^m/RA^n , the generator e_i goes to zero. This means that the generator e_i is redundant, and when we delete it, we get a matrix of relations in which the i th row of R and the j th column of R are deleted. \square

The matrices P and Q are often products of elementary operations. One should be careful that rows of zeros cannot be eliminated. For example, the 2×1 matrix

$$R_1 = \begin{pmatrix} 4 \\ 0 \end{pmatrix}$$

gives the single relation

$$4e_1 = 0,$$

but the second generator e_2 cannot be eliminated. This matrix presents the module $\mathbb{Z}/4\mathbb{Z} \times \mathbb{Z}$. On the other hand, the 1×2 matrix

$$R_2 = \begin{pmatrix} 4 & 0 \end{pmatrix}$$

gives two relations

$$\begin{aligned} 4e_1 &= 0, \\ 0 &= 0, \end{aligned}$$

so the second generator can be eliminated and R_2 presents the module $\mathbb{Z}/4\mathbb{Z}$.

The rules of Proposition 24.8 make it possible to simplify a presentation matrix quite a lot in some cases. For example, consider the relation matrix

$$R = \begin{pmatrix} 3 & 8 & 7 & 9 \\ 2 & 4 & 6 & 6 \\ 1 & 2 & 2 & 1 \end{pmatrix}.$$

By subtracting 2 times row 3 from row 2 and subtracting 3 times row 3 from row 1, we get

$$\begin{pmatrix} 0 & 2 & 1 & 6 \\ 0 & 0 & 2 & 4 \\ 1 & 2 & 2 & 1 \end{pmatrix}.$$

After deleting column 1 and row 3, we get

$$\begin{pmatrix} 2 & 1 & 6 \\ 0 & 2 & 4 \end{pmatrix}.$$

By subtracting 2 times row 1 from row 2, we get

$$\begin{pmatrix} 2 & 1 & 6 \\ -4 & 0 & -8 \end{pmatrix}.$$

After deleting column 2 and row 1, we get

$$\begin{pmatrix} -4 & -8 \end{pmatrix}.$$

By subtracting 2 times column 1 from column 2, we get

$$\begin{pmatrix} -4 & 0 \end{pmatrix}.$$

Finally, we can drop the second column and we get

$$(4),$$

which shows that R presents the module $\mathbb{Z}/4\mathbb{Z}$.

Unfortunately a submodule of a free module of finite dimension is not necessarily finitely generated but, by Proposition 24.5, if A is a PID, then any submodule of a finitely generated module is finitely generated. This property actually characterizes Noetherian rings. To prove it, we need a slightly different version of Proposition 24.2.

Proposition 24.9. *Let $f: E \rightarrow F$ be a linear map between two A -modules E and F .*

- (1) *Given any set of generators (v_1, \dots, v_r) of $\text{Im}(f)$, for any r vectors $u_1, \dots, u_r \in E$ such that $f(u_i) = v_i$ for $i = 1, \dots, r$, if U is the finitely generated submodule of E generated by (u_1, \dots, u_r) , then the module E is the sum*

$$E = \text{Ker}(f) + U.$$

Consequently, if both $\text{Ker}(f)$ and $\text{Im}(f)$ are finitely generated, then E is finitely generated.

- (2) *If E is finitely generated, then so is $\text{Im}(f)$.*

Proof. (1) Pick any $w \in E$, write $f(w)$ over the generators (v_1, \dots, v_r) of $\text{Im}(f)$ as $f(w) = a_1v_1 + \dots + a_rv_r$, and let $u = a_1u_1 + \dots + a_ru_r$. Observe that

$$\begin{aligned} f(w - u) &= f(w) - f(u) \\ &= a_1v_1 + \dots + a_rv_r - (a_1f(u_1) + \dots + a_rf(u_r)) \\ &= a_1v_1 + \dots + a_rv_r - (a_1v_1 + \dots + a_rv_r) \\ &= 0. \end{aligned}$$

Therefore, $h = w - u \in \text{Ker}(f)$, and since $w = h + u$ with $h \in \text{Ker}(f)$ and $u \in U$, we have $E = \text{Ker}(f) + U$, as claimed. If $\text{Ker}(f)$ is also finitely generated, by taking the union of a finite set of generators for $\text{Ker}(f)$ and (v_1, \dots, v_r) , we obtain a finite set of generators for E .

(2) If (u_1, \dots, u_n) generate E , it is obvious that $(f(u_1), \dots, f(u_n))$ generate $\text{Im}(f)$. \square

Theorem 24.10. *A ring A is Noetherian iff every submodule N of a finitely generated A -module M is itself finitely generated.*

Proof. First, assume that every submodule N of a finitely generated A -module M is itself finitely generated. The ring A is a module over itself and it is generated by the single element 1. Furthermore, every submodule of A is an ideal, so the hypothesis implies that every ideal in A is finitely generated, which shows that A is Noetherian.

Now, assume A is Noetherian. First, observe that it is enough to prove the theorem for the finitely generated free modules A^n (with $n \geq 1$). Indeed, assume that we proved for every $n \geq 1$ that every submodule of A^n is finitely generated. If M is any finitely generated A -module, then there is a surjection $\varphi: A^n \rightarrow M$ for some n (where n is the number of elements of a finite generating set for M). Given any submodule N of M , $L = \varphi^{-1}(N)$ is a submodule of A^n . Since A^n is finitely generated, the submodule N of A^n is finitely generated, and then $N = \varphi(L)$ is finitely generated.

It remains to prove the theorem for $M = A^n$. We proceed by induction on n . For $n = 1$, a submodule N of A is an ideal, and since A is Noetherian, N is finitely generated. For the induction step where $n > 1$, consider the projection $\pi: A^n \rightarrow A^{n-1}$ given by

$$\pi(a_1, \dots, a_n) = (a_1, \dots, a_{n-1}).$$

The kernel of π is the module

$$\text{Ker}(\pi) = \{(0, \dots, 0, a_n) \in A^n \mid a_n \in A\} \approx A.$$

For any submodule N of A^n , let $\varphi: N \rightarrow A^{n-1}$ be the restriction of π to N . Since $\varphi(N)$ is a submodule of A^{n-1} , by the induction hypothesis, $\text{Im}(\varphi) = \varphi(N)$ is finitely generated. Also, $\text{Ker}(\varphi) = N \cap \text{Ker}(\pi)$ is a submodule of $\text{Ker}(\pi) \approx A$, and thus $\text{Ker}(\varphi)$ is isomorphic to an ideal of A , and thus is finitely generated (since A is Noetherian). Since both $\text{Im}(\varphi)$ and $\text{Ker}(\varphi)$ are finitely generated, by Proposition 24.9, the submodule N is also finitely generated. \square

As a consequence of Theorem 24.10, every finitely generated A -module over a Noetherian ring A is finitely presented, because if $\varphi: A^n \rightarrow M$ is a surjection onto the finitely generated module M , then $\text{Ker}(\varphi)$ is finitely generated. In particular, if A is a PID, then every finitely generated module is finitely presented.

If the ring A is not Noetherian, then there exist finitely generated A -modules that are not finitely presented. This is not so easy to prove.

We will prove in Theorem 25.14 that if A is a Euclidean ring, and more generally in Theorem 25.17 if A is a PID, then a matrix R can “diagonalized” as

$$R = QDP^{-1}$$

where D is a diagonal matrix. It follows from Proposition 24.8 that every finitely generated module M over a PID has a presentation with m generators and r relations of the form

$$\alpha_i e_i = 0,$$

where $\alpha_i \neq 0$ and $\alpha_1 \mid \alpha_2 \mid \cdots \mid \alpha_r$, which shows that M is isomorphic to the direct sum

$$M \approx A^{m-r} \oplus A/(\alpha_1 A) \oplus \cdots \oplus A/(\alpha_r A).$$

This is a version of Theorem 24.32 that will be proved in Section 24.7.

24.3 Tensor Products of Modules over a Commutative Ring

It is possible to define tensor products of modules over a ring, just as in Section 23.1, and the results of this section continue to hold. The results of Section 23.3 also continue to hold since they are based on the universal mapping property. However, the results of Section 23.2 on bases generally fail, except for free modules. Similarly, the results of Section 23.4 on duality generally fail. Tensor algebras can be defined for modules, as in Section 23.5. Symmetric tensor and alternating tensors can be defined for modules but again, results involving bases generally fail.

Tensor products of modules have some unexpected properties. For example, if p and q are relatively prime integers, then

$$\mathbb{Z}/p\mathbb{Z} \otimes_{\mathbb{Z}} \mathbb{Z}/q\mathbb{Z} = (0).$$

This is because, by Bezout’s identity, there are $a, b \in \mathbb{Z}$ such that

$$ap + bq = 1,$$

so, for all $x \in \mathbb{Z}/p\mathbb{Z}$ and all $y \in \mathbb{Z}/q\mathbb{Z}$, we have

$$\begin{aligned} x \otimes y &= ap(x \otimes y) + bq(x \otimes y) \\ &= a(px \otimes y) + b(x \otimes qy) \\ &= a(0 \otimes y) + b(x \otimes 0) \\ &= 0. \end{aligned}$$

It is possible to salvage certain properties of tensor products holding for vector spaces by restricting the class of modules under consideration. For example, *projective modules* have a pretty good behavior w.r.t. tensor products.

A free A -module F , is a module that has a basis (*i.e.*, there is a family, $(e_i)_{i \in I}$, of linearly independent vectors in F that span F). Projective modules have many equivalent characterizations. Here is one that is best suited for our needs:

Definition 24.6. An A -module, P , is *projective* if it is a summand of a free module, that is, if there is a free A -module, F , and some A -module, Q , so that

$$F = P \oplus Q.$$

Given any A -module, M , we let $M^* = \text{Hom}_A(M, A)$ be its *dual*. We have the following proposition:

Proposition 24.11. *For any finitely-generated projective A -modules, P , and any A -module, Q , we have the isomorphisms:*

$$\begin{aligned} P^{**} &\cong P \\ \text{Hom}_A(P, Q) &\cong P^* \otimes_A Q. \end{aligned}$$

Proof sketch. We only consider the second isomorphism. Since P is projective, we have some A -modules, P_1, F , with

$$P \oplus P_1 = F,$$

where F is some free module. Now, we know that for any A -modules, U, V, W , we have

$$\text{Hom}_A(U \oplus V, W) \cong \text{Hom}_A(U, W) \coprod \text{Hom}_A(V, W) \cong \text{Hom}_A(U, W) \oplus \text{Hom}_A(V, W),$$

so

$$P^* \oplus P_1^* \cong F^*, \quad \text{Hom}_A(P, Q) \oplus \text{Hom}_A(P_1, Q) \cong \text{Hom}_A(F, Q).$$

By tensoring with Q and using the fact that tensor distributes w.r.t. coproducts, we get

$$(P^* \otimes_A Q) \oplus (P_1^* \otimes_A Q) \cong (P^* \oplus P_1^*) \otimes_A Q \cong F^* \otimes_A Q.$$

Now, the proof of Proposition 23.9 goes through because F is free and finitely generated, so

$$\alpha_\otimes: (P^* \otimes_A Q) \oplus (P_1^* \otimes_A Q) \cong F^* \otimes_A Q \longrightarrow \text{Hom}_A(F, Q) \cong \text{Hom}_A(P, Q) \oplus \text{Hom}_A(P_1, Q)$$

is an isomorphism and as α_\otimes maps $P^* \otimes_A Q$ to $\text{Hom}_A(P, Q)$, it yields an isomorphism between these two spaces. \square

The isomorphism $\alpha_{\otimes}: P^* \otimes_A Q \cong \text{Hom}_A(P, Q)$ of Proposition 24.11 is still given by

$$\alpha_{\otimes}(u^* \otimes f)(x) = u^*(x)f, \quad u^* \in P^*, f \in Q, x \in P.$$

It is convenient to introduce the *evaluation map*, $\text{Ev}_x: P^* \otimes_A Q \rightarrow Q$, defined for every $x \in P$ by

$$\text{Ev}_x(u^* \otimes f) = u^*(x)f, \quad u^* \in P^*, f \in Q.$$

We will need the following generalization of part (4) of Proposition 23.7.

Proposition 24.12. *Given any two families of A -modules $(M_i)_{i \in I}$ and $(N_j)_{j \in J}$ (where I and J are finite index sets), we have an isomorphism*

$$\left(\bigoplus_{i \in I} M_i \right) \otimes \left(\bigoplus_{j \in J} M_j \right) \approx \bigoplus_{(i,j) \in I \times J} (M_i \otimes N_j).$$

Proposition 24.12 also holds for infinite index sets.

Proposition 24.13. *Let M and N be two A -module with N a free module, and pick any basis (v_1, \dots, v_n) for N . Then, every element of $M \otimes N$ can expressed in a unique way as a sum of the form*

$$u_1 \otimes v_1 + \dots + u_n \otimes v_n, \quad u_i \in M,$$

so that $M \otimes N$ is isomorphic to M^n (as an A -module).

Proof. Since N is free with basis (v_1, \dots, v_n) , we have an isomorphism

$$N \approx Av_1 \oplus \dots \oplus Av_n.$$

By Proposition 24.12, we obtain an isomorphism

$$M \otimes N \approx M \otimes (Av_1 \oplus \dots \oplus Av_n) \approx (M \otimes Av_1) \oplus \dots \oplus (M \otimes Av_n).$$

Because (v_1, \dots, v_n) is a basis of N , each v_j is torsion-free so the map $a \mapsto av_j$ is an isomorphism of A onto Av_j , and because $M \otimes A \approx M$, we have the isomorphism

$$M \otimes N \approx (M \otimes A) \oplus \dots \oplus (M \otimes A) \approx M \oplus \dots \oplus M = M^n,$$

as claimed. □

Proposition 24.13 also holds for an infinite basis $(v_j)_{j \in J}$ of N . Obviously, a version of Proposition 24.13 also holds if M is free and N is arbitrary.

The next proposition will be also be needed.

Proposition 24.14. *Given any A -module M and any ideal \mathfrak{a} in A , there is an isomorphism*

$$(A/\mathfrak{a}) \otimes_A M \approx M/\mathfrak{a}M$$

given by the map $(\bar{a} \otimes u) \mapsto au \pmod{\mathfrak{a}M}$, for all $\bar{a} \in A/\mathfrak{a}$ and all $u \in M$.

Sketch of proof. Consider the map $\varphi: (A/\mathfrak{a}) \times M \rightarrow M/\mathfrak{a}M$ given by

$$\varphi(\bar{a}, u) = au \pmod{\mathfrak{a}M}$$

for all $\bar{a} \in A/\mathfrak{a}$ and all $u \in M$. It is immediately checked that φ is well-defined because $au \pmod{\mathfrak{a}M}$ does not depend on the representative $a \in A$ chosen in the equivalence class \bar{a} , and φ is bilinear. Therefore, φ induces a linear map $\varphi: (A/\mathfrak{a}) \otimes M \rightarrow M/\mathfrak{a}M$, such that $\varphi(\bar{a} \otimes u) = au \pmod{\mathfrak{a}M}$. We also define the map $\psi: M \rightarrow (A/\mathfrak{a}) \otimes M$ by

$$\psi(u) = \bar{1} \otimes u.$$

Since $\mathfrak{a}M$ is generated by vectors of the form au with $a \in \mathfrak{a}$ and $u \in M$, and since

$$\psi(au) = \bar{1} \otimes au = \bar{a} \otimes u = 0 \otimes u = 0,$$

we see that $\mathfrak{a}M \subseteq \text{Ker}(\psi)$, so ψ induces a linear map $\psi: M/\mathfrak{a}M \rightarrow (A/\mathfrak{a}) \otimes M$. We have

$$\begin{aligned} \psi(\varphi(\bar{a} \otimes u)) &= \psi(au) \\ &= \bar{1} \otimes au \\ &= \bar{a} \otimes u \end{aligned}$$

and

$$\begin{aligned} \varphi(\psi(u)) &= \varphi(\bar{1} \otimes u) \\ &= 1u \\ &= u, \end{aligned}$$

which shows that φ and ψ are mutual inverses. □

24.4 Extension of the Ring of Scalars

The need to extend the ring of scalars arises, in particular when dealing with eigenvalues. First, we need to define how to restrict scalar multiplication to a subring. The situation is that we have two rings A and B , a B -module M , and a ring homomorphism $\rho: A \rightarrow B$. The special case that arises often is that A is a subring of B (B could be a field) and ρ is the inclusion map. Then, we can make M into an A -module by defining the scalar multiplication $\cdot: A \times M \rightarrow M$ as follows:

$$a \cdot x = \rho(a)x, \quad \text{for all } a \in A \text{ and all } x \in M.$$

This A -module is denoted by $\rho_*(M)$. In particular, viewing B as B -module, we obtain the A -module $\rho_*(B)$.

Now, we can describe the process of scalar extension. Given any A -module M , we make $\rho_*(B) \otimes_A M$ into a (left) B -module as follows: for every $\beta \in B$, let $\mu_\beta: \rho_*(B) \times M \rightarrow \rho_*(B) \otimes_A M$ be given by

$$\mu_\beta(\beta', x) = (\beta\beta') \otimes x.$$

The map μ_β is bilinear so it induces a linear map $\mu_\beta: \rho_*(B) \otimes_A M \rightarrow \rho_*(B) \otimes_A M$ such that

$$\mu_\beta(\beta' \otimes x) = (\beta\beta') \otimes x.$$

If we define the scalar multiplication $\cdot: B \times (\rho_*(B) \otimes_A M) \rightarrow \rho_*(B) \otimes_A M$ by

$$\beta \cdot z = \mu_\beta(z), \quad \text{for all } \beta \in B \text{ and all } z \in \rho_*(B) \otimes_A M,$$

then it is easy to check that the axioms M1, M2, M3, M4 hold. Let us check M2 and M3. We have

$$\begin{aligned} \mu_{\beta_1+\beta_2}(\beta' \otimes x) &= (\beta_1 + \beta_2)\beta' \otimes x \\ &= (\beta_1\beta' + \beta_2\beta') \otimes x \\ &= \beta_1\beta' \otimes x + \beta_2\beta' \otimes x \\ &= \mu_{\beta_1}(\beta' \otimes x) + \mu_{\beta_2}(\beta' \otimes x) \end{aligned}$$

and

$$\begin{aligned} \mu_{\beta_1\beta_2}(\beta' \otimes x) &= \beta_1\beta_2\beta' \otimes x \\ &= \mu_{\beta_1}(\beta_2\beta' \otimes x) \\ &= \mu_{\beta_1}(\mu_{\beta_2}(\beta' \otimes x)). \end{aligned}$$

With the scalar multiplication by elements of B given by

$$\beta \cdot (\beta' \otimes x) = (\beta\beta') \otimes x,$$

the tensor product $\rho_*(B) \otimes_A M$ is a B -module denoted by $\rho^*(M)$, or $M_{(B)}$ when ρ is the inclusion of A into B . The B -module $\rho^*(M)$ is sometimes called the *module induced from M by extension to B of the ring of scalars through ρ* .

The above process can also be applied to linear maps. We have the following proposition whose proof is given in Bourbaki [14] (Chapter II, Section 5, Proposition 1).

Proposition 24.15. *Given a ring homomorphism $\rho: A \rightarrow B$ and given any A -module M , the map $\varphi: M \rightarrow \rho_*(\rho^*(M))$ given by $\varphi(x) = 1 \otimes x$ is A -linear and $\varphi(M)$ spans the B -module $\rho^*(M)$. For every B -module N , and for every A -linear map $f: M \rightarrow \rho_*(N)$, there is a unique B -linear map $\bar{f}: \rho^*(M) \rightarrow N$ such that*

$$\bar{f} \circ \varphi = f,$$

or equivalently,

$$\bar{f}(1 \otimes x) = f(x), \quad \text{for all } x \in M.$$

As a consequence of Proposition 24.15, we obtain the following result.

Proposition 24.16. *Given a ring homomorphism $\rho: A \rightarrow B$, for any two A -modules M and N , for every A -linear map $f: M \rightarrow N$, there is a unique B -linear map $\bar{f}: \rho^*(M) \rightarrow \rho^*(N)$ (also denoted by $\rho^*(f)$) given by*

$$\bar{f} = \text{id}_B \otimes f,$$

such that the following diagram commutes:

$$\begin{array}{ccc} M & \xrightarrow{\varphi_M} & \rho^*(M) \\ f \downarrow & & \downarrow \bar{f} \\ N & \xrightarrow{\varphi_N} & \rho^*(N) \end{array}$$

Proof. Apply Proposition 24.16 to the A -linear map $\varphi_N \circ f$. □

If S spans the module M , it is clear that $\varphi(S)$ spans $\rho^*(M)$. In particular, if M is finitely generated, so is $\rho^*(M)$. Bases of M also extend to bases of $\rho^*(M)$.

Proposition 24.17. *Given a ring homomorphism $\rho: A \rightarrow B$, for any A -modules M , if (u_1, \dots, u_n) is a basis of M , then $(\varphi(u_1), \dots, \varphi(u_n))$ is a basis of $\rho^*(M)$, where φ is the A -linear map given by $\varphi(x) = 1 \otimes x$. Furthermore, if ρ is injective, then so is φ .*

Proof. The first assertion follows immediately from Proposition 24.13, since it asserts that every element z of $\rho^*(M) = \rho_*(B) \otimes_A M$ can be written in a unique way as

$$z = b_1 \otimes u_1 + \dots + b_n \otimes u_n = b_1(1 \otimes u_1) + \dots + b_n(1 \otimes u_n),$$

and $\varphi(u_i) = 1 \otimes u_i$. Next, if ρ is injective, by definition of the scalar multiplication in the A -module $\rho_*(\rho^*(M))$, we have $\varphi(a_1 u_1 + \dots + a_n u_n) = 0$ iff

$$\rho(a_1)\varphi(u_1) + \dots + \rho(a_n)\varphi(u_n) = 0,$$

and since $(\varphi(u_1), \dots, \varphi(u_n))$ is a basis of $\rho^*(M)$, we must have $\rho(a_i) = 0$ for $i = 1, \dots, n$, which (by injectivity of ρ) implies that $a_i = 0$ for $i = 1, \dots, n$. Therefore, φ is injective. □

In particular, if A is a subring of B , then ρ is the inclusion map and Proposition 24.17 shows that a basis of M becomes a basis of $M_{(B)}$ and that M is embedded into $M_{(B)}$. It is also easy to see that if M and N are two free A -modules and $f: M \rightarrow N$ is a linear map represented by the matrix X with respect to some bases (u_1, \dots, u_n) of M and (v_1, \dots, v_m) of N , then the B -linear map \bar{f} is also represented by the matrix X over the bases $(\varphi(u_1), \dots, \varphi(u_n))$ and $(\varphi(v_1), \dots, \varphi(v_m))$.

Proposition 24.17 yields another proof of the fact that any two bases of a free A -modules have the same cardinality. Indeed, if \mathfrak{m} is a maximal ideal in the ring A , then we have the quotient ring homomorphism $\pi: A \rightarrow A/\mathfrak{m}$, and we get the A/\mathfrak{m} -module $\pi^*(M)$. If M is

free, any basis (u_1, \dots, u_n) of M becomes the basis $(\varphi(u_1), \dots, \varphi(u_n))$ of $\pi^*(M)$; but A/\mathfrak{m} is a field, so the dimension n is uniquely determined. This argument also applies to an infinite basis $(u_i)_{i \in I}$. Observe that by Proposition 24.14, we have an isomorphism

$$\pi^*(M) = (A/\mathfrak{m}) \otimes_A M \approx M/\mathfrak{m}M,$$

so $M/\mathfrak{m}M$ is a vector space over the field A/\mathfrak{m} , which is the argument used in Theorem 24.1.

Proposition 24.18. *Given a ring homomorphism $\rho: A \rightarrow B$, for any two A -modules M and N , there is a unique isomorphism*

$$\rho^*(M) \otimes_B \rho^*(N) \approx \rho^*(M \otimes_A N),$$

such that $(1 \otimes u) \otimes (1 \otimes v) \mapsto 1 \otimes (u \otimes v)$, for all $u \in M$ and all $v \in N$.

The proof uses identities from Proposition 23.7. It is not hard but it requires a little gymnastic; a good exercise for the reader.

24.5 The Torsion Module Associated With An Endomorphism

We saw in Section 5.7 that given a linear map $f: E \rightarrow E$ from a K -vector space E into itself, we can define a scalar multiplication $\cdot: K[X] \times E \rightarrow E$ that makes E into a $K[X]$ -module. If E is finite-dimensional, this $K[X]$ -module denoted by E_f is a torsion module, and the main results of this chapter yield important direct sum decompositions of E into subspaces invariant under f .

Recall that given any polynomial $p(X) = a_0X^n + a_1X^{n-1} + \dots + a_n$ with coefficients in the field K , we define the *linear map* $p(f): E \rightarrow E$ by

$$p(f) = a_0f^n + a_1f^{n-1} + \dots + a_n\text{id},$$

where $f^k = f \circ \dots \circ f$, the k -fold composition of f with itself. Note that

$$p(f)(u) = a_0f^n(u) + a_1f^{n-1}(u) + \dots + a_nu,$$

for every vector $u \in E$. Then, we define the scalar multiplication $\cdot: K[X] \times E \rightarrow E$ by polynomials as follows: for every polynomial $p(X) \in K[X]$, for every $u \in E$,

$$p(X) \cdot u = p(f)(u).^3$$

³If necessary to avoid confusion, we use the notion $p(X) \cdot_f u$ instead of $p(X) \cdot u$.

It is easy to verify that this scalar multiplication satisfies the axioms M1, M2, M3, M4:

$$\begin{aligned} p \cdot (u + v) &= p \cdot u + p \cdot v \\ (p + q) \cdot u &= p \cdot u + q \cdot u \\ (pq) \cdot u &= p \cdot (q \cdot u) \\ 1 \cdot u &= u, \end{aligned}$$

for all $p, q \in K[X]$ and all $u, v \in E$. Thus, with this new scalar multiplication, E is a $K[X]$ -module denoted by E_f .

If $p = \lambda$ is just a scalar in K (a polynomial of degree 0), then

$$\lambda \cdot u = (\lambda \text{id})(u) = \lambda u,$$

which means that K acts on E by scalar multiplication as before. If $p(X) = X$ (the monomial X), then

$$X \cdot u = f(u).$$

Since K is a field, the ring $K[X]$ is a PID.

If E is finite-dimensional, say of dimension n , since K is a subring of $K[X]$ and since E is finitely generated over K , the $K[X]$ -module E_f is finitely generated over $K[X]$. Furthermore, E_f is a torsion module. This follows from the Cayley-Hamilton Theorem (Theorem 5.16), but this can also be shown in an elementary fashion as follows. The space $\text{Hom}(E, E)$ of linear maps of E into itself is a vector space of dimension n^2 , therefore the $n^2 + 1$ linear maps

$$\text{id}, f, f^2, \dots, f^{n^2}$$

are linearly dependent, which yields a nonzero polynomial q such that $q(f) = 0$.

We can now translate notions defined for modules into notions for endomorphisms of vector spaces.

1. To say that U is a submodule of E_f means that U is a subspace of E invariant under f ; that is, $f(U) \subseteq U$.
2. To say that V is a cyclic submodule of E_f means that there is some vector $u \in V$, such that V is spanned by $(u, f(u), \dots, f^k(u), \dots)$. If E has finite dimension n , then V is spanned by $(u, f(u), \dots, f^k(u))$ for some $k \leq n - 1$. We say that V is a *cyclic subspace for f with generator u* . Sometimes, V is denoted by $Z(u; f)$.
3. To say that the ideal $\mathfrak{a} = (p(X))$ (with $p(X)$ a monic polynomial) is the annihilator of the submodule V means that $p(f)(u) = 0$ for all $u \in V$, and we call p the *minimal polynomial* of V .

4. Suppose E_f is cyclic and let $\mathfrak{a} = (q)$ be its annihilator, where

$$q(X) = X^n + a_{n-1}X^{n-1} + \cdots + a_1X + a_0.$$

Then, there is some vector u such that $(u, f(u), \dots, f^k(u))$ span E_f , and because q is the minimal polynomial of E_f , we must have $k = n - 1$. The fact that $q(f) = 0$ implies that

$$f^n(u) = -a_0u - a_1f(u) - \cdots - a_{n-1}f^{n-1}(u),$$

and so f is represented by the following matrix known as the *companion matrix* of $q(X)$:

$$U = \begin{pmatrix} 0 & 0 & 0 & \cdots & 0 & -a_0 \\ 1 & 0 & 0 & \cdots & 0 & -a_1 \\ 0 & 1 & 0 & \cdots & 0 & -a_2 \\ \vdots & \ddots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \ddots & 0 & -a_{n-2} \\ 0 & 0 & 0 & \cdots & 1 & -a_{n-1} \end{pmatrix}.$$

It is an easy exercise to prove that the characteristic polynomial $\chi_U(X)$ of U gives back $q(X)$:

$$\chi_U(X) = q(X).$$

We will need the following proposition to characterize when two linear maps are similar.

Proposition 24.19. *Let $f: E \rightarrow E$ and $f': E' \rightarrow E'$ be two linear maps over the vector spaces E and E' . A linear map $g: E \rightarrow E'$ can be viewed as a linear map between the $K[X]$ -modules E_f and $E_{f'}$ iff*

$$g \circ f = f' \circ g.$$

Proof. First, suppose g is $K[X]$ -linear. Then, we have

$$g(p \cdot_f u) = p \cdot_{f'} g(u)$$

for all $p \in K[X]$ and all $u \in E$, so for $p = X$ we get

$$g(p \cdot_f u) = g(X \cdot_f u) = g(f(u))$$

and

$$p \cdot_{f'} g(u) = X \cdot_{f'} g(u) = f'(g(u)),$$

which means that $g \circ f = f' \circ g$.

Conversely, if $g \circ f = f' \circ g$, we prove by induction that

$$g \circ f^n = f'^n \circ g, \quad \text{for all } n \geq 1.$$

Indeed, we have

$$\begin{aligned}
 g \circ f^{n+1} &= g \circ f^n \circ f \\
 &= f'^n \circ g \circ f \\
 &= f'^n \circ f' \circ g \\
 &= f'^{n+1} \circ g,
 \end{aligned}$$

establishing the induction step. It follows that for any polynomial $p(X) = \sum_{k=0}^n a_k X^k$, we have

$$\begin{aligned}
 g(p(X) \cdot_f u) &= g\left(\sum_{k=0}^n a_k f^k(u)\right) \\
 &= \sum_{k=0}^n a_k g \circ f^k(u) \\
 &= \sum_{k=0}^n a_k f'^k \circ g(u) \\
 &= \left(\sum_{k=0}^n a_k f'^k\right)(g(u)) \\
 &= p(X) \cdot_{f'} g(u),
 \end{aligned}$$

so, g is indeed $K[X]$ -linear. □

Definition 24.7. We say that the linear maps $f: E \rightarrow E$ and $f': E' \rightarrow E'$ are *similar* iff there is an isomorphism $g: E \rightarrow E'$ such that

$$f' = g \circ f \circ g^{-1},$$

or equivalently,

$$g \circ f = f' \circ g.$$

Then, Proposition 24.19 shows the following fact:

Proposition 24.20. *With notation of Proposition 24.19, two linear maps f and f' are similar iff g is an isomorphism between E_f and $E'_{f'}$.*

Later on, we will see that the isomorphism of finitely generated torsion modules can be characterized in terms of invariant factors, and this will be translated into a characterization of similarity of linear maps in terms of so-called similarity invariants. If f and f' are represented by matrices A and A' over bases of E and E' , then f and f' are similar iff the matrices A and A' are similar (there is an invertible matrix P such that $A' = PAP^{-1}$). Similar matrices (and endomorphisms) have the same characteristic polynomial.

It turns out that there is a useful relationship between E_f and the module $K[X] \otimes_K E$. Observe that the map $\cdot : K[X] \times E \rightarrow E$ given by

$$p \cdot u = p(f)(u)$$

is K -bilinear, so it yields a K -linear map $\sigma : K[X] \otimes_K E \rightarrow E$ such that

$$\sigma(p \otimes u) = p \cdot u = p(f)(u).$$

We know from Section 24.4 that $K[X] \otimes_K E$ is a $K[X]$ -module (obtained from the inclusion $K \subseteq K[X]$), which we will denote by $E[X]$. Since E is a vector space, $E[X]$ is a free $K[X]$ -module, and if (u_1, \dots, u_n) is a basis of E , then $(1 \otimes u_1, \dots, 1 \otimes u_n)$ is a basis of $E[X]$.

The free $K[X]$ -module $E[X]$ is not as complicated as it looks. Over the basis $(1 \otimes u_1, \dots, 1 \otimes u_n)$, every element $z \in E[X]$ can be written uniquely as

$$z = p_1(1 \otimes u_1) + \dots + p_n(1 \otimes u_n) = p_1 \otimes u_1 + \dots + p_n \otimes u_n,$$

where p_1, \dots, p_n are polynomials in $K[X]$. For notational simplicity, we may write

$$z = p_1 u_1 + \dots + p_n u_n,$$

where p_1, \dots, p_n are viewed as coefficients in $K[X]$. With this notation, we see that $E[X]$ is isomorphic to $(K[X])^n$, which is easy to understand.

Observe that σ is $K[X]$ -linear, because

$$\begin{aligned} \sigma(q(p \otimes u)) &= \sigma((qp) \otimes u) \\ &= (qp) \cdot u \\ &= q(f)(p(f)(u)) \\ &= q \cdot (p(f)(u)) \\ &= q \cdot \sigma(p \otimes u). \end{aligned}$$

Therefore, σ is a linear map of $K[X]$ -modules, $\sigma : E[X] \rightarrow E_f$. Using our simplified notation, if $z = p_1 u_1 + \dots + p_n u_n \in E[X]$, then

$$\sigma(z) = p_1(f)(u_1) + \dots + p_n(f)(u_n),$$

which amounts to plugging f for X and evaluating. Similarly, f is a $K[X]$ -linear map of E_f , because

$$\begin{aligned} f(p \cdot u) &= f(p(f)(u)) \\ &= (fp(f))(u) \\ &= p(f)(f(u)) \\ &= p \cdot f(u), \end{aligned}$$

where we used the fact that $fp(f) = p(f)f$ because $p(f)$ is a polynomial in f . By Proposition 24.16, the linear map $f: E \rightarrow E$ induces a $K[X]$ -linear map $\bar{f}: E[X] \rightarrow E[X]$ such that

$$\bar{f}(p \otimes u) = p \otimes f(u).$$

Observe that we have

$$f(\sigma(p \otimes u)) = f(p(f)(u)) = p(f)(f(u))$$

and

$$\sigma(\bar{f}(p \otimes u)) = \sigma(p \otimes f(u)) = p(f)(f(u)),$$

so we get

$$\sigma \circ \bar{f} = f \circ \sigma. \quad (*)$$

Using our simplified notation,

$$\bar{f}(p_1 u_1 + \cdots + p_n u_n) = p_1 f(u_1) + \cdots + p_n f(u_n).$$

Define the $K[X]$ -linear map $\psi: E[X] \rightarrow E[X]$ by

$$\psi(p \otimes u) = (Xp) \otimes u - p \otimes f(u).$$

Observe that $\psi = X1_{E[X]} - \bar{f}$, which we abbreviate as $X1 - \bar{f}$. Using our simplified notation

$$\psi(p_1 u_1 + \cdots + p_n u_n) = Xp_1 u_1 + \cdots + Xp_n u_n - (p_1 f(u_1) + \cdots + p_n f(u_n)).$$

It should be noted that everything we did in Section 24.5 applies to modules over a commutative ring A , except for the statements that assume that $A[X]$ is a PID. So, if M is an A -module, we can define the $A[X]$ -modules M_f and $M[X] = A[X] \otimes_A M$, except that M_f is generally not a torsion module, and all the results showed above hold. Then, we have the following remarkable result.

Theorem 24.21. (*The Characteristic Sequence*) *Let A be a ring and let E be an A -module. The following sequence of $A[X]$ -linear maps is exact:*

$$0 \longrightarrow E[X] \xrightarrow{\psi} E[X] \xrightarrow{\sigma} E_f \longrightarrow 0.$$

This means that ψ is injective, σ is surjective, and that $\text{Im}(\psi) = \text{Ker}(\sigma)$. As a consequence, E_f is isomorphic to the quotient of $E[X]$ by $\text{Im}(X1 - \bar{f})$.

Proof. Because $\sigma(1 \otimes u) = u$ for all $u \in E$, the map σ is surjective. We have

$$\begin{aligned} \sigma(X(p \otimes u)) &= X \cdot \sigma(p \otimes u) \\ &= f(\sigma(p \otimes u)), \end{aligned}$$

which shows that

$$\sigma \circ X1 = f \circ \sigma = \sigma \circ \bar{f},$$

using $(*)$. This implies that

$$\begin{aligned}\sigma \circ \psi &= \sigma \circ (X1 - \bar{f}) \\ &= \sigma \circ X1 - \sigma \circ \bar{f} \\ &= \sigma \circ \bar{f} - \sigma \circ \bar{f} = 0,\end{aligned}$$

and thus, $\text{Im}(\psi) \subseteq \text{Ker}(\sigma)$. It remains to prove that $\text{Ker}(\sigma) \subseteq \text{Im}(\psi)$.

Since the monomials X^k form a basis of $A[X]$, by Proposition 24.13 (with the roles of M and N exchanged), every $z \in E[X] = A[X] \otimes_A E$ has a unique expression as

$$z = \sum_k X^k \otimes u_k,$$

for a family (u_k) of finite support of $u_k \in E$. If $z \in \text{Ker}(\sigma)$, then

$$0 = \sigma(z) = \sum_k f^k(u_k),$$

which allows us to write

$$\begin{aligned}z &= \sum_k X^k \otimes u_k - 1 \otimes 0 \\ &= \sum_k X^k \otimes u_k - 1 \otimes \left(\sum_k f^k(u_k) \right) \\ &= \sum_k (X^k \otimes u_k - 1 \otimes f^k(u_k)) \\ &= \sum_k (X^k(1 \otimes u_k) - \bar{f}^k(1 \otimes u_k)) \\ &= \sum_k (X^k 1 - \bar{f}^k)(1 \otimes u_k).\end{aligned}$$

Now, $X1$ and \bar{f} commute, since

$$\begin{aligned}(X1 \circ \bar{f})(p \otimes u) &= (X1)(p \otimes f(u)) \\ &= (Xp) \otimes f(u)\end{aligned}$$

and

$$\begin{aligned}(\bar{f} \circ X1)(p \otimes u) &= \bar{f}((Xp) \otimes u) \\ &= (Xp) \otimes f(u),\end{aligned}$$

so we can write

$$X^k 1 - \bar{f}^k = (X1 - \bar{f}) \left(\sum_{j=0}^{k-1} (X1)^j \bar{f}^{k-j-1} \right),$$

and

$$z = (X1 - \bar{f}) \left(\sum_k \left(\sum_{j=0}^{k-1} (X1)^j \bar{f}^{k-j-1} \right) (1 \otimes u_k) \right),$$

which shows that $z = \psi(y)$ for some $y \in E[X]$.

Finally, we prove that ψ is injective as follows. We have

$$\begin{aligned} \psi(z) &= \psi \left(\sum_k X^k \otimes u_k \right) \\ &= (X1 - \bar{f}) \left(\sum_k X^k \otimes u_k \right) \\ &= \sum_k X^{k+1} \otimes (u_k - f(u_{k+1})), \end{aligned}$$

where (u_k) is a family of finite support of $u_k \in E$. If $\psi(z) = 0$, then

$$\sum_k X^{k+1} \otimes (u_k - f(u_{k+1})) = 0,$$

and because the X^k form a basis of $A[X]$, we must have

$$u_k - f(u_{k+1}) = 0, \quad \text{for all } k.$$

Since (u_k) has finite support, there is a largest k , say $m+1$ so that $u_{m+1} = 0$, and then from

$$u_k = f(u_{k+1}),$$

we deduce that $u_k = 0$ for all k . Therefore, $z = 0$, and ψ is injective. \square

Remark: The exact sequence of Theorem 24.21 yields a *presentation* of M_f .

Since $A[X]$ is a free A -module, $A[X] \otimes_A M$ is a free A -module, but $A[X] \otimes_A M$ is generally not a free $A[X]$ -module. However, if M is a free module, then $M[X]$ is a free $A[X]$ -module, since if $(u_i)_{i \in I}$ is a basis for M , then $(1 \otimes u_i)_{i \in I}$ is a basis for $M[X]$. This allows us to define the characteristic polynomial $\chi_f(X)$ of an endomorphism of a free module M as

$$\chi_f(X) = \det(X1 - \bar{f}).$$

Note that to have a correct definition, we need to define the determinant of a linear map allowing the indeterminate X as a scalar, and this is what the definition of $M[X]$ achieves (among other things). Theorem 24.21 can be used to quick a short proof of the Cayley-Hamilton Theorem, see Bourbaki [14] (Chapter III, Section 8, Proposition 20). Proposition 5.10 is still the crucial ingredient of the proof.

We now develop the theory necessary to understand the structure of finitely generated modules over a PID.

24.6 Torsion Modules over a PID; The Primary Decomposition

We begin by considering modules over a product ring obtained from a direct decomposition, as in Definition 21.3. In this section and the next, we closely follow Bourbaki [15] (Chapter VII). Let A be a commutative ring and let $(\mathfrak{b}_1, \dots, \mathfrak{b}_n)$ be ideals in A such that there is an isomorphism $A \approx A/\mathfrak{b}_1 \times \cdots \times A/\mathfrak{b}_n$. From Theorem 21.16 part (b), there exist some elements e_1, \dots, e_n of A such that

$$\begin{aligned} e_i^2 &= e_i \\ e_i e_j &= 0, \quad i \neq j \\ e_1 + \cdots + e_n &= 1_A, \end{aligned}$$

and $\mathfrak{b}_i = (1_A - e_i)A$, for $i, j = 1, \dots, n$.

Given an A -module M with $A \approx A/\mathfrak{b}_1 \times \cdots \times A/\mathfrak{b}_n$, let M_i be the subset of M annihilated by \mathfrak{b}_i ; that is,

$$M_i = \{x \in M \mid bx = 0, \text{ for all } b \in \mathfrak{b}_i\}.$$

Because \mathfrak{b}_i is an ideal, each M_i is a submodule of M . Observe that if $\lambda, \mu \in A$, $b \in \mathfrak{b}_i$, and if $\lambda - \mu = b$, then for any $x \in M_i$, since $bx = 0$,

$$\lambda x = (\mu + b)x = \mu x + bx = \mu x,$$

so M_i can be viewed as a A/\mathfrak{b}_i -module.

Proposition 24.22. *Given a ring $A \approx A/\mathfrak{b}_1 \times \cdots \times A/\mathfrak{b}_n$ as above, the A -module M is the direct sum*

$$M = M_1 \oplus \cdots \oplus M_n,$$

where M_i is the submodule of M annihilated by \mathfrak{b}_i .

Proof. For $i = 1, \dots, n$, let $p_i: M \rightarrow M$ be the map given by

$$p_i(x) = e_i x, \quad x \in M.$$

The map p_i is clearly linear, and because of the properties satisfied by the e_i s, we have

$$\begin{aligned} p_i^2 &= p_i \\ p_i p_j &= 0, \quad i \neq j \\ p_1 + \cdots + p_n &= \text{id}. \end{aligned}$$

This shows that the p_i are projections, and by Proposition 4.6 (which also holds for modules), we have a direct sum

$$M = p_1(M) \oplus \cdots \oplus p_n(M) = e_1 M \oplus \cdots \oplus e_n M.$$

It remains to show that $M_i = e_i M$. Since $(1 - e_i)e_i = e_i - e_i^2 = e_i - e_i = 0$, we see that $e_i M$ is annihilated by $\mathfrak{b}_i = (1 - e_i)A$. Furthermore, for $i \neq j$, for any $x \in M$, we have $(1 - e_i)e_j x = (e_j - e_i e_j)x = e_j x$, so no nonzero element of $e_j M$ is annihilated by $1 - e_i$, and thus not annihilated by \mathfrak{b}_i . It follows that $e_i M = M_i$, as claimed. \square

Given an A -module M , for any nonzero $\alpha \in A$, let

$$M(\alpha) = \{x \in M \mid \alpha x = 0\},$$

the submodule of M annihilated by α . If α divides β , then $M(\alpha) \subseteq M(\beta)$, so we can define

$$M_\alpha = \bigcup_{n \geq 1} M(\alpha^n) = \{x \in M \mid (\exists n \geq 1)(\alpha^n x = 0)\},$$

the submodule of M consisting of all elements of M annihilated by some power of α . If N is any submodule of M , it is clear that

$$N_\alpha = M \cap M_\alpha.$$

Recall that in a PID, an irreducible element is also called a *prime element*.

Definition 24.8. If A is a PID and p is a prime element in A , we say that a module M is *p -primary* if $M = M_p$.

Proposition 24.23. Let M be module over a PID A . For every nonzero $\alpha \in A$, if

$$\alpha = up_1^{n_1} \cdots p_r^{n_r}$$

is a factorization of α into prime factors (where u is a unit), then the module $M(\alpha)$ annihilated by α is the direct sum

$$M(\alpha) = M(p_1^{n_1}) \oplus \cdots \oplus M(p_r^{n_r}).$$

Furthermore, the projection from $M(\alpha)$ onto $M(p_i^{n_i})$ is of the form $x \mapsto \gamma_i x$, for some $\gamma_i \in A$, and

$$M(p_i^{n_i}) = M(\alpha) \cap M_{p_i}.$$

Proof. First, observe that since $M(\alpha)$ is annihilated by α , we can view $M(\alpha)$ as a $A/(\alpha)$ -module. By the Chinese Remainder Theorem (Theorem 21.15) applied to the ideals $(up_1^{n_1}) = (p_1^{n_1}), (p_2^{n_2}), \dots, (p_r^{n_r})$, we have an isomorphism

$$A/(\alpha) \approx A/(p_1^{n_1}) \times \cdots \times A/(p_r^{n_r}).$$

Since we also have isomorphisms

$$A/(p_i^{n_i}) \approx (A/(\alpha))/((p_i^{n_i})/(\alpha)),$$

we can apply Proposition 24.22, and we get a direct sum

$$M(\alpha) = N_1 \oplus \cdots \oplus N_r,$$

where N_i is the $A/(\alpha)$ -submodule of $M(\alpha)$ annihilated by $(p_i^{n_i})/(\alpha)$, and the projections onto the N_i are of the form stated in the proposition. However, N_i is just the A -module $M(p_i^{n_i})$ annihilated by $p_i^{n_i}$, because every nonzero element of $(p_i^{n_i})/(\alpha)$ is an equivalence class modulo (α) of the form $\overline{ap_i^{n_i}}$ for some nonzero $a \in A$, and by definition, $x \in N_i$ iff

$$0 = \overline{ap_i^{n_i}} x = ap_i^{n_i} x, \quad \text{for all } a \in A - \{0\},$$

in particular for $a = 1$, which implies that $x \in M(p_i^{n_i})$.

The inclusion $M(p_i^{n_i}) \subseteq M(\alpha) \cap M_{p_i}$ is clear. Conversely, pick $x \in M(\alpha) \cap M_{p_i}$, which means that $\alpha x = 0$ and $p_i^s x = 0$ for some $s \geq 1$. If $s < n_i$, we are done, so assume $s \geq n_i$. Since $p_i^{n_i}$ is a gcd of α and p_i^s , by Bezout, we can write

$$p_i^{n_i} = \lambda p_i^s + \mu \alpha$$

for some $\lambda, \mu \in A$, and then $p_i^{n_i} x = \lambda p_i^s x + \mu \alpha x = 0$, which shows that $x \in M(p_i^{n_i})$, as desired. \square

Recall that if M is a torsion module over a ring A which is an integral domain, then every finite set of elements x_1, \dots, x_n in M is annihilated by $a = a_1 \cdots a_n$, where each a_i annihilates x_i .

Since A is a PID, we can pick a set P of irreducible elements of A such that every nonzero nonunit of A has a unique factorization up to a unit. Then, we have the following structure theorem for torsion modules which holds even for modules that are not finitely generated.

Theorem 24.24. (*Primary Decomposition Theorem*) *Let M be a torsion-module over a PID. For every irreducible element $p \in P$, let M_p be the submodule of M annihilated by some power of p . Then, M is the (possibly infinite) direct sum*

$$M = \bigoplus_{p \in P} M_p.$$

Proof. Since M is a torsion-module, for every $x \in M$, there is some $\alpha \in A$ such that $x \in M(\alpha)$. By Proposition 24.23, if $\alpha = up_1^{n_1} \cdots p_r^{n_r}$ is a factorization of α into prime factors (where u is a unit), then the module $M(\alpha)$ is the direct sum

$$M(\alpha) = M(p_1^{n_1}) \oplus \cdots \oplus M(p_r^{n_r}).$$

This means that x can be written as

$$x = \sum_{p \in P} x_p, \quad x_p \in M_p,$$

with only finitely many x_p nonzero. If

$$\sum_{p \in P} x_p = \sum_{p \in P} y_p$$

for all $p \in P$, with only finitely many x_p and y_p nonzero, then x_p and y_p are annihilated by some common nonzero element $a \in A$, so $x_p, y_p \in M(a)$. By Proposition 24.23, we must have $x_p = y_p$ for all p , which proves that we have a direct sum. \square

It is clear that if p and p' are two irreducible elements such that $p = up'$ for some unit u , then $M_p = M_{p'}$. Therefore, M_p only depends on the ideal (p) .

Definition 24.9. Given a torsion-module M over a PID, the modules M_p associated with irreducible elements in P are called the *p-primary components* of M .

The p -primary components of a torsion module uniquely determine the module, as shown by the next proposition.

Proposition 24.25. *Two torsion modules M and N over a PID are isomorphic iff for every irreducible element $p \in P$, the p -primary components M_p and N_p of M and N are isomorphic.*

Proof. Let $f: M \rightarrow N$ be an isomorphism. For any $p \in P$, we have $x \in M_p$ iff $p^k x = 0$ for some $k \geq 1$, so

$$0 = f(p^k x) = p^k f(x),$$

which shows that $f(x) \in N_p$. Therefore, f restricts to a linear map $f|_{M_p}$ from M_p to N_p . Since f is an isomorphism, we also have a linear map $f^{-1}: N \rightarrow M$, and our previous reasoning shows that f^{-1} restricts to a linear map $f^{-1}|_{N_p}$ from N_p to M_p . But, $f|_{M_p}$ and $f^{-1}|_{N_p}$ are mutual inverses, so M_p and N_p are isomorphic.

Conversely, if $M_p \approx N_p$ for all $p \in P$, by Theorem 24.24, we get an isomorphism between $M = \bigoplus_{p \in P} M_p$ and $N = \bigoplus_{p \in P} N_p$. \square

In view of Proposition 24.25, the direct sum of Theorem 24.24 in terms of its p -primary components is called the *canonical primary decomposition* of M .

If M is a finitely generated torsion-module, then Theorem 24.24 takes the following form.

Theorem 24.26. *(Primary Decomposition Theorem for finitely generated torsion modules) Let M be a finitely generated torsion-module over a PID A . If $\text{Ann}(M) = (a)$ and if $a = up_1^{n_1} \cdots p_r^{n_r}$ is a factorization of a into prime factors, then M is the finite direct sum*

$$M = \bigoplus_{i=1}^r M(p_i^{n_i}).$$

Furthermore, the projection of M over $M(p_i^{n_i})$ is of the form $x \mapsto \gamma_i x$, for some $\gamma_i \in A$.

Proof. This is an immediate consequence of Proposition 24.23. \square

In particular, Theorem 24.26 applies when $A = \mathbb{Z}$. In this case, M is a finitely generated torsion abelian group, and the theorem says that such a group is the direct sum of a finite number of groups whose elements have order some power of a prime number p .

Theorem 24.24 has several useful corollaries.

Proposition 24.27. *If M is a torsion module over a PID, for every submodule N of M , we have a direct sum*

$$N = \bigoplus_{p \in P} N \cap M_p.$$

Proof. It is easily verified that $N \cap M_p$ is the p -primary component of N . \square

Proposition 24.28. *If M is a torsion module over a PID, a submodule N of M is a direct factor of M iff N_p is a direct factor of M_p for every irreducible element $p \in A$.*

Proof. This is because if N and N' are two submodules of M , we have $M = N \oplus N'$ iff, by Proposition 24.27, $M_p = N_p \oplus N'_p$ for every irreducible elements $p \in A$. \square

An A -module M is said to be *semi-simple* iff for every submodule N of M , there is some submodule N' of M such that $M = N \oplus N'$.

Proposition 24.29. *Let A be a PID which is not a field, and let M be any A -module. Then, M is semi-simple iff it is a torsion module and if $M_p = M(p)$ for every irreducible element $p \in A$ (in other words, if $x \in M$ is annihilated by a power of p , then it is already annihilated by p).*

Proof. Assume that M is semi-simple. Let $x \in M$ and pick any irreducible element $p \in A$. Then, the submodule pAx has a supplement N such that

$$M = pAx \oplus N,$$

so we can write $x = pax + y$, for some $y \in N$ and some $a \in A$. But then,

$$y = (1 - pa)x,$$

and since p is irreducible, p is not a unit, so $1 - pa \neq 0$. Observe that

$$p(1 - ap)x = py \in pAx \cap N = (0).$$

Since $p(1 - ap) \neq 0$, x is a torsion element, and thus M is a torsion module. The above argument shows that

$$p(1 - ap)x = 0,$$

which implies that $px = ap^2x$, and by induction,

$$px = a^n p^{n+1}x, \quad \text{for all } n \geq 1.$$

If we pick x in M_p , then there is some $m \geq 1$ such that $p^m x = 0$, and we conclude that

$$px = 0.$$

Therefore, $M_p = M(p)$, as claimed.

Conversely, assume that M is a torsion-module and that $M_p = M(p)$ for every irreducible element $p \in A$. By Proposition 24.28, it is sufficient to prove that a module annihilated by a an irreducible element is semi-simple. This is because such a module is a vector space over the field $A/(p)$ (recall that in a PID, an ideal (p) is maximal iff p is irreducible), and in a vector space, every subspace has a supplement. \square

Theorem 24.26 shows that a finitely generated torsion module is a direct sum of p -primary modules M_p . We can do better. In the next section, we show that each primary module M_p is the direct sum of cyclic modules of the form $A/(p^n)$.

24.7 Finitely Generated Modules over a PID; Invariant Factor Decomposition

There are several ways of obtaining the decomposition of a finitely generated module as a direct sum of cyclic modules. One way to proceed is to first use the Primary Decomposition Theorem and then to show how each primary module M_p is the direct sum of cyclic modules of the form $A/(p^n)$. This is the approach followed by Lang [67] (Chapter III, section 7), among others. We prefer to use a proposition that produces a particular basis for a submodule of a finitely generated free module, because it yields more information. This is the approach followed in Dummitt and Foote [32] (Chapter 12) and Bourbaki [15] (Chapter VII). The proof that we present is due to Pierre Samuel.

Proposition 24.30. *Let F be a finitely generated free module over a PID A , and let M be any submodule of F . Then, M is a free module and there is a basis (e_1, \dots, e_n) of F , some $q \leq n$, and some nonzero elements $a_1, \dots, a_q \in A$, such that $(a_1 e_1, \dots, a_q e_q)$ is a basis of M and a_i divides a_{i+1} for all i , with $1 \leq i \leq q - 1$.*

Proof. The proposition is trivial when $M = \{0\}$, thus assume that M is nontrivial. Pick some basis (u_1, \dots, u_n) for F . Let $L(F, A)$ be the set of linear forms on F . For any $f \in L(F, A)$, it is immediately verified that $f(M)$ is an ideal in A . Thus, $f(M) = a_h A$, for some $a_h \in A$, since every ideal in A is a principal ideal. Since A is a PID, any nonempty family of ideals in A has a maximal element, so let f be a linear map such that $a_h A$ is a maximal ideal in A . Let $\pi_i: F \rightarrow A$ be the i -th projection, i.e., π_i is defined such that $\pi_i(x_1 u_1 + \dots + x_n u_n) = x_i$.

It is clear that π_i is a linear map, and since M is nontrivial, one of the $\pi_i(M)$ is nontrivial, and $a_h \neq 0$. There is some $e' \in M$ such that $f(e') = a_h$.

We claim that, for every $g \in L(F, A)$, the element $a_h \in A$ divides $g(e')$.

Indeed, if d is the gcd of a_h and $g(e')$, by the Bézout identity, we can write

$$d = ra_h + sg(e'),$$

for some $r, s \in A$, and thus

$$d = rf(e') + sg(e') = (rf + sg)(e').$$

However, $rf + sg \in L(F, A)$, and thus,

$$a_h A \subseteq dA \subseteq (rf + sg)(M),$$

since d divides a_h , and by maximality of $a_h A$, we must have $a_h A = dA$, which implies that $d = a_h$, and thus, a_h divides $g(e')$. In particular, a_h divides each $\pi_i(e')$ and let $\pi_i(e') = a_h b_i$, with $b_i \in A$.

Let $e = b_1 u_1 + \cdots + b_n u_n$. Note that

$$e' = \pi_1(e')u_1 + \cdots + \pi_n(e')u_n = a_h b_1 u_1 + \cdots + a_h b_n u_n,$$

and thus, $e' = a_h e$. Since $a_h = f(e') = f(a_h e) = a_h f(e)$, and since $a_h \neq 0$, we must have $f(e) = 1$.

Next, we claim that

$$F = Ae \oplus f^{-1}(0)$$

and

$$M = Ae' \oplus (M \cap f^{-1}(0)),$$

with $e' = a_h e$.

Indeed, every $x \in F$ can be written as

$$x = f(x)e + (x - f(x)e),$$

and since $f(e) = 1$, we have $f(x - f(x)e) = f(x) - f(x)f(e) = f(x) - f(x) = 0$. Thus, $F = Ae + f^{-1}(0)$. Similarly, for any $x \in M$, we have $f(x) = ra_h$, for some $r \in A$, and thus,

$$x = f(x)e + (x - f(x)e) = ra_h e + (x - f(x)e) = re' + (x - f(x)e),$$

we still have $x - f(x)e \in f^{-1}(0)$, and clearly, $x - f(x)e = x - ra_h e = x - re' \in M$, since $e' \in M$. Thus, $M = Ae' + (M \cap f^{-1}(0))$.

To prove that we have a direct sum, it is enough to prove that $Ae \cap f^{-1}(0) = \{0\}$. For any $x = re \in Ae$, if $f(x) = 0$, then $f(re) = rf(e) = r = 0$, since $f(e) = 1$ and, thus, $x = 0$. Therefore, the sums are direct sums.

We can now prove that M is a free module by induction on the size, q , of a maximal linearly independent family for M .

If $q = 0$, the result is trivial. Otherwise, since

$$M = Ae' \oplus (M \cap f^{-1}(0)),$$

it is clear that $M \cap f^{-1}(0)$ is a submodule of F and that every maximal linearly independent family in $M \cap f^{-1}(0)$ has at most $q - 1$ elements. By the induction hypothesis, $M \cap f^{-1}(0)$ is a free module, and by adding e' to a basis of $M \cap f^{-1}(0)$, we obtain a basis for M , since the sum is direct.

The second part is shown by induction on the dimension n of F .

The case $n = 0$ is trivial. Otherwise, since

$$F = Ae \oplus f^{-1}(0),$$

and since, by the previous argument, $f^{-1}(0)$ is also free, $f^{-1}(0)$ has dimension $n - 1$. By the induction hypothesis applied to its submodule $M \cap f^{-1}(0)$, there is a basis (e_2, \dots, e_n) of $f^{-1}(0)$, some $q \leq n$, and some nonzero elements $a_2, \dots, a_q \in A$, such that, (a_2e_2, \dots, a_qe_q) is a basis of $M \cap f^{-1}(0)$, and a_i divides a_{i+1} for all i , with $2 \leq i \leq q - 1$. Let $e_1 = e$, and $a_1 = a_h$, as above. It is clear that (e_1, \dots, e_n) is a basis of F , and that (a_1e_1, \dots, a_qe_q) is a basis of M , since the sums are direct, and $e' = a_1e_1 = a_h e$. It remains to show that a_1 divides a_2 . Consider the linear map $g: F \rightarrow A$ such that $g(e_1) = g(e_2) = 1$, and $g(e_i) = 0$, for all i , with $3 \leq i \leq n$. We have $a_h = a_1 = g(a_1e_1) = g(e') \in g(M)$, and thus $a_h A \subseteq g(M)$. Since $a_h A$ is maximal, we must have $g(M) = a_h A = a_1 A$. Since $a_2 = g(a_2e_2) \in g(M)$, we have $a_2 \in a_1 A$, which shows that a_1 divides a_2 . \square

We need the following basic proposition.

Proposition 24.31. *For any commutative ring A , if F is a free A -module and if (e_1, \dots, e_n) is a basis of F , for any elements $a_1, \dots, a_n \in A$, there is an isomorphism*

$$F/(Aa_1e_1 \oplus \dots \oplus Aa_ne_n) \approx (A/a_1A) \oplus \dots \oplus (A/a_nA).$$

Proof. Let $\sigma: F \rightarrow A/(a_1A) \oplus \dots \oplus A/(a_nA)$ be the linear map given by

$$\sigma(x_1e_1 + \dots + x_ne_n) = (\bar{x}_1, \dots, \bar{x}_n),$$

where \bar{x}_i is the equivalence class of x_i in A/a_iA . The map σ is clearly surjective, and its kernel consists of all vectors $x_1e_1 + \dots + x_ne_n$ such that $x_i \in a_iA$, for $i = 1, \dots, n$, which means that

$$\text{Ker}(\sigma) = Aa_1e_1 \oplus \dots \oplus Aa_ne_n.$$

Since $M/\text{Ker}(\sigma)$ is isomorphic to $\text{Im}(\sigma)$, we get the desired isomorphism. \square

We can now prove the existence part of the structure theorem for finitely generated modules over a PID.

Theorem 24.32. *Let M be a finitely generated nontrivial A -module, where A a PID. Then, M is isomorphic to a direct sum of cyclic modules*

$$M \approx A/\mathfrak{a}_1 \oplus \cdots \oplus A/\mathfrak{a}_m,$$

where the \mathfrak{a}_i are proper ideals of A (possibly zero) such that

$$\mathfrak{a}_1 \subseteq \mathfrak{a}_2 \subseteq \cdots \subseteq \mathfrak{a}_m \neq A.$$

More precisely, if $\mathfrak{a}_1 = \cdots = \mathfrak{a}_r = (0)$ and $(0) \neq \mathfrak{a}_{r+1} \subseteq \cdots \subseteq \mathfrak{a}_m \neq A$, then

$$M \approx A^r \oplus (A/\mathfrak{a}_{r+1} \oplus \cdots \oplus A/\mathfrak{a}_m),$$

where $A/\mathfrak{a}_{r+1} \oplus \cdots \oplus A/\mathfrak{a}_m$ is the torsion submodule of M . The module M is free iff $r = m$, and a torsion-module iff $r = 0$. In the latter case, the annihilator of M is \mathfrak{a}_1 .

Proof. Since M is finitely generated and nontrivial, there is a surjective homomorphism $\varphi: A^n \rightarrow M$ for some $n \geq 1$, and M is isomorphic to $A^n/\text{Ker}(\varphi)$. Since $\text{Ker}(\varphi)$ is a submodule of the free module A^n , by Proposition 24.30, $\text{Ker}(\varphi)$ is a free module and there is a basis (e_1, \dots, e_n) of A^n and some nonzero elements a_1, \dots, a_q ($q \leq n$) such that $(a_1 e_1, \dots, a_q e_q)$ is a basis of $\text{Ker}(\varphi)$ and $a_1 \mid a_2 \mid \cdots \mid a_q$. Let $a_{q+1} = \dots = a_n = 0$.

By Proposition 24.31, we have an isomorphism

$$A^n/\text{Ker}(\varphi) \approx A/a_1 A \oplus \cdots \oplus A/a_n A.$$

Whenever a_i is unit, the factor $A/a_i A = (0)$, so we can weed out the units. Let $r = n - q$, and let $s \in \mathbb{N}$ be the smallest index such that a_{s+1} is not a unit. Note that $s = 0$ means that there are no units. Also, as $M \neq (0)$, $s < n$. Then,

$$M \approx A^n/\text{Ker}(\varphi) \approx A/a_{s+1} A \oplus \cdots \oplus A/a_n A.$$

Let $m = r + q - s = n - s$. Then, we have the sequence

$$\underbrace{a_{s+1}, \dots, a_q}_{q-s}, \underbrace{a_{q+1}, \dots, a_n}_{r=n-q},$$

where $a_{s+1} \mid a_{s+2} \mid \cdots \mid a_q$ are nonzero and nonunits and $a_{q+1} = \cdots = a_n = 0$, so we define the m ideals \mathfrak{a}_i as follows:

$$\mathfrak{a}_i = \begin{cases} (0) & \text{if } 1 \leq i \leq r \\ a_{r+q+1-i} A & \text{if } r+1 \leq i \leq m. \end{cases}$$

With these definitions, the ideals \mathfrak{a}_i are proper ideals and we have

$$\mathfrak{a}_i \subseteq \mathfrak{a}_{i+1}, \quad i = 1, \dots, m-1.$$

When $r = 0$, since $a_{s+1} \mid a_{s+2} \mid \cdots \mid a_n$, it is clear that $\mathfrak{a}_1 = a_n A$ is the annihilator of M . The other statements of the theorem are clear. \square

The natural number r is called the *free rank* or *Betti number* of the module M . The generators $\alpha_1, \dots, \alpha_m$ of the ideals $\mathfrak{a}_1, \dots, \mathfrak{a}_m$ (defined up to a unit) are often called the *invariant factors* of M (in the notation of Theorem 24.32, the generators of the ideals $\mathfrak{a}_1, \dots, \mathfrak{a}_m$ are denoted by a_q, \dots, a_{s+1} , $s \leq q$).

As corollaries of Theorem 24.32, we obtain again the following facts established in Section 24.1:

1. A finitely generated module over a PID is the direct sum of its torsion module and a free module.
2. A finitely generated torsion-free module over a PID is free.

It turns out that the ideals $\mathfrak{a}_1 \subseteq \mathfrak{a}_2 \subseteq \dots \subseteq \mathfrak{a}_m \neq A$ are uniquely determined by the module M . Uniqueness proofs found in most books tend to be intricate and not very intuitive. The shortest proof that we are aware of is from Bourbaki [15] (Chapter VII, Section 4), and uses wedge products.

The following preliminary results are needed.

Proposition 24.33. *If A is a commutative ring and if $\mathfrak{a}_1, \dots, \mathfrak{a}_m$ are ideals of A , then there is an isomorphism*

$$A/\mathfrak{a}_1 \otimes \dots \otimes A/\mathfrak{a}_m \approx A/(\mathfrak{a}_1 + \dots + \mathfrak{a}_m).$$

Sketch of proof. We proceed by induction on m . For $m = 2$, we define the map $\varphi: A/\mathfrak{a}_1 \times A/\mathfrak{a}_2 \rightarrow A/(\mathfrak{a}_1 + \mathfrak{a}_2)$ by

$$\varphi(\bar{a}, \bar{b}) = ab \pmod{\mathfrak{a}_1 + \mathfrak{a}_2}.$$

It is well-defined because if $a' = a + a_1$ and $b' = b + a_2$ with $a_1 \in \mathfrak{a}_1$ and $a_2 \in \mathfrak{a}_2$, then

$$a'b' = (a + a_1)(b + a_2) = ab + ba_1 + aa_2 + a_1a_2,$$

and so

$$a'b' \equiv ab \pmod{\mathfrak{a}_1 + \mathfrak{a}_2}.$$

It is also clear that this map is bilinear, so it induces a linear map $\varphi: A/\mathfrak{a}_1 \otimes A/\mathfrak{a}_2 \rightarrow A/(\mathfrak{a}_1 + \mathfrak{a}_2)$ such that $\varphi(\bar{a} \otimes \bar{b}) = ab \pmod{\mathfrak{a}_1 + \mathfrak{a}_2}$.

Next, observe that any arbitrary tensor

$$\bar{a}_1 \otimes \bar{b}_1 + \dots + \bar{a}_n \otimes \bar{b}_n$$

in $A/\mathfrak{a}_1 \otimes A/\mathfrak{a}_2$ can be rewritten as

$$\bar{1} \otimes (\overline{a_1 b_1} + \dots + \overline{a_n b_n}),$$

which is of the form $\bar{1} \otimes \bar{s}$, with $s \in A$. We can use this fact to show that φ is injective and surjective, and thus an isomorphism.

For example, if $\varphi(\bar{1} \otimes \bar{s}) = 0$, because $\varphi(\bar{1} \otimes \bar{s}) = s \pmod{\mathfrak{a}_1 + \mathfrak{a}_2}$, we have $s \in \mathfrak{a}_1 + \mathfrak{a}_2$, so we can write $s = a + b$ with $a \in \mathfrak{a}_1$ and $b \in \mathfrak{a}_2$. Then

$$\begin{aligned}\bar{1} \otimes \bar{s} &= \bar{1} \otimes \overline{a + b} \\ &= \bar{1} \otimes (\bar{a} + \bar{b}) \\ &= \bar{1} \otimes \bar{a} + \bar{1} \otimes \bar{b} \\ &= \bar{a} \otimes \bar{1} + \bar{1} \otimes \bar{b} \\ &= 0 + 0 = 0,\end{aligned}$$

since $a \in \mathfrak{a}_1$ and $b \in \mathfrak{a}_2$, which proves injectivity. \square

Recall that the exterior algebra of an A -module M is defined by

$$\bigwedge M = \bigoplus_{k \geq 0} \bigwedge^k(M).$$

Proposition 24.34. *If A is a commutative ring, then for any n modules M_i , there is an isomorphism*

$$\bigwedge \left(\bigoplus_{i=1}^n M_i \right) \approx \bigotimes_{i=1}^n \bigwedge M_i.$$

A proof can be found in Bourbaki [14] (Chapter III, Section 7, No 7, Proposition 10).

Proposition 24.35. *Let A be a commutative ring and let $\mathfrak{a}_1, \dots, \mathfrak{a}_n$ be n ideals of A . If the module M is the direct sum of n cyclic modules*

$$M = A/\mathfrak{a}_1 \oplus \cdots \oplus A/\mathfrak{a}_n,$$

then for every $p > 0$, the exterior power $\bigwedge^p M$ is isomorphic to the direct sum of the modules A/\mathfrak{a}_H , where H ranges over all subsets $H \subseteq \{1, \dots, n\}$ with p elements, and with

$$\mathfrak{a}_H = \sum_{h \in H} \mathfrak{a}_h.$$

Proof. If u_i is the image of 1 in A/\mathfrak{a}_i , then A/\mathfrak{a}_i is equal to Au_i . By Proposition 24.34, we have

$$\bigwedge M \approx \bigotimes_{i=1}^n \bigwedge(Au_i).$$

We also have

$$\bigwedge(Au_i) = \bigoplus_{k \geq 0} \bigwedge^k(Au_i) \approx A \oplus Au_i,$$

since $au_i \wedge bu_i = 0$, and it follows that

$$\bigwedge^p M \approx \bigoplus_{\substack{H \subseteq \{1, \dots, n\} \\ |H|=p}} (Au_{k_1}) \otimes \cdots \otimes (Au_{k_p}).$$

However, by Proposition 24.33, we have

$$(Au_{k_1}) \otimes \cdots \otimes (Au_{k_p}) = A/\mathfrak{a}_{k_1} \otimes \cdots \otimes A/\mathfrak{a}_{k_p} \approx A/(\mathfrak{a}_{k_1} + \cdots + \mathfrak{a}_{k_p}) = A/\mathfrak{a}_H.$$

Therefore,

$$\bigwedge^p M \approx \bigoplus_{\substack{H \subseteq \{1, \dots, n\} \\ |H|=p}} A/\mathfrak{a}_H,$$

as claimed. \square

When the ideals \mathfrak{a}_i form a chain of inclusions $\mathfrak{a}_1 \subseteq \cdots \subseteq \mathfrak{a}_n$, we get the following remarkable result.

Proposition 24.36. *Let A be a commutative ring and let $\mathfrak{a}_1, \dots, \mathfrak{a}_n$ be n ideals of A such that $\mathfrak{a}_1 \subseteq \mathfrak{a}_2 \subseteq \cdots \subseteq \mathfrak{a}_n$. If the module M is the direct sum of n cyclic modules*

$$M = A/\mathfrak{a}_1 \oplus \cdots \oplus A/\mathfrak{a}_n,$$

then for every p with $1 \leq p \leq n$, the ideal \mathfrak{a}_p is the annihilator of the exterior power $\bigwedge^p M$. If $\mathfrak{a}_n \neq A$, then $\bigwedge^p M \neq (0)$ for $p = 1, \dots, n$, and $\bigwedge^p M = (0)$ for $p > n$.

Proof. With the notation of Proposition 24.35, we have $\mathfrak{a}_H = \mathfrak{a}_{\max(H)}$, where $\max(H)$ is the greatest element in the set H . Since $\max(H) \geq p$ for any subset with p elements and since $\max(H) = p$ when $H = \{1, \dots, p\}$, we see that

$$\mathfrak{a}_p = \bigcap_{\substack{H \subseteq \{1, \dots, n\} \\ |H|=p}} \mathfrak{a}_H.$$

By Proposition 24.35, we have

$$\bigwedge^p M \approx \bigoplus_{\substack{H \subseteq \{1, \dots, n\} \\ |H|=p}} A/\mathfrak{a}_H$$

which proves that \mathfrak{a}_p is indeed the annihilator of $\bigwedge^p M$. The rest is clear. \square

Proposition 24.36 immediately implies the following crucial fact.

Proposition 24.37. *Let A be a commutative ring and let $\mathfrak{a}_1, \dots, \mathfrak{a}_m$ be m ideals of A and $\mathfrak{a}'_1, \dots, \mathfrak{a}'_n$ be n ideals of A such that $\mathfrak{a}_1 \subseteq \mathfrak{a}_2 \subseteq \dots \subseteq \mathfrak{a}_m \neq A$ and $\mathfrak{a}'_1 \subseteq \mathfrak{a}'_2 \subseteq \dots \subseteq \mathfrak{a}'_n \neq A$. If we have an isomorphism*

$$A/\mathfrak{a}_1 \oplus \dots \oplus A/\mathfrak{a}_m \approx A/\mathfrak{a}'_1 \oplus \dots \oplus A/\mathfrak{a}'_n,$$

then $m = n$ and $\mathfrak{a}_i = \mathfrak{a}'_i$ for $i = 1, \dots, n$.

Proposition 24.37 yields the uniqueness of the decomposition in Theorem 24.32.

Theorem 24.38. (*Invariant Factors Decomposition*) *Let M be a finitely generated nontrivial A -module, where A a PID. Then, M is isomorphic to a direct sum of cyclic modules*

$$M \approx A/\mathfrak{a}_1 \oplus \dots \oplus A/\mathfrak{a}_m,$$

where the \mathfrak{a}_i are proper ideals of A (possibly zero) such that

$$\mathfrak{a}_1 \subseteq \mathfrak{a}_2 \subseteq \dots \subseteq \mathfrak{a}_m \neq A.$$

More precisely, if $\mathfrak{a}_1 = \dots = \mathfrak{a}_r = (0)$ and $(0) \neq \mathfrak{a}_{r+1} \subseteq \dots \subseteq \mathfrak{a}_m \neq A$, then

$$M \approx A^r \oplus (A/\mathfrak{a}_{r+1} \oplus \dots \oplus A/\mathfrak{a}_m),$$

where $A/\mathfrak{a}_{r+1} \oplus \dots \oplus A/\mathfrak{a}_m$ is the torsion submodule of M . The module M is free iff $r = m$, and a torsion-module iff $r = 0$. In the latter case, the annihilator of M is \mathfrak{a}_1 . Furthermore, the integer r and ideals $\mathfrak{a}_1 \subseteq \mathfrak{a}_2 \subseteq \dots \subseteq \mathfrak{a}_m \neq A$ are uniquely determined by M .

Proof. By Theorem 24.7, since $M_{\text{tor}} = A/\mathfrak{a}_{r+1} \oplus \dots \oplus A/\mathfrak{a}_m$, we know that the dimension r of the free summand only depends on M . The uniqueness of the sequence of ideals follows from Proposition 24.37. \square

In view of the uniqueness part of Theorem 24.38, we make the following definition.

Definition 24.10. Given a finitely generated module M over a PID A as in Theorem 24.38, the ideals $\mathfrak{a}_i = \alpha_i A$ are called the *invariant factors* of M . The generators α_i of these ideals (uniquely defined up to a unit) are also called the *invariant factors* of M .

Proposition 24.30 can be sharpened as follows:

Proposition 24.39. *Let F be a finitely generated free module over a PID A , and let M be any submodule of F . Then, M is a free module and there is a basis (e_1, \dots, e_n) of F , some $q \leq n$, and some nonzero elements $a_1, \dots, a_q \in A$, such that $(a_1 e_1, \dots, a_q e_q)$ is a basis of M and a_i divides a_{i+1} for all i , with $1 \leq i \leq q-1$. Furthermore, the free module M' with basis (e_1, \dots, e_q) and the ideals $a_1 A, \dots, a_q A$ are uniquely determined by M ; the quotient module M'/M is the torsion module of F/M , and we have an isomorphism*

$$M'/M \approx A/a_1 A \oplus \dots \oplus A/a_q A.$$

Proof. Since $a_i \neq 0$ for $i = 1, \dots, q$, observe that

$$M' = \{x \in F \mid (\exists \beta \in A, \beta \neq 0)(\beta x \in M)\},$$

which shows that M'/M is the torsion module of F/M . Therefore, M' is uniquely determined. Since

$$M = Aa_1e_1 \oplus \cdots \oplus Aa_qe_q,$$

by Proposition 24.31 we have an isomorphism

$$M'/M \approx A/a_1A \oplus \cdots \oplus A/a_qA.$$

Now, it is possible that the first s elements a_i are units, in which case $A/a_iA = (0)$, so we can eliminate such factors and we get

$$M'/M \approx A/a_{s+1}A \oplus \cdots \oplus A/a_qA,$$

with $a_qA \subseteq a_{q-1}A \subseteq \cdots \subseteq a_{s+1}A \neq A$. By Proposition 24.37, $q - s$ and the ideals a_jA are uniquely determined for $j = s + 1, \dots, q$, and since $a_1A = \cdots = a_sA = A$, the q ideals a_iA are uniquely determined. \square

The ideals a_1A, \dots, a_qA of Proposition 24.39 are called the *invariant factors of M with respect to F* . They *should not be confused* with the invariant factors of a module M .

It turns out that a_1, \dots, a_q can also be computed in terms of gcd's of minors of a certain matrix. Recall that if X is an $m \times n$ matrix, then a $k \times k$ minor of X is the determinant of any $k \times k$ matrix obtained by picking k columns of X , and then k rows from these k columns.

Proposition 24.40. *Let F be a free module of finite dimension over a PID, (u_1, \dots, u_n) be a basis of F , M be a submodule of F , and (x_1, \dots, x_m) be a set of generators of M . If a_1A, \dots, a_qA are the invariant factors of M with respect to F as in Proposition 24.39, then for $k = 1, \dots, q$, the product $a_1 \cdots a_k$ is a gcd of the $k \times k$ minors of the $n \times m$ matrix X whose columns are the coordinates of the x_j over the u_i .*

Proof. Proposition 24.30 shows that $M \subseteq a_1F$. Consequently, the coordinates of any element of M are multiples of a_1 . On the other hand, we know that there is a linear form f for which a_1A is a maximal ideal and some $e' \in M$ such that $f(e') = a_1$. If we write e' as a linear combination of the x_i , we see that a_1 belongs to the ideal spanned by the coordinates of the x_i over the basis (u_1, \dots, u_n) . Since these coordinates are all multiples of a_1 , it follows that a_1 is their gcd, which proves the case $k = 1$.

For any $k \geq 2$, consider the exterior power $\bigwedge^k M$. Using the notation of the proof of Proposition 24.30, the module M has the basis (a_1e_1, \dots, a_qe_q) , so $\bigwedge^k M$ has a basis consisting of elements of the form

$$a_{i_1}e_{i_1} \wedge \cdots \wedge a_{i_k}e_{i_k} = a_{i_1} \cdots a_{i_k} e_{i_1} \wedge \cdots \wedge e_{i_k},$$

for all sequences (i_1, \dots, i_k) such that $1 \leq i_1 < i_2 < \dots < i_k \leq q$. However, the vectors $e_{i_1} \wedge \dots \wedge e_{i_k}$ form a basis of $\bigwedge^k F$. Thus, the map from $\bigwedge^k M$ into $\bigwedge^k F$ induced by the inclusion $M \subseteq F$ defines an isomorphism of $\bigwedge^k M$ onto the submodule of $\bigwedge^k F$ having the elements $a_{i_1} \dots a_{i_k} e_{i_1} \wedge \dots \wedge e_{i_k}$ as a basis. Since a_j is a multiple of the a_i for $i < j$, the products $a_{i_1} \dots a_{i_k}$ are all multiples of $\delta_k = a_1 \dots a_k$, and one of these is equal to δ_k . The reasoning used for $k = 1$ shows that δ_k is a gcd of the set of coordinates of any spanning set of $\bigwedge^k M$ over any basis of $\bigwedge^k F$. If we pick as basis of $\bigwedge^k F$ the wedge products $u_{i_1} \wedge \dots \wedge u_{i_k}$, and as generators of $\bigwedge^k M$ the wedge products $x_{i_1} \wedge \dots \wedge x_{i_k}$, it is easy to see that the coordinates of the $x_{i_1} \wedge \dots \wedge x_{i_k}$ are indeed determinants which are the $k \times k$ minors of the matrix X . \square

Proposition 24.40 yields a_1, \dots, a_q (up to units) as follows: First, a_1 is a gcd of the entries in X . Having computed a_1, \dots, a_k , let $b_k = a_1 \dots a_k$, compute $b_{k+1} = a_1 \dots a_k a_{k+1}$ as a gcd of all the $(k+1) \times (k+1)$ minors of X , and then a_{k+1} is obtained by dividing b_{k+1} by b_k (recall that a PID is an integral domain).

We also have the following interesting result about linear maps between free modules over a PID.

Proposition 24.41. *Let A be a PID, let F be a free module of dimension n , F' be a free module of dimension m , and $f: F \rightarrow F'$ be a linear map from F to F' . Then, there exist a basis (e_1, \dots, e_n) of F , a basis (e'_1, \dots, e'_m) of F' , and some nonzero elements $\alpha_1, \dots, \alpha_r \in A$ such that*

$$f(e_i) = \begin{cases} \alpha_i e'_i & \text{if } 1 \leq i \leq r \\ 0 & \text{if } r+1 \leq i \leq n, \end{cases}$$

and $\alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_r$. Furthermore, the ideals $\alpha_1 A, \dots, \alpha_r A$ are the invariant factors of $f(F)$ with respect F' .

Proof. Let F_0 be the kernel of f . Since $M' = f(F)$ is a submodule of the free module F' , it is free, and similarly F_0 is free as a submodule of the free module F (by Proposition 24.30). By Proposition 24.2, we have

$$F = F_0 \oplus F_1,$$

where F_1 is a free module, and the restriction of f to F_1 is an isomorphism onto $M' = f(F)$. Proposition 24.39 applied to F' and M' yields a basis (e'_1, \dots, e'_m) of F' such that $(\alpha_1 e'_1, \dots, \alpha_r e'_r)$ is a basis of M' , where $\alpha_1 A, \dots, \alpha_r A$ are the invariant factors for M' with respect to F' . Since the restriction of f to F_1 is an isomorphism, there is a basis (e_1, \dots, e_r) of F_1 such that

$$f(e_i) = \alpha_i e'_i, \quad i = 1, \dots, r.$$

We can extend this basis to a basis of F by picking a basis of F_0 (a free module), which yields the desired result. \square

The matrix version of Proposition 24.41 is the following proposition.

Proposition 24.42. *If X is an $m \times n$ matrix of rank r over a PID A , then there exist some invertible $n \times n$ matrix P , some invertible $m \times m$ matrix Q , and a $m \times n$ matrix D of the form*

$$D = \begin{pmatrix} \alpha_1 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & \alpha_2 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & \alpha_r & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & 0 & \cdots & 0 \end{pmatrix}$$

for some nonzero $\alpha_i \in A$, such that

- (1) $\alpha_1 \mid \alpha_2 \mid \cdots \mid \alpha_r$,
- (2) $X = QDP^{-1}$, and
- (3) The α_i s are uniquely determined up to a unit.

The ideals $\alpha_1 A, \dots, \alpha_r A$ are called the *invariant factors* of the matrix X . Recall that two $m \times n$ matrices X and Y are *equivalent* iff

$$Y = QXP^{-1},$$

for some invertible matrices, P and Q . Then, Proposition 24.42 implies the following fact.

Proposition 24.43. *Two $m \times n$ matrices X and Y are equivalent iff they have the same invariant factors.*

If X is the matrix of a linear map $f: F \rightarrow F'$ with respect to some basis (u_1, \dots, u_n) of F and some basis (u'_1, \dots, u'_m) of F' , then the columns of X are the coordinates of the $f(u_j)$ over the u'_i , where the $f(u_j)$ generate $f(F)$, so Proposition 24.40 applies and yields the following result:

Proposition 24.44. *If X is a $m \times n$ matrix of rank r over a PID A , and if $\alpha_1 A, \dots, \alpha_r A$ are its invariant factors, then α_1 is a gcd of the entries in X , and for $k = 2, \dots, r$, the product $\alpha_1 \cdots \alpha_k$ is a gcd of all $k \times k$ minors of X .*

There are algorithms for converting a matrix X over a PID to the form $X = QDP^{-1}$ as described in Proposition 24.42. For Euclidean domains, this can be achieved by using the elementary row and column operations $P(i, k)$, $E_{i,j;\beta}$, and $E_{i,\lambda}$ described in Chapter 6, where we require the scalar λ used in $E_{i,\lambda}$ to be a unit. For an arbitrary PID, another kind of elementary matrix (containing some 2×2 submatrix in addition to diagonal entries) is needed. These procedures involve computing gcd's and use the Bezout identity to mimic

division. Such methods are presented in Serre [96], Jacobson [59], and Van Der Waerden [112], and sketched in Artin [4]. We describe and justify several of these methods in Section 25.4.

From Section 24.2, we know that a submodule of a finitely generated module over a PID is finitely presented. Therefore, in Proposition 24.39, the submodule M of the free module F is finitely presented by some matrix R with a number of rows equal to the dimension of F . Using Theorem 25.17, the matrix R can be diagonalized as $R = QDP^{-1}$ where D is a diagonal matrix. Then, the columns of Q form a basis (e_1, \dots, e_n) of F , and since $RP = QD$, the nonzero columns of RP form the basis (a_1e_1, \dots, a_qe_q) of M . When the ring A is a Euclidean domain, Theorem 25.14 shows that P and Q are products of elementary row and column operations. In particular, when $A = \mathbb{Z}$, in which cases our \mathbb{Z} -modules are abelian groups, we can find P and Q using Euclidean division.

In this case, a finitely generated submodule M of \mathbb{Z}^n is called a *lattice*. It is given as the set of integral linear combinations of a finite set of integral vectors.

Here is an example taken from Artin [4] (Chapter 12, Section 4). Let F be the free \mathbb{Z} -module \mathbb{Z}^2 , and let M be the lattice generated by the columns of the matrix

$$R = \begin{pmatrix} 2 & -1 \\ 1 & 2 \end{pmatrix}.$$

The columns (u_1, u_2) of R are linearly independent, but they are not a basis of \mathbb{Z}^2 . For example, in order to obtain e_1 as a linear combination of these columns, we would need to solve the linear system

$$\begin{aligned} 2x - y &= 1 \\ x + 2y &= 0. \end{aligned}$$

From the second equation, we get $x = -2y$, which yields

$$-5y = 1.$$

But, $y = -1/5$ is not an integer. We leave it as an exercise to check that

$$\begin{pmatrix} 1 & 0 \\ -3 & 1 \end{pmatrix} \begin{pmatrix} 2 & -1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 5 \end{pmatrix},$$

which means that

$$\begin{pmatrix} 2 & -1 \\ 1 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 3 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 5 \end{pmatrix} \begin{pmatrix} 2 & -1 \\ -1 & 1 \end{pmatrix},$$

so $R = QDP^{-1}$ with

$$Q = \begin{pmatrix} 1 & 0 \\ 3 & 1 \end{pmatrix}, \quad D = \begin{pmatrix} 1 & 0 \\ 0 & 5 \end{pmatrix}, \quad P = \begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix}.$$

The new basis (u'_1, u'_2) for \mathbb{Z}^2 consists of the columns of Q and the new basis for M consists of the columns $(u'_1, 5u'_2)$ of QD , where

$$QD = \begin{pmatrix} 1 & 0 \\ 3 & 5 \end{pmatrix}.$$

A picture of the lattice and its generators (u_1, u_2) and of the same lattice with the new basis $(u'_1, 5u'_2)$ is shown in Figure 24.1, where the lattice points are displayed as stars.

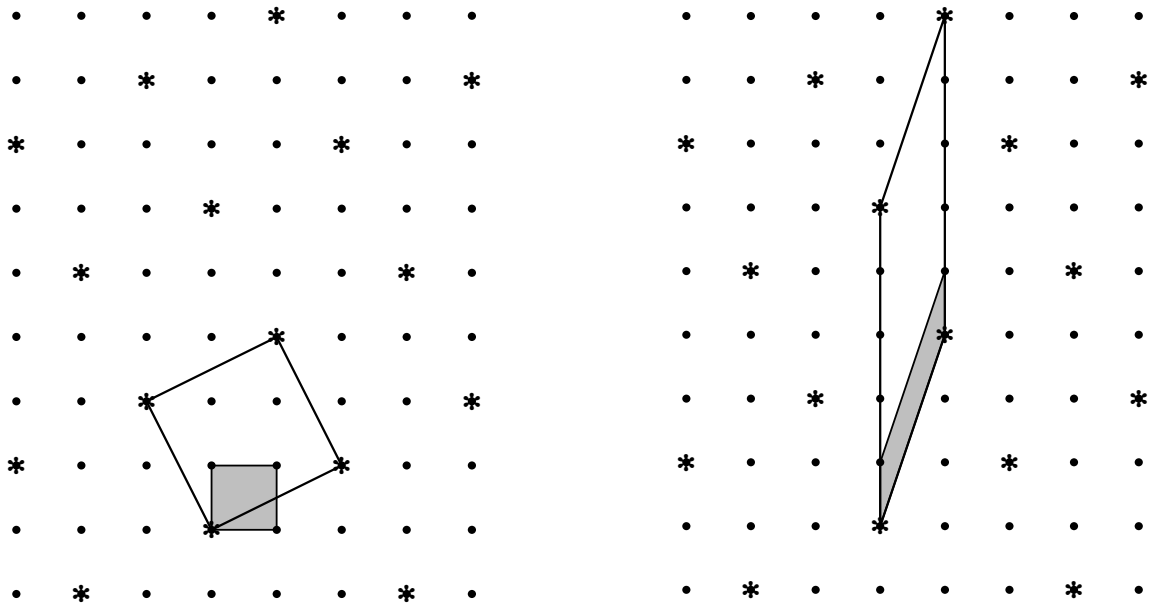


Figure 24.1: Diagonalization applied to a lattice

The invariant factor decomposition of a finitely generated module M over a PID A given by Theorem 24.38 says that

$$M_{\text{tor}} \approx A/\mathfrak{a}_{r+1} \oplus \cdots \oplus A/\mathfrak{a}_m,$$

a direct sum of cyclic modules, with $(0) \neq \mathfrak{a}_{r+1} \subseteq \cdots \subseteq \mathfrak{a}_m \neq A$. Using the Chinese Remainder Theorem (Theorem 21.15), we can further decompose each module $A/\alpha_i A$ into a direct sum of modules of the form $A/p^n A$, where p is a prime in A .

Theorem 24.45. (*Elementary Divisors Decomposition*) *Let M be a finitely generated non-trivial A -module, where A a PID. Then, M is isomorphic to the direct sum $A^r \oplus M_{\text{tor}}$, where A^r is a free module and where the torsion module M_{tor} is a direct sum of cyclic modules of the form $A/p_i^{n_{i,j}}$, for some primes $p_1, \dots, p_t \in A$ and some positive integers $n_{i,j}$, such that for each $i = 1, \dots, t$, there is a sequence of integers*

$$1 \leq \underbrace{n_{i,1}, \dots, n_{i,1}}_{m_{i,1}} < \underbrace{n_{i,2}, \dots, n_{i,2}}_{m_{i,2}} < \cdots < \underbrace{n_{i,s_i}, \dots, n_{i,s_i}}_{m_{i,s_i}},$$

with $s_i \geq 1$, and where $n_{i,j}$ occurs $m_{i,j} \geq 1$ times, for $j = 1, \dots, s_i$. Furthermore, the irreducible elements p_i and the integers $r, t, n_{i,j}, s_i, m_{i,j}$ are uniquely determined.

Proof. By Theorem 24.38, we already know that $M \approx A^r \oplus M_{\text{tor}}$, where r is uniquely determined, and where

$$M_{\text{tor}} \approx A/\mathfrak{a}_{r+1} \oplus \cdots \oplus A/\mathfrak{a}_m,$$

a direct sum of cyclic modules, with $(0) \neq \mathfrak{a}_{r+1} \subseteq \cdots \subseteq \mathfrak{a}_m \neq A$. Then, each \mathfrak{a}_i is a principal ideal of the form $\alpha_i A$, where $\alpha_i \neq 0$ and α_i is not a unit. Using the Chinese Remainder Theorem (Theorem 21.15), if we factor α_i into prime factors as

$$\alpha_i = up_1^{k_1} \cdots p_h^{k_h},$$

with $k_j \geq 1$, we get an isomorphism

$$A/\alpha_i A \approx A/p_1^{k_1} A \oplus \cdots \oplus A/p_h^{k_h} A.$$

This implies that M_{tor} is the direct sum of modules of the form $A/p_i^{n_{i,j}}$, for some primes $p_i \in A$.

To prove uniqueness, observe that the p_i -primary component of M_{tor} is the direct sum

$$(A/p_i^{n_{i,1}} A)^{m_{i,1}} \oplus \cdots \oplus (A/p_i^{n_{i,s_i}} A)^{m_{i,s_i}},$$

and these are uniquely determined. Since $n_{i,1} < \cdots < n_{i,s_i}$, we have

$$p_i^{n_{i,s_i}} A \subseteq \cdots \subseteq p_i^{n_{i,1}} A \neq A,$$

Proposition 24.37 implies that the irreducible elements p_i and $n_{i,j}, s_i$, and $m_{i,j}$ are unique. \square

In view of Theorem 24.45, we make the following definition.

Definition 24.11. Given a finitely generated module M over a PID A as in Theorem 24.45, the ideals $p_i^{n_{i,j}} A$ are called the *elementary divisors* of M , and the $m_{i,j}$ are their *multiplicities*. The ideal (0) is also considered to be an elementary divisor and r is its multiplicity.

Remark: Theorem 24.45 shows how the elementary divisors are obtained from the invariant factors: the elementary divisors are the prime power factors of the invariant factors.

Conversely, we can get the invariant factors from the elementary divisors. We may assume that M is a torsion module. Let

$$m = \max_{1 \leq i \leq t} \{m_{i,1} + \cdots + m_{i,s_i}\},$$

and construct the $t \times m$ matrix $C = (c_{ij})$ whose i th row is the sequence

$$\underbrace{n_{i,s_i}, \dots, n_{i,s_i}}_{m_{i,s_i}}, \dots, \underbrace{n_{i,2}, \dots, n_{i,2}}_{m_{i,2}}, \underbrace{n_{i,1}, \dots, n_{i,1}}_{m_{i,1}}, 0, \dots, 0,$$

padded with 0's if necessary to make it of length m . Then, the j th invariant factor is

$$\alpha_j A = p_1^{c_{1j}} p_2^{c_{2j}} \cdots p_t^{c_{t,j}} A.$$

Observe that because the last column contains at least one prime, the α_i are not units, and $\alpha_m \mid \alpha_{m-1} \mid \cdots \mid \alpha_1$, so that $\alpha_1 A \subseteq \cdots \subseteq \alpha_{m-1} A \subseteq \alpha_m A \neq A$, as desired.

From a computational point of view, finding the elementary divisors is usually practically impossible, because it requires factoring. For example, if $A = K[X]$ where K is a field, such as $K = \mathbb{R}$ or $K = \mathbb{C}$, factoring amounts to finding the roots of a polynomial, but by Galois theory, in general, this is not algorithmically doable. On the other hand, the invariant factors can be computed using elementary row and column operations.

It can also be shown that A and the modules of the form $A/p^n A$ are indecomposable (with $n > 0$). A module M is said to be *indecomposable* if M cannot be written as a direct sum of two nonzero proper submodules of M . For a proof, see Bourbaki [15] (Chapter VII, Section 4, No. 8, Proposition 8). Theorem 24.45 shows that a finitely generated module over a PID is a direct sum of indecomposable modules.

We will now apply the structure theorems for finitely generated (torsion) modules to the $K[X]$ -module E_f associated with an endomorphism f on a vector space E .

Chapter 25

The Rational Canonical Form and Other Normal Forms

25.1 The Rational Canonical Form

Let E be a finite-dimensional vector space over a field K , and let $f: E \rightarrow E$ be an endomorphism of E . We know from Section 24.5 that there is a $K[X]$ -module E_f associated with f , and that M_f is a finitely generated torsion module over the PID $K[X]$. In this chapter, we show how Theorems from Sections 24.6 and 24.7 yield important results about the structure of the linear map f .

Recall that the annihilator of a subspace V is an ideal (p) uniquely defined by a monic polynomial p called the *minimal polynomial* of V .

Our first result is obtained by translating the primary decomposition theorem, Theorem 24.26. It is not too surprising that we obtain again Theorem 22.9!

Theorem 25.1. (*Primary Decomposition Theorem*) Let $f: E \rightarrow E$ be a linear map on the finite-dimensional vector space E over the field K . Write the minimal polynomial m of f as

$$m = p_1^{r_1} \cdots p_k^{r_k},$$

where the p_i are distinct irreducible monic polynomials over K , and the r_i are positive integers. Let

$$W_i = \text{Ker}(p_i(f)^{r_i}), \quad i = 1, \dots, k.$$

Then

- (a) $E = W_1 \oplus \cdots \oplus W_k$.
- (b) Each W_i is invariant under f and the projection from W onto W_i is given by a polynomial in f .
- (c) The minimal polynomial of the restriction $f|W_i$ of f to W_i is $p_i^{r_i}$.

Next, we apply the Invariant Factors Decomposition Theorem, Theorem 24.38, to E_f . This theorem says that E_f is isomorphic to a direct sum

$$E_f \approx K[X]/(p_1) \oplus \cdots \oplus K[X]/(p_m)$$

of $m \leq n$ cyclic modules, where the p_j are uniquely determined monic polynomials of degree at least 1, such that

$$p_m \mid p_{m-1} \mid \cdots \mid p_1.$$

Each cyclic module $K[X]/(p_i)$ is isomorphic to a cyclic subspace for f , say V_i , whose minimal polynomial is p_i .

It is customary to renumber the polynomials p_i as follows. The n polynomials q_1, \dots, q_n are defined by:

$$q_j(X) = \begin{cases} 1 & \text{if } 1 \leq j \leq n-m \\ p_{n-j+1}(X) & \text{if } n-m+1 \leq j \leq n. \end{cases}$$

Then, we see that

$$q_1 \mid q_2 \mid \cdots \mid q_n,$$

where the first $n-m$ polynomials are equal to 1, and we have the direct sum

$$E = E_1 \oplus \cdots \oplus E_n,$$

where E_i is a cyclic subspace for f whose minimal polynomial is q_i . In particular, $E_i = (0)$ for $i = 1, \dots, n-m$. Theorem 24.38 also says that the minimal polynomial of f is $q_n = p_1$. We sum all this up in the following theorem.

Theorem 25.2. (*Cyclic Decomposition Theorem, First Version*) *Let $f: E \rightarrow E$ be an endomorphism on a K -vector space of dimension n . There exist n monic polynomials $q_1, \dots, q_n \in K[X]$ such that*

$$q_1 \mid q_2 \mid \cdots \mid q_n,$$

and E is the direct sum

$$E = E_1 \oplus \cdots \oplus E_n$$

of cyclic subspaces $E_i = Z(u_i; f)$ for f , such that the minimal polynomial of the restriction of f to E_i is q_i . The polynomials q_i satisfying the above conditions are unique, and q_n is the minimal polynomial of f .

In view of translation point (4) at the beginning of Section 24.5, we know that over the basis

$$(u_i, f(u_i), \dots, f^{n_i-1}(u_i))$$

of the cyclic subspace $E_i = Z(u_i; f)$, with $n_i = \deg(q_i)$, the matrix of the restriction of f to E_i is the *companion matrix* of $p_i(X)$, of the form

$$\begin{pmatrix} 0 & 0 & 0 & \cdots & 0 & -a_0 \\ 1 & 0 & 0 & \cdots & 0 & -a_1 \\ 0 & 1 & 0 & \cdots & 0 & -a_2 \\ \vdots & \ddots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \ddots & 0 & -a_{n_i-2} \\ 0 & 0 & 0 & \cdots & 1 & -a_{n_i-1} \end{pmatrix}.$$

If we put all these bases together, we obtain a block matrix whose blocks are of the above form. Therefore, we proved the following result.

Theorem 25.3. (*Rational Canonical Form, First Version*) Let $f: E \rightarrow E$ be an endomorphism on a K -vector space of dimension n . There exist n monic polynomials $q_1, \dots, q_n \in K[X]$ such that

$$q_1 \mid q_2 \mid \cdots \mid q_n,$$

with $q_1 = \cdots = q_{n-m} = 1$, and a basis of E such that the matrix X of f is a block matrix of the form

$$X = \begin{pmatrix} A_{n-m+1} & 0 & \cdots & 0 & 0 \\ 0 & A_{n-m+2} & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & A_{n-1} & 0 \\ 0 & 0 & \cdots & 0 & A_n \end{pmatrix},$$

where each A_i is the companion matrix of q_i . The polynomials q_i satisfying the above conditions are unique, and q_n is the minimal polynomial of f .

Definition 25.1. A matrix X as in Theorem 25.3 is called a matrix in *rational form*. The polynomials q_1, \dots, q_n arising in Theorems 25.2 and 25.3 are called the *similarity invariants* (or *invariant factors*) of f .

Theorem 25.3 shows that every matrix is similar to a matrix in rational form. Such a matrix is unique.

By Proposition 24.20, two linear maps f and f' are similar iff there is an isomorphism between E_f and $E'_{f'}$, and thus by the uniqueness part of Theorem 24.38, iff they have the same similarity invariants q_1, \dots, q_n .

Proposition 25.4. If E and E' are two finite-dimensional vector spaces and if $f: E \rightarrow E$ and $f': E' \rightarrow E'$ are two linear maps, then f and f' are similar iff they have the same similarity invariants.

The effect of extending the field K to a field L is the object of the next proposition.

Proposition 25.5. *Let $f: E \rightarrow E$ be a linear map on a K -vector space E , and let (q_1, \dots, q_n) be the similarity invariants of f . If L is a field extension of K (which means that $K \subseteq L$), and if $E_{(L)} = L \otimes_K E$ is the vector space obtained by extending the scalars, and $f_{(L)} = 1_L \otimes f$ the linear map of $E_{(L)}$ induced by f , then the similarity invariants of $f_{(L)}$ are (q_1, \dots, q_n) viewed as polynomials in $L[X]$.*

Proof. We know that E_f is isomorphic to the direct sum

$$E_f \approx K[X]/(q_1 K[X]) \oplus \cdots \oplus K[X]/(q_n K[X]),$$

so by tensoring with $L[X]$ and using Propositions 24.12 and 23.7, we get

$$\begin{aligned} L[X] \otimes_{K[X]} E_f &\approx L[X] \otimes_{K[X]} (K[X]/(q_1 K[X]) \oplus \cdots \oplus K[X]/(q_n K[X])) \\ &\approx L[X] \otimes_{K[X]} (K[X]/(q_1 K[X])) \oplus \cdots \oplus L[X] \otimes_{K[X]} (K[X]/(q_n K[X])) \\ &\approx (K[X]/(q_1 K[X])) \otimes_{K[X]} L[X] \oplus \cdots \oplus (K[X]/(q_n K[X])) \otimes_{K[X]} L[X]. \end{aligned}$$

However, by Proposition 24.14, we have isomorphisms

$$(K[X]/(q_i K[X])) \otimes_{K[X]} L[X] \approx L[X]/(q_i L[X]),$$

so we get

$$L[X] \otimes_{K[X]} E_f \approx L[X]/(q_1 L[X]) \oplus \cdots \oplus L[X]/(q_n L[X]).$$

Since E_f is a $K[X]$ -module, the $L[X]$ module $L[X] \otimes_{K[X]} E_f$ is the module obtained from E_f by the ring extension $K[X] \subseteq L[X]$, and since f is a $K[X]$ -linear map of E_f , it becomes $f_{(L[X])}$ on $L[X] \otimes_{K[X]} E_f$, which is the same as $f_{(L)}$ viewed as an L -linear map of the space $E_{(L)} = L \otimes_K E$, so $L[X] \otimes_{K[X]} E_f$ is actually isomorphic to $E_{(L)f_{(L)}}$, and we have

$$E_{(L)f_{(L)}} \approx L[X]/(q_1 L[X]) \oplus \cdots \oplus L[X]/(q_n L[X]),$$

which shows that (q_1, \dots, q_n) are the similarity invariants of $f_{(L)}$. □

Proposition justifies the terminology “invariant” in similarity invariants. Indeed, under a field extension $K \subseteq L$, the similarity invariants of $f_{(L)}$ remain the same. This is not true of the elementary divisors, which depend on the field; indeed, an irreducible polynomial $p \in K[X]$ may split over $L[X]$. Since q_n is the minimal polynomial of f , the above reasoning also shows that the minimal polynomial of $f_{(L)}$ remains the same under a field extension.

Proposition 25.5 has the following corollary.

Proposition 25.6. *Let K be a field and let $L \supseteq K$ be a field extension of K . For any two square matrices X and Y over K , if there is an invertible matrix Q over L such that $Y = QXQ^{-1}$, then there is an invertible matrix P over K such that $Y = PXP^{-1}$.*

Recall from Theorem 24.21 that the sequence of $K[X]$ -linear maps

$$0 \longrightarrow E[X] \xrightarrow{\psi} E[X] \xrightarrow{\sigma} E_f \longrightarrow 0$$

is exact, and as a consequence, E_f is isomorphic to the quotient of $E[X]$ by $\text{Im}(X1 - \bar{f})$. Furthermore, because E is a vector space, $E[X]$ is a free module with basis $(1 \otimes u_1, \dots, 1 \otimes u_n)$, where (u_1, \dots, u_n) is a basis of E . By Theorem 24.38, we have an isomorphism

$$E_f \approx K[X]/(q_1 K[X]) \oplus \cdots \oplus K[X]/(q_n K[X]),$$

and by Proposition 24.39, $E[X]/\text{Im}(X1 - \bar{f})$ is isomorphic to a direct sum

$$E[X]/\text{Im}(X1 - \bar{f}) \approx K[X]/(p_1 K[X]) \oplus \cdots \oplus K[X]/(p_m K[X]),$$

where p_1, \dots, p_m are the invariant factors of $\text{Im}(X1 - \bar{f})$ with respect to $E[X]$. Since $E[X] \approx E[X]/\text{Im}(X1 - \bar{f})$, by the uniqueness part of Theorem 24.38 and because the polynomials are monic, we must have $m = n$ and $p_i = q_i$, for $i = 1, \dots, n$. Therefore, we proved the following crucial fact:

Proposition 25.7. *For any linear map $f: E \rightarrow E$ over a K -vector space E of dimension n , the similarity invariants of f are equal to the invariant factors of $\text{Im}(X1 - \bar{f})$ with respect to $E[X]$.*

Proposition 25.7 is the key to computing the similarity invariants of a linear map. This can be done using a procedure to convert $XI - U$ to its *Smith normal form*. Propositions 25.7 and 24.44 yield the following result.

Proposition 25.8. *For any linear map $f: E \rightarrow E$ over a K -vector space E of dimension n , if (q_1, \dots, q_n) are the similarity invariants of f , for any matrix U representing f with respect to any basis, then for $k = 1, \dots, n$ the product*

$$d_k(X) = q_1(X) \cdots q_k(X)$$

is the gcd of the $k \times k$ -minors of the matrix $XI - U$.

Note that the matrix $XI - U$ is nonother than the matrix that yields the characteristic polynomial $\chi_f(X) = \det(XI - U)$ of f .

Proposition 25.9. *For any linear map $f: E \rightarrow E$ over a K -vector space E of dimension n , if (q_1, \dots, q_n) are the similarity invariants of f , then the following properties hold:*

(1) *If $\chi_f(X)$ is the characteristic polynomial of f , then*

$$\chi_f(X) = q_1(X) \cdots q_n(X).$$

- (2) The minimal polynomial $m(X) = q_n(X)$ of f divides the characteristic polynomial $\chi_f(X)$ of f .
- (3) The characteristic polynomial $\chi_f(X)$ divides $m(X)^n$.
- (4) E is cyclic for f iff $m(X) = \chi(X)$.

Proof. Property (1) follows from Proposition 25.8 for $k = n$. It also follows from Theorem 25.3 and the fact that for the companion matrix associated with q_i , the characteristic polynomial of this matrix is also q_i . Property (2) is obvious from (1). Since each q_i divides q_{i+1} , each q_i divides q_n , so their product $\chi_f(X)$ divides $m(X)^n = q_n(X)^n$. The last condition says that $q_1 = \cdots = q_{n-1} = 1$, which means that E_f has a single summand. \square

Observe that Proposition 25.9 yields another proof of the Cayley–Hamilton Theorem. It also implies that a linear map is nilpotent iff its characteristic polynomial is X^n .

25.2 The Rational Canonical Form, Second Version

Let us now translate the Elementary Divisors Decomposition Theorem, Theorem 24.45, in terms of E_f . We obtain the following result.

Theorem 25.10. (*Cyclic Decomposition Theorem, Second Version*) Let $f: E \rightarrow E$ be an endomorphism on a K -vector space of dimension n . Then, E is the direct sum of cyclic subspaces $E_j = Z(u_j; f)$ for f , such that the minimal polynomial of E_j is of the form $p_i^{n_{i,j}}$, for some irreducible monic polynomials $p_1, \dots, p_t \in K[X]$ and some positive integers $n_{i,j}$, such that for each $i = 1, \dots, t$, there is a sequence of integers

$$1 \leq \underbrace{n_{i,1}, \dots, n_{i,1}}_{m_{i,1}} < \underbrace{n_{i,2}, \dots, n_{i,2}}_{m_{i,2}} < \cdots < \underbrace{n_{i,s_i}, \dots, n_{i,s_i}}_{m_{i,s_i}},$$

with $s_i \geq 1$, and where $n_{i,j}$ occurs $m_{i,j} \geq 1$ times, for $j = 1, \dots, s_i$. Furthermore, the monic polynomials p_i and the integers $r, t, n_{i,j}, s_i, m_{i,j}$ are uniquely determined.

Note that there are $\mu = \sum m_{i,j}$ cyclic subspaces $Z(u_j; f)$. Using bases for the cyclic subspaces $Z(u_j; f)$ as in Theorem 25.3, we get the following theorem.

Theorem 25.11. (*Rational Canonical Form, Second Version*) Let $f: E \rightarrow E$ be an endomorphism on a K -vector space of dimension n . There exist t distinct irreducible monic polynomials $p_1, \dots, p_t \in K[X]$ and some positive integers $n_{i,j}$, such that for each $i = 1, \dots, t$, there is a sequence of integers

$$1 \leq \underbrace{n_{i,1}, \dots, n_{i,1}}_{m_{i,1}} < \underbrace{n_{i,2}, \dots, n_{i,2}}_{m_{i,2}} < \cdots < \underbrace{n_{i,s_i}, \dots, n_{i,s_i}}_{m_{i,s_i}},$$

with $s_i \geq 1$, and where $n_{i,j}$ occurs $m_{i,j} \geq 1$ times, for $j = 1, \dots, s_i$, and there is a basis of E such that the matrix X of f is a block matrix of the form

$$X = \begin{pmatrix} A_1 & 0 & \cdots & 0 & 0 \\ 0 & A_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & A_{\mu-1} & 0 \\ 0 & 0 & \cdots & 0 & A_\mu \end{pmatrix},$$

where each A_j is the companion matrix of some $p_i^{n_{i,j}}$, and $\mu = \sum m_{i,j}$. The monic polynomials p_1, \dots, p_t and the integers $r, t, n_{i,j}, s_i, m_{i,j}$ are uniquely determined

The polynomials $p_i^{n_{i,j}}$ are called the *elementary divisors* of f (and X). These polynomials are factors of the minimal polynomial.

As we pointed earlier, unlike the similarity invariants, the elementary divisors may change when we pass to a field extension.

We will now consider the special case where all the irreducible polynomials p_i are of the form $X - \lambda_i$; that is, when are the eigenvalues of f belong to K . In this case, we find again the Jordan form.

25.3 The Jordan Form Revisited

In this section, we assume that all the roots of the minimal polynomial of f belong to K . This will be the case if K is algebraically closed. The irreducible polynomials p_i of Theorem 25.10 are the polynomials $X - \lambda_i$, for the distinct eigenvalues λ_i of f . Then, each cyclic subspace $Z(u_j; f)$ has a minimal polynomial of the form $(X - \lambda)^m$, for some eigenvalue λ of f and some $m \geq 1$. It turns out that by choosing a suitable basis for the cyclic subspace $Z(u_j; f)$, the matrix of the restriction of f to $Z(u_j; f)$ is a Jordan block.

Proposition 25.12. *Let E be a finite-dimensional K -vector space and let $f: E \rightarrow E$ be a linear map. If E is a cyclic $K[X]$ -module and if $(X - \lambda)^n$ is the minimal polynomial of f , then there is a basis of E of the form*

$$((f - \lambda \text{id})^{n-1}(u), (f - \lambda \text{id})^{n-2}(u), \dots, (f - \lambda \text{id})(u), u),$$

for some $u \in E$. With respect to this basis, the matrix of f is the Jordan block

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

Proof. Since E is a cyclic $K[X]$ -module, there is some $u \in E$ so that E is generated by $u, f(u), f^2(u), \dots$, which means that every vector in E is of the form $p(f)(u)$, for some polynomial, $p(X)$. We claim that $u, f(u), \dots, f^{n-2}(u), f^{n-1}(u)$ generate E , which implies that the dimension of E is at most n .

This is because if $p(X)$ is any polynomial of degree at least n , then we can divide $p(X)$ by $(X - \lambda)^n$, obtaining

$$p = (X - \lambda)^n q + r,$$

where $0 \leq \deg(r) < n$, and as $(X - \lambda)^n$ annihilates E , we get

$$p(f)(u) = r(f)(u),$$

which means that every vector of the form $p(f)(u)$ with $p(X)$ of degree $\geq n$ is actually a linear combination of $u, f(u), \dots, f^{n-2}(u), f^{n-1}(u)$.

We claim that the vectors

$$u, (f - \lambda \text{id})(u), \dots, (f - \lambda \text{id})^{n-2}(u), (f - \lambda \text{id})^{n-1}(u)$$

are linearly independent. Indeed, if we had a nontrivial linear combination

$$a_0(f - \lambda \text{id})^{n-1}(u) + a_1(f - \lambda \text{id})^{n-2}(u) + \dots + a_{n-2}(f - \lambda \text{id})(u) + a_{n-1}u = 0,$$

then the polynomial

$$a_0(X - \lambda)^{n-1} + a_1(X - \lambda)^{n-2} + \dots + a_{n-2}(X - \lambda) + a_{n-1}$$

of degree at most $n - 1$ would annihilate E , contradicting the fact that $(X - \lambda)^n$ is the minimal polynomial of f (and thus, of smallest degree). Consequently, as the dimension of E is at most n ,

$$((f - \lambda \text{id})^{n-1}(u), (f - \lambda \text{id})^{n-2}(u), \dots, (f - \lambda \text{id})(u), u),$$

is a basis of E and since $u, f(u), \dots, f^{n-2}(u), f^{n-1}(u)$ span E ,

$$(u, f(u), \dots, f^{n-2}(u), f^{n-1}(u))$$

is also a basis of E .

Let us see how f acts on the basis

$$((f - \lambda \text{id})^{n-1}(u), (f - \lambda \text{id})^{n-2}(u), \dots, (f - \lambda \text{id})(u), u).$$

If we write $f = f - \lambda \text{id} + \lambda \text{id}$, as $(f - \lambda \text{id})^n$ annihilates E , we get

$$f((f - \lambda \text{id})^{n-1}(u)) = (f - \lambda \text{id})^n(u) + \lambda(f - \lambda \text{id})^{n-1}(u) = \lambda(f - \lambda \text{id})^{n-1}(u)$$

and

$$f((f - \lambda \text{id})^k(u)) = (f - \lambda \text{id})^{k+1}(u) + \lambda(f - \lambda \text{id})^k(u), \quad 0 \leq k \leq n - 2.$$

But this means precisely that the matrix of f in this basis is the Jordan block $J_n(\lambda)$. \square

Combining Theorem 25.11 and Proposition 25.12, we obtain a strong version of the Jordan form.

Theorem 25.13. (*Jordan Canonical Form*) *Let E be finite-dimensional K -vector space. The following properties are equivalent:*

- (1) *The eigenvalues of f all belong to K .*
- (2) *There is a basis of E in which the matrix of f is upper (or lower) triangular.*
- (3) *There exist a basis of E in which the matrix A of f is Jordan matrix. Furthermore, the number of Jordan blocks $J_r(\lambda)$ appearing in A , for fixed r and λ , is uniquely determined by f .*

Proof. The implication (1) \implies (3) follows from Theorem 25.11 and Proposition 25.12. The implications (3) \implies (2) and (2) \implies (1) are trivial. \square

Compared to Theorem 22.16, the new ingredient is the uniqueness assertion in (3), which is not so easy to prove.

Observe that the minimal polynomial of f is the least common multiple of the polynomials $(X - \lambda)^r$ associated with the Jordan blocks $J_r(\lambda)$ appearing in A , and the characteristic polynomial of A is the product of these polynomials.

We now return to the problem of computing effectively the similarity invariants of a matrix A . By Proposition 25.7, this is equivalent to computing the invariant factors of $XI - A$. In principle, this can be done using Proposition 24.42. A procedure to do this effectively for the ring $A = K[X]$ is to convert $XI - A$ to its Smith normal form. This will also yield the rational canonical form for A .

25.4 The Smith Normal Form

The Smith normal form is the special case of Proposition 24.42 applied to the PID $K[X]$ where K is a field, but it also says that the matrices P and Q are products of elementary matrices. It turns out that such a result holds for any Euclidean ring, and the proof is basically the same.

Recall from Definition 20.9 that a *Euclidean ring* is an integral domain A such that there exists a function $\sigma: A \rightarrow \mathbb{N}$ with the following property: For all $a, b \in A$ with $b \neq 0$, there are some $q, r \in A$ such that

$$a = bq + r \quad \text{and} \quad \sigma(r) < \sigma(b).$$

Note that the pair (q, r) is not necessarily unique.

We make use of the elementary row and column operations $P(i, k)$, $E_{i,j;\beta}$, and $E_{i,\lambda}$ described in Chapter 6, where we require the scalar λ used in $E_{i,\lambda}$ to be a unit.

Theorem 25.14. *If M is an $m \times n$ matrix over a Euclidean ring A , then there exist some invertible $n \times n$ matrix P and some invertible $m \times m$ matrix Q , where P and Q are products of elementary matrices, and a $m \times n$ matrix D of the form*

$$D = \begin{pmatrix} \alpha_1 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & \alpha_2 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & \alpha_r & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & 0 & \cdots & 0 \end{pmatrix}$$

for some nonzero $\alpha_i \in A$, such that

- (1) $\alpha_1 \mid \alpha_2 \mid \cdots \mid \alpha_r$, and
- (2) $M = QDP^{-1}$.

Proof. We follow Jacobson's proof [59] (Chapter 3, Theorem 3.8). We proceed by induction on $m + n$.

If $m = n = 1$, let $P = (1)$ and $Q = (1)$.

For the induction step, if $M = 0$, let $P = I_n$ and $Q = I_m$. If $M \neq 0$, the strategy is to apply a sequence of elementary transformations that converts M to a matrix of the form

$$M' = \begin{pmatrix} \alpha_1 & 0 & \cdots & 0 \\ 0 & & & \\ \vdots & & Y & \\ 0 & & & \end{pmatrix}$$

where Y is a $(m-1) \times (n-1)$ -matrix such that α_1 divides every entry in Y . Then, we proceed by induction on Y . To find M' , we perform the following steps.

Step 1. Pick some nonzero entry a_{ij} in M such that $\sigma(a_{ij})$ is minimal. Then permute column j and column 1, and permute row i and row 1, to bring this entry in position $(1, 1)$. We denote this new matrix again by M .

Step 2a.

If $m = 1$ go to Step 2b.

If $m > 1$, then there are two possibilities:

- (i) M is of the form

$$\begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ 0 & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{m2} & \cdots & a_{mn} \end{pmatrix}.$$

If $n = 1$, stop; else go to Step 2b.

(ii) There is some nonzero entry a_{i1} ($i > 1$) below a_{11} in the first column.

(a) If there is some entry a_{k1} in the first column such that a_{11} does not divide a_{k1} , then pick such an entry (say, with the smallest index i such that $\sigma(a_{i1})$ is minimal), and divide a_{k1} by a_{11} ; that is, find b_k and b_{k1} such that

$$a_{k1} = a_{11}b_k + b_{k1}, \quad \text{with } \sigma(b_{k1}) < \sigma(a_{11}).$$

Subtract b_k times row 1 from row k and permute row k and row 1, to obtain a matrix of the form

$$M = \begin{pmatrix} b_{k1} & b_{k2} & \cdots & b_{kn} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix}.$$

Go back to Step 2a.

(b) If a_{11} divides every (nonzero) entry a_{i1} for $i \geq 2$, say $a_{i1} = a_{11}b_i$, then subtract b_i times row 1 from row i for $i = 2, \dots, m$; go to Step 2b.

Observe that whenever we return to the beginning of Step 2a, we have $\sigma(b_{k1}) < \sigma(a_{11})$. Therefore, after a finite number of steps, we must exit Step 2a with a matrix in which all entries in column 1 but the first are zero and go to Step 2b.

Step 2b.

This step is reached only if $n > 1$ and if the only nonzero entry in the first column is a_{11} .

(a) If M is of the form

$$\begin{pmatrix} a_{11} & 0 & \cdots & 0 \\ 0 & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{m2} & \cdots & a_{mn} \end{pmatrix}$$

and $m = 1$ stop; else go to Step 3.

(b) If there is some entry a_{1k} in the first row such that a_{11} does not divide a_{1k} , then pick such an entry (say, with the smallest index j such that $\sigma(a_{1j})$ is minimal), and divide a_{1k} by a_{11} ; that is, find b_k and b_{1k} such that

$$a_{1k} = a_{11}b_k + b_{1k}, \quad \text{with } \sigma(b_{1k}) < \sigma(a_{11}).$$

Subtract b_k times column 1 from column k and permute column k and column 1, to obtain a matrix of the form

$$M = \begin{pmatrix} b_{1k} & a_{k2} & \cdots & a_{kn} \\ b_{2k} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ b_{mk} & a_{m2} & \cdots & a_{mn} \end{pmatrix}.$$

Go back to Step 2b.

(c) If a_{11} divides every (nonzero) entry a_{1j} for $j \geq 2$, say $a_{1j} = a_{11}b_j$, then subtract b_j times column 1 from column j for $j = 2, \dots, n$; go to Step 3.

As in Step 2a, whenever we return to the beginning of Step 2b, we have $\sigma(b_{1k}) < \sigma(a_{11})$. Therefore, after a finite number of steps, we must exit Step 2b with a matrix in which all entries in row 1 but the first are zero.

Step 3. This step is reached only if the only nonzero entry in the first row is a_{11} .

(i) If

$$M = \begin{pmatrix} a_{11} & 0 & \cdots & 0 \\ 0 & & & \\ \vdots & & Y & \\ 0 & & & \end{pmatrix}$$

go to Step 4.

(ii) If Step 2b ruined column 1 which now contains some nonzero entry below a_{11} , go back to Step 2a.

We perform a sequence of alternating steps between Step 2a and Step 2b. Because the σ -value of the $(1, 1)$ -entry strictly decreases whenever we reenter Step 2a and Step 2b, such a sequence must terminate with a matrix of the form

$$M = \begin{pmatrix} a_{11} & 0 & \cdots & 0 \\ 0 & & & \\ \vdots & & Y & \\ 0 & & & \end{pmatrix}$$

Step 4. If a_{11} divides all entries in Y , stop.

Otherwise, there is some column, say j , such that a_{11} does not divide some entry a_{ij} , so add the j th column to the first column. This yields a matrix of the form

$$M = \begin{pmatrix} a_{11} & 0 & \cdots & 0 \\ b_{2j} & & & \\ \vdots & & Y & \\ b_{mj} & & & \end{pmatrix}$$

where the i th entry in column 1 is nonzero, so go back to Step 2a,

Again, since the σ -value of the $(1, 1)$ -entry strictly decreases whenever we reenter Step 2a and Step 2b, such a sequence must terminate with a matrix of the form

$$M' = \begin{pmatrix} \alpha_1 & 0 & \cdots & 0 \\ 0 & & & \\ \vdots & & Y & \\ 0 & & & \end{pmatrix}$$

where α_1 divides every entry in Y . Then, we apply the induction hypothesis to Y . \square

If the PID A is the polynomial ring $K[X]$ where K is a field, the α_i are nonzero polynomials, so we can apply row operations to normalize their leading coefficients to be 1. We obtain the following theorem.

Theorem 25.15. (*Smith Normal Form*) *If M is an $m \times n$ matrix over the polynomial ring $K[X]$, where K is a field, then there exist some invertible $n \times n$ matrix P and some invertible $m \times m$ matrix Q , where P and Q are products of elementary matrices with entries in $K[X]$, and a $m \times n$ matrix D of the form*

$$D = \begin{pmatrix} q_1 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & q_2 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & q_r & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & 0 & \cdots & 0 \end{pmatrix}$$

for some nonzero monic polynomials $q_i \in k[X]$, such that

- (1) $q_1 \mid q_2 \mid \cdots \mid q_r$, and
- (2) $M = QDP^{-1}$.

In particular, if we apply Theorem 25.15 to a matrix M of the form $M = XI - A$, where A is a square matrix, then $\det(XI - A) = \chi_A(X)$ is never zero, and since $XI - A = QDP^{-1}$ with P, Q invertible, all the entries in D must be nonzero and we obtain the following result showing that the similarity invariants of A can be computed using elementary operations.

Theorem 25.16. *If A is an $n \times n$ matrix over the field K , then there exist some invertible $n \times n$ matrices P and Q , where P and Q are products of elementary matrices with entries in $K[X]$, and a $n \times n$ matrix D of the form*

$$D = \begin{pmatrix} 1 & \cdots & 0 & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 1 & 0 & 0 & \cdots & 0 \\ 0 & \cdots & 0 & q_1 & 0 & \cdots & 0 \\ 0 & \cdots & 0 & 0 & q_2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & 0 & 0 & \cdots & q_m \end{pmatrix}$$

for some nonzero monic polynomials $q_i \in k[X]$ of degree ≥ 1 , such that

- (1) $q_1 \mid q_2 \mid \cdots \mid q_m$,
- (2) q_1, \dots, q_m are the similarity invariants of A , and
- (3) $XI - A = QDP^{-1}$.

The matrix D in Theorem 25.16 is often called *Smith normal form* of A , even though this is confusing terminology since D is really the Smith normal form of $XI - A$.

Of course, we know from previous work that in Theorem 25.15, the $\alpha_1, \dots, \alpha_r$ are unique, and that in Theorem 25.16, the q_1, \dots, q_m are unique. This can also be proved using some simple properties of minors, but we leave it as an exercise (for help, look at Jacobson [59], Chapter 3, Theorem 3.9).

The rational canonical form of A can also be obtained from Q^{-1} and D , but first, let us consider the generalization of Theorem 25.15 to PID's that are not necessarily Euclidean rings.

We need to find a “norm” that assigns a natural number $\sigma(a)$ to any nonzero element of a PID A , in such a way that $\sigma(a)$ decreases whenever we return to Step 2a and Step 2b. Since a PID is a UFD, we use the number

$$\sigma(a) = k_1 + \cdots + k_r$$

of prime factors in the factorization of a nonunit element

$$a = up_1^{k_1} \cdots p_r^{k_r},$$

and we set

$$\sigma(u) = 0$$

if u is a unit.

We can't divide anymore, but we can find gcd's and use Bezout to mimic division. The key ingredient is this: for any two nonzero elements $a, b \in A$, if a does not divide b then let $d \neq 0$ be a gcd of a and b . By Bezout, there exist $x, y \in A$ such that

$$ax + by = d.$$

We can also write $a = td$ and $b = -sd$, for some $s, t \in A$, so that $tdx - sdy = d$, which implies that

$$tx - sy = 1,$$

since A is an integral domain. Observe that

$$\begin{pmatrix} t & -s \\ -y & x \end{pmatrix} \begin{pmatrix} x & s \\ y & t \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix},$$

which shows that both matrices on the left of the equation are invertible, and so is the transpose of the second one,

$$\begin{pmatrix} x & y \\ s & t \end{pmatrix}$$

(they all have determinant 1). We also have

$$as + bt = tds - sdt = 0,$$

so

$$\begin{pmatrix} x & y \\ s & t \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} d \\ 0 \end{pmatrix}$$

and

$$\begin{pmatrix} a & b \end{pmatrix} \begin{pmatrix} x & s \\ y & t \end{pmatrix} = \begin{pmatrix} d & 0 \end{pmatrix}.$$

Because a does not divide b , their gcd d has strictly fewer prime factors than a , so

$$\sigma(d) < \sigma(a).$$

Using matrices of the form

$$\begin{pmatrix} x & y & 0 & 0 & \cdots & 0 \\ s & t & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & \cdots & 1 \end{pmatrix}$$

with $xt - ys = 1$, we can modify Steps 2a and Step 2b to obtain the following theorem.

Theorem 25.17. *If M is an $m \times n$ matrix over a PID A , then there exist some invertible $n \times n$ matrix P and some invertible $m \times m$ matrix Q , where P and Q are products of elementary matrices and matrices of the form*

$$\begin{pmatrix} x & y & 0 & 0 & \cdots & 0 \\ s & t & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & \cdots & 1 \end{pmatrix}$$

with $xt - ys = 1$, and a $m \times n$ matrix D of the form

$$D = \begin{pmatrix} \alpha_1 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & \alpha_2 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & \alpha_r & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & 0 & \cdots & 0 \end{pmatrix}$$

for some nonzero $\alpha_i \in A$, such that

(1) $\alpha_1 \mid \alpha_2 \mid \cdots \mid \alpha_r$, and

(2) $M = QDP^{-1}$.

Proof sketch. In Step 2a, if a_{11} does not divide a_{k1} , then first permute row 2 and row k (if $k \neq 2$). Then, if we write $a = a_{11}$ and $b = a_{k1}$, if d is a gcd of a and b and if x, y, s, t are determined as explained above, multiply on the left by the matrix

$$\begin{pmatrix} x & y & 0 & 0 & \cdots & 0 \\ s & t & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & \cdots & 1 \end{pmatrix}$$

to obtain a matrix of the form

$$\begin{pmatrix} d & a_{12} & \cdots & a_{1n} \\ 0 & a_{22} & \cdots & a_{2n} \\ a_{31} & a_{32} & \cdots & a_{3n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix}$$

with $\sigma(d) < \sigma(a_{11})$. Then, go back to Step 2a.

In Step 2b, if a_{11} does not divide a_{1k} , then first permute column 2 and column k (if $k \neq 2$). Then, if we write $a = a_{11}$ and $b = a_{1k}$, if d is a gcd of a and b and if x, y, s, t are determined as explained above, multiply on the right by the matrix

$$\begin{pmatrix} x & s & 0 & 0 & \cdots & 0 \\ y & t & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & \cdots & 1 \end{pmatrix}$$

to obtain a matrix of the form

$$\begin{pmatrix} d & 0 & a_{13} & \cdots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \cdots & a_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & a_{m3} & \cdots & a_{mn} \end{pmatrix}$$

with $\sigma(d) < \sigma(a_{11})$. Then, go back to Step 2b. The other steps remain the same. Whenever we return to Step 2a or Step 2b, the σ -value of the $(1, 1)$ -entry strictly decreases, so the whole procedure terminates. \square

We conclude this section by explaining how the rational canonical form of a matrix A can be obtained from the canonical form QDP^{-1} of $XI - A$.

Let $f: E \rightarrow E$ be a linear map over a K -vector space of dimension n . Recall from Theorem 24.21 (see Section 24.5) that as a $K[X]$ -module, E_f is the image of the free module $E[X]$ by the map $\sigma: E[X] \rightarrow E_f$, where $E[X]$ consists of all linear combinations of the form

$$p_1 e_1 + \cdots + p_n e_n,$$

where (e_1, \dots, e_n) is a basis of E and $p_1, \dots, p_n \in K[X]$ are polynomials, and σ is given by

$$\sigma(p_1 e_1 + \cdots + p_n e_n) = p_1(f)(e_1) + \cdots + p_n(f)(e_n).$$

Furthermore, the kernel of σ is equal to the image of the map $\psi: E[X] \rightarrow E[X]$, where

$$\psi(p_1 e_1 + \cdots + p_n e_n) = Xp_1 e_1 + \cdots + Xp_n e_n - (p_1 f(e_1) + \cdots + p_n f(e_n)).$$

The matrix A is the representation of a linear map f over the canonical basis (e_1, \dots, e_n) of $E = K^n$, and $XI - A$ is the matrix of ψ with respect to the basis (e_1, \dots, e_n) (over $K[X]$). What Theorem 25.16 tells us is that there are $K[X]$ -bases (u_1, \dots, u_n) and (v_1, \dots, v_n) of E_f with respect to which the matrix of ψ is D . Then

$$\begin{aligned} \psi(u_i) &= v_i, & i &= 1, \dots, n-m, \\ \psi(u_{n-m+i}) &= q_i v_{n-m+i}, & i &= 1, \dots, m, \end{aligned}$$

and because $\text{Im}(\psi) = \text{Ker}(\sigma)$, this implies that

$$\sigma(v_i) = 0, \quad i = 1, \dots, n-m.$$

Consequently, $w_1 = \sigma(v_{n-m+1}), \dots, w_m = \sigma(v_n)$ span E_f as a $K[X]$ -module, with $w_i \in E$, and we have

$$M_f = K[X]w_1 \oplus \cdots \oplus K[X]w_m,$$

where $K[X]w_i \approx K[X]/(q_i)$ as a cyclic $K[X]$ -module. Since $\text{Im}(\psi) = \text{Ker}(\sigma)$, we have

$$0 = \sigma(\psi(u_{n-m+i})) = \sigma(q_i v_{n-m+i}) = q_i \sigma(v_{n-m+i}) = q_i w_i,$$

so as a K -vector space, the cyclic subspace $Z(w_i; f) = K[X]w_i$ has q_i as annihilator, and by a remark from Section 24.5, it has the basis (over K)

$$(w_i, f(w_i), \dots, f^{n_i-1}(w_i)), \quad n_i = \deg(q_i).$$

Furthermore, over this basis, the restriction of f to $Z(w_i; f)$ is represented by the companion matrix of q_i . By putting all these bases together, we obtain a block matrix which is the canonical rational form of f (and A).

Now, $XI - A = QDP^{-1}$ is the matrix of ψ with respect to the canonical basis (e_1, \dots, e_n) (over $K[X]$), and D is the matrix of ψ with respect to the bases (u_1, \dots, u_n) and (v_1, \dots, v_n) (over $K[X]$), which tells us that the columns of Q consist of the coordinates (in $K[X]$) of the basis vectors (v_1, \dots, v_n) with respect to the basis (e_1, \dots, e_n) . Therefore, the coordinates (in K) of the vectors (w_1, \dots, w_m) spanning E_f over $K[X]$, where $w_i = \sigma(v_{n-m+i})$, are obtained by substituting the matrix A for X in the coordinates of the columns vectors of Q , and evaluating the resulting expressions.

Since

$$D = Q^{-1}(XI - A)P,$$

the matrix D is obtained from A by a sequence of elementary row operations whose product is Q^{-1} and a sequence of elementary column operations whose product is P . Therefore, to compute the vectors w_1, \dots, w_m from A , we simply have to figure out how to construct Q from the sequence of elementary row operations that yield Q^{-1} . The trick is to use column operations to gather a product of row operations in reverse order.

Indeed, if Q^{-1} is the product of elementary row operations

$$Q^{-1} = E_k \cdots E_2 E_1,$$

then

$$Q = E_1^{-1} E_2^{-1} \cdots E_k^{-1}.$$

Now, row operations operate on the left and column operations operate on the right, so the product $E_1^{-1} E_2^{-1} \cdots E_k^{-1}$ can be computed from left to right as a sequence of column operations.

Let us review the meaning of the elementary row and column operations $P(i, k)$, $E_{i,j;\beta}$, and $E_{i,\lambda}$.

1. As a row operation, $P(i, k)$ permutes row i and row k .
2. As a column operation, $P(i, k)$ permutes column i and column k .
3. The inverse of $P(i, k)$ is $P(i, k)$ itself.
4. As a row operation, $E_{i,j;\beta}$ adds β times row j to row i .

5. As a column operation, $E_{i,j;\beta}$ adds β times column i to column j (note the switch in the indices).
6. The inverse of $E_{i,j;\beta}$ is $E_{i,j;-\beta}$.
7. As a row operation, $E_{i,\lambda}$ multiplies row i by λ .
8. As a column operation, $E_{i,\lambda}$ multiplies column i by λ .
9. The inverse of $E_{i,\lambda}$ is $E_{i,\lambda^{-1}}$.

Given a square matrix A (over K), the row and column operations applied to $XI - A$ in converting it to its Smith normal form may involve coefficients that are polynomials and it is necessary to explain what is the action of an operation $E_{i,j;\beta}$ in this case. If the coefficient β in $E_{i,j;\beta}$ is a polynomial over K , as a row operation, the action of $E_{i,j;\beta}$ on a matrix X is to multiply the j th row of M by the matrix $\beta(A)$ obtained by substituting the matrix A for X and then to add the resulting vector to row i . Similarly, as a column operation, the action of $E_{i,j;\beta}$ on a matrix X is to multiply the i th column of M by the matrix $\beta(A)$ obtained by substituting the matrix A for X and then to add the resulting vector to column j . An algorithm to compute the rational canonical form of a matrix can now be given. We apply the elementary column operations E_i^{-1} for $i = 1, \dots, k$, starting with the identity matrix.

Algorithm for Converting an $n \times n$ matrix to Rational Canonical Form

While applying elementary row and column operations to compute the Smith normal form D of $XI - A$, keep track of the row operations and perform the following steps:

1. Let $P' = I_n$, and for every elementary row operation E do the following:
 - (a) If $E = P(i, k)$, permute column i and column k of P' .
 - (b) If $E = E_{i,j;\beta}$, multiply the i th column of P' by the matrix $\beta(A)$ obtained by substituting the matrix A for X , and then subtract the resulting vector from column j .
 - (c) If $E = E_{i,\lambda}$ where $\lambda \in K$, then multiply the i th column of P' by λ^{-1} .
2. When step (1) terminates, the first $n - m$ columns of P' are zero and the last m are linearly independent. For $i = 1, \dots, m$, multiply the $(n - m + i)$ th column w_i of P' successively by I, A^1, A^2, A^{n_i-1} , where n_i is the degree of the polynomial q_i (appearing in D), and form the $n \times n$ matrix P consisting of the vectors

$$w_1, Aw_1, \dots, A^{n_1-1}w_1, w_2, Aw_2, \dots, A^{n_2-1}w_2, \dots, w_m, Aw_m, \dots, A^{n_m-1}w_m.$$

Then, $P^{-1}AP$ is the canonical rational form of A .

Here is an example taken from Dummit and Foote [32] (Chapter 12, Section 12.2). Let A be the matrix

$$A = \begin{pmatrix} 1 & 2 & -4 & 4 \\ 2 & -1 & 4 & -8 \\ 1 & 0 & 1 & -2 \\ 0 & 1 & -2 & 3 \end{pmatrix}.$$

One should check that the following sequence of row and column operations produces the Smith normal form D of $XI - A$:

$$\begin{array}{llllll} \text{row } P(1, 3) & \text{row } E_{1,-1} & \text{row } E_{2,1;2} & \text{row } E_{3,1;-(X-1)} & \text{column } E_{1,3;X-1} & \text{column } E_{1,4;2} \\ \text{row } P(2, 4) & \text{row } E_{2,-1} & \text{row } E_{3,2;2} & \text{row } E_{4,2;-(X+1)} & \text{column } E_{2,3;2} & \text{column } E_{2,4;X-3}, \end{array}$$

with

$$D = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & (X-1)^2 & 0 \\ 0 & 0 & 0 & (X-1)^2 \end{pmatrix}.$$

Then, applying Step 1 of the above algorithm, we get the sequence of column operations:

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow{P(1,3)} \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow{E_{1,-1}} \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow{E_{2,1,-2}} \begin{pmatrix} 0 & 0 & 1 & 0 \\ -2 & 1 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow{E_{3,1,A-I}} \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \xrightarrow{P(2,4)} \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \xrightarrow{E_{2,-1}} \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \xrightarrow{E_{3,2,-2}} \begin{pmatrix} 0 & -2 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \end{pmatrix} \xrightarrow{E_{4,2;A+I}} \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} = P'.$$

Step 2 of the algorithm yields the vectors

$$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \quad A \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \\ 1 \\ 0 \end{pmatrix}, \quad \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}, \quad A \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 2 \\ -1 \\ 0 \\ 1 \end{pmatrix},$$

so we get

$$P = \begin{pmatrix} 1 & 1 & 0 & 2 \\ 0 & 2 & 1 & -1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

We find that

$$P^{-1} = \begin{pmatrix} 1 & 0 & -1 & -2 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & -2 & 1 \\ 0 & 0 & 0 & 1 \end{pmatrix},$$

and thus, the rational canonical form of A is

$$P^{-1}AP = \begin{pmatrix} 0 & -1 & 0 & 0 \\ 1 & 2 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & 1 & 2 \end{pmatrix}.$$

Chapter 26

Topology

26.1 Metric Spaces and Normed Vector Spaces

This chapter contains a review of basic topological concepts. First, metric spaces are defined. Next, normed vector spaces are defined. Closed and open sets are defined, and their basic properties are stated. The general concept of a topological space is defined. The closure and the interior of a subset are defined. The subspace topology and the product topology are defined. Continuous maps and homeomorphisms are defined. Limits of sequences are defined. Continuous linear maps and multilinear maps are defined and studied briefly. The chapter ends with the definition of a normed affine space.

Most spaces considered in this book have a topological structure given by a metric or a norm, and we first review these notions. We begin with metric spaces. Recall that $\mathbb{R}_+ = \{x \in \mathbb{R} \mid x \geq 0\}$.

Definition 26.1. A *metric space* is a set E together with a function $d: E \times E \rightarrow \mathbb{R}_+$, called a *metric*, or *distance*, assigning a nonnegative real number $d(x, y)$ to any two points $x, y \in E$, and satisfying the following conditions for all $x, y, z \in E$:

$$(D1) \quad d(x, y) = d(y, x). \quad (\text{symmetry})$$

$$(D2) \quad d(x, y) \geq 0, \text{ and } d(x, y) = 0 \text{ iff } x = y. \quad (\text{positivity})$$

$$(D3) \quad d(x, z) \leq d(x, y) + d(y, z). \quad (\text{triangle inequality})$$

Geometrically, condition (D3) expresses the fact that in a triangle with vertices x, y, z , the length of any side is bounded by the sum of the lengths of the other two sides. From (D3), we immediately get

$$|d(x, y) - d(y, z)| \leq d(x, z).$$

Let us give some examples of metric spaces. Recall that the *absolute value* $|x|$ of a real number $x \in \mathbb{R}$ is defined such that $|x| = x$ if $x \geq 0$, $|x| = -x$ if $x < 0$, and for a complex number $x = a + ib$, by $|x| = \sqrt{a^2 + b^2}$.

Example 26.1.

1. Let $E = \mathbb{R}$, and $d(x, y) = |x - y|$, the absolute value of $x - y$. This is the so-called natural metric on \mathbb{R} .
2. Let $E = \mathbb{R}^n$ (or $E = \mathbb{C}^n$). We have the *Euclidean metric*

$$d_2(x, y) = (|x_1 - y_1|^2 + \cdots + |x_n - y_n|^2)^{\frac{1}{2}},$$

the distance between the points (x_1, \dots, x_n) and (y_1, \dots, y_n) .

3. For every set E , we can define the *discrete metric*, defined such that $d(x, y) = 1$ iff $x \neq y$, and $d(x, x) = 0$.
4. For any $a, b \in \mathbb{R}$ such that $a < b$, we define the following sets:

$$[a, b] = \{x \in \mathbb{R} \mid a \leq x \leq b\}, \quad (\text{closed interval})$$

$$]a, b[= \{x \in \mathbb{R} \mid a < x < b\}, \quad (\text{open interval})$$

$$[a, b[= \{x \in \mathbb{R} \mid a \leq x < b\}, \quad (\text{interval closed on the left, open on the right})$$

$$]a, b] = \{x \in \mathbb{R} \mid a < x \leq b\}, \quad (\text{interval open on the left, closed on the right})$$

Let $E = [a, b]$, and $d(x, y) = |x - y|$. Then, $([a, b], d)$ is a metric space.

We will need to define the notion of proximity in order to define convergence of limits and continuity of functions. For this, we introduce some standard “small neighborhoods.”

Definition 26.2. Given a metric space E with metric d , for every $a \in E$, for every $\rho \in \mathbb{R}$, with $\rho > 0$, the set

$$B(a, \rho) = \{x \in E \mid d(a, x) \leq \rho\}$$

is called the *closed ball of center a and radius ρ* , the set

$$B_0(a, \rho) = \{x \in E \mid d(a, x) < \rho\}$$

is called the *open ball of center a and radius ρ* , and the set

$$S(a, \rho) = \{x \in E \mid d(a, x) = \rho\}$$

is called the *sphere of center a and radius ρ* . It should be noted that ρ is finite (i.e., not $+\infty$). A subset X of a metric space E is *bounded* if there is a closed ball $B(a, \rho)$ such that $X \subseteq B(a, \rho)$.

Clearly, $B(a, \rho) = B_0(a, \rho) \cup S(a, \rho)$.

Example 26.2.

1. In $E = \mathbb{R}$ with the distance $|x - y|$, an open ball of center a and radius ρ is the open interval $]a - \rho, a + \rho[$.
2. In $E = \mathbb{R}^2$ with the Euclidean metric, an open ball of center a and radius ρ is the set of points inside the disk of center a and radius ρ , excluding the boundary points on the circle.
3. In $E = \mathbb{R}^3$ with the Euclidean metric, an open ball of center a and radius ρ is the set of points inside the sphere of center a and radius ρ , excluding the boundary points on the sphere.

One should be aware that intuition can be misleading in forming a geometric image of a closed (or open) ball. For example, if d is the discrete metric, a closed ball of center a and radius $\rho < 1$ consists only of its center a , and a closed ball of center a and radius $\rho \geq 1$ consists of the entire space!



If $E = [a, b]$, and $d(x, y) = |x - y|$, as in Example 26.1, an open ball $B_0(a, \rho)$, with $\rho < b - a$, is in fact the interval $[a, a + \rho[$, which is closed on the left.

We now consider a very important special case of metric spaces, normed vector spaces. Normed vector spaces have already been defined in Chapter 7 (Definition 7.1) but for the reader's convenience we repeat the definition.

Definition 26.3. Let E be a vector space over a field K , where K is either the field \mathbb{R} of reals, or the field \mathbb{C} of complex numbers. A *norm on E* is a function $\|\cdot\|: E \rightarrow \mathbb{R}_+$, assigning a nonnegative real number $\|u\|$ to any vector $u \in E$, and satisfying the following conditions for all $x, y, z \in E$:

$$(N1) \quad \|x\| \geq 0, \text{ and } \|x\| = 0 \text{ iff } x = 0. \quad (\text{positivity})$$

$$(N2) \quad \|\lambda x\| = |\lambda| \|x\|. \quad (\text{scaling})$$

$$(N3) \quad \|x + y\| \leq \|x\| + \|y\|. \quad (\text{triangle inequality})$$

A vector space E together with a norm $\|\cdot\|$ is called a *normed vector space*.

From (N3), we easily get

$$||\|x\| - \|y\|| \leq \|x - y\|.$$

Given a normed vector space E , if we define d such that

$$d(x, y) = \|x - y\|,$$

it is easily seen that d is a metric. Thus, every normed vector space is immediately a metric space. Note that the metric associated with a norm is invariant under translation, that is,

$$d(x + u, y + u) = d(x, y).$$

For this reason, we can restrict ourselves to open or closed balls of center 0.

Examples of normed vector spaces were given in Example 7.1. We repeat the most important examples.

Example 26.3. Let $E = \mathbb{R}^n$ (or $E = \mathbb{C}^n$). There are three standard norms. For every $(x_1, \dots, x_n) \in E$, we have the norm $\|x\|_1$, defined such that,

$$\|x\|_1 = |x_1| + \dots + |x_n|,$$

we have the *Euclidean norm* $\|x\|_2$, defined such that,

$$\|x\|_2 = (|x_1|^2 + \dots + |x_n|^2)^{\frac{1}{2}},$$

and the *sup-norm* $\|x\|_\infty$, defined such that,

$$\|x\|_\infty = \max\{|x_i| \mid 1 \leq i \leq n\}.$$

More generally, we define the ℓ_p -norm (for $p \geq 1$) by

$$\|x\|_p = (|x_1|^p + \dots + |x_n|^p)^{1/p}.$$

We proved in Proposition 7.1 that the ℓ_p -norms are indeed norms. One should work out what are the open balls in \mathbb{R}^2 for $\|\cdot\|_1$ and $\|\cdot\|_\infty$.

In a normed vector space, we define a closed ball or an open ball of radius ρ as a closed ball or an open ball of center 0. We may use the notation $B(\rho)$ and $B_0(\rho)$.

We will now define the crucial notions of open sets and closed sets, and of a topological space.

Definition 26.4. Let E be a metric space with metric d . A subset $U \subseteq E$ is an *open set* in E if either $U = \emptyset$, or for every $a \in U$, there is some open ball $B_0(a, \rho)$ such that, $B_0(a, \rho) \subseteq U$.¹ A subset $F \subseteq E$ is a *closed set* in E if its complement $E - F$ is open in E .

The set E itself is open, since for every $a \in E$, every open ball of center a is contained in E . In $E = \mathbb{R}^n$, given n intervals $[a_i, b_i]$, with $a_i < b_i$, it is easy to show that the open n -cube

$$\{(x_1, \dots, x_n) \in E \mid a_i < x_i < b_i, 1 \leq i \leq n\}$$

¹Recall that $\rho > 0$.

is an open set. In fact, it is possible to find a metric for which such open n -cubes are open balls! Similarly, we can define the closed n -cube

$$\{(x_1, \dots, x_n) \in E \mid a_i \leq x_i \leq b_i, 1 \leq i \leq n\},$$

which is a closed set.

The open sets satisfy some important properties that lead to the definition of a topological space.

Proposition 26.1. *Given a metric space E with metric d , the family \mathcal{O} of all open sets defined in Definition 26.4 satisfies the following properties:*

- (O1) *For every finite family $(U_i)_{1 \leq i \leq n}$ of sets $U_i \in \mathcal{O}$, we have $U_1 \cap \dots \cap U_n \in \mathcal{O}$, i.e., \mathcal{O} is closed under finite intersections.*
- (O2) *For every arbitrary family $(U_i)_{i \in I}$ of sets $U_i \in \mathcal{O}$, we have $\bigcup_{i \in I} U_i \in \mathcal{O}$, i.e., \mathcal{O} is closed under arbitrary unions.*
- (O3) *$\emptyset \in \mathcal{O}$, and $E \in \mathcal{O}$, i.e., \emptyset and E belong to \mathcal{O} .*

Furthermore, for any two distinct points $a \neq b$ in E , there exist two open sets U_a and U_b such that, $a \in U_a$, $b \in U_b$, and $U_a \cap U_b = \emptyset$.

Proof. It is straightforward. For the last point, letting $\rho = d(a, b)/3$ (in fact $\rho = d(a, b)/2$ works too), we can pick $U_a = B_0(a, \rho)$ and $U_b = B_0(b, \rho)$. By the triangle inequality, we must have $U_a \cap U_b = \emptyset$. \square

The above proposition leads to the very general concept of a topological space.



One should be careful that, in general, the family of open sets is not closed under infinite intersections. For example, in \mathbb{R} under the metric $|x - y|$, letting $U_n =] - 1/n, +1/n[$, each U_n is open, but $\bigcap_n U_n = \{0\}$, which is not open.

26.2 Topological Spaces

Motivated by Proposition 26.1, a topological space is defined in terms of a family of sets satisfying the properties of open sets stated in that proposition.

Definition 26.5. Given a set E , a *topology on E* (or a *topological structure on E*), is defined as a family \mathcal{O} of subsets of E called *open sets*, and satisfying the following three properties:

- (1) For every finite family $(U_i)_{1 \leq i \leq n}$ of sets $U_i \in \mathcal{O}$, we have $U_1 \cap \dots \cap U_n \in \mathcal{O}$, i.e., \mathcal{O} is closed under finite intersections.
- (2) For every arbitrary family $(U_i)_{i \in I}$ of sets $U_i \in \mathcal{O}$, we have $\bigcup_{i \in I} U_i \in \mathcal{O}$, i.e., \mathcal{O} is closed under arbitrary unions.

(3) $\emptyset \in \mathcal{O}$, and $E \in \mathcal{O}$, i.e., \emptyset and E belong to \mathcal{O} .

A set E together with a topology \mathcal{O} on E is called a *topological space*. Given a topological space (E, \mathcal{O}) , a subset F of E is a *closed set* if $F = E - U$ for some open set $U \in \mathcal{O}$, i.e., F is the complement of some open set.



It is possible that an open set is also a closed set. For example, \emptyset and E are both open and closed. When a topological space contains a proper nonempty subset U which is both open and closed, the space E is said to be *disconnected*.

A topological space (E, \mathcal{O}) is said to satisfy the *Hausdorff separation axiom* (or T_2 -separation axiom) if for any two distinct points $a \neq b$ in E , there exist two open sets U_a and U_b such that, $a \in U_a$, $b \in U_b$, and $U_a \cap U_b = \emptyset$. When the T_2 -separation axiom is satisfied, we also say that (E, \mathcal{O}) is a *Hausdorff space*.

As shown by Proposition 26.1, any metric space is a topological Hausdorff space, the family of open sets being in fact the family of arbitrary unions of open balls. Similarly, any normed vector space is a topological Hausdorff space, the family of open sets being the family of arbitrary unions of open balls. The topology \mathcal{O} consisting of all subsets of E is called the *discrete topology*.

Remark: Most (if not all) spaces used in analysis are Hausdorff spaces. Intuitively, the Hausdorff separation axiom says that there are enough “small” open sets. Without this axiom, some counter-intuitive behaviors may arise. For example, a sequence may have more than one limit point (or a compact set may not be closed). Nevertheless, non-Hausdorff topological spaces arise naturally in algebraic geometry. But even there, some substitute for separation is used.

One of the reasons why topological spaces are important is that the definition of a topology only involves a certain family \mathcal{O} of sets, and not **how** such family is generated from a metric or a norm. For example, different metrics or different norms can define the same family of open sets. Many topological properties only depend on the family \mathcal{O} and not on the specific metric or norm. But the fact that a topology is definable from a metric or a norm is important, because it usually implies nice properties of a space. All our examples will be spaces whose topology is defined by a metric or a norm.

By taking complements, we can state properties of the closed sets dual to those of Definition 26.5. Thus, \emptyset and E are closed sets, and the closed sets are closed under finite unions and arbitrary intersections.

It is also worth noting that the Hausdorff separation axiom implies that for every $a \in E$, the set $\{a\}$ is closed. Indeed, if $x \in E - \{a\}$, then $x \neq a$, and so there exist open sets U_a and U_x such that $a \in U_a$, $x \in U_x$, and $U_a \cap U_x = \emptyset$. Thus, for every $x \in E - \{a\}$, there is an open set U_x containing x and contained in $E - \{a\}$, showing by (O3) that $E - \{a\}$ is open, and thus that the set $\{a\}$ is closed.

Given a topological space (E, \mathcal{O}) , given any subset A of E , since $E \in \mathcal{O}$ and E is a closed set, the family $\mathcal{C}_A = \{F \mid A \subseteq F, F \text{ a closed set}\}$ of closed sets containing A is nonempty, and since any arbitrary intersection of closed sets is a closed set, the intersection $\bigcap \mathcal{C}_A$ of the sets in the family \mathcal{C}_A is the smallest closed set containing A . By a similar reasoning, the union of all the open subsets contained in A is the largest open set contained in A .

Definition 26.6. Given a topological space (E, \mathcal{O}) , given any subset A of E , the smallest closed set containing A is denoted by \overline{A} , and is called the *closure*, or *adherence* of A . A subset A of E is *dense in E* if $\overline{A} = E$. The largest open set contained in A is denoted by $\overset{\circ}{A}$, and is called the *interior* of A . The set $\text{Fr } A = \overline{A} \cap \overline{E - A}$ is called the *boundary* (or *frontier*) of A . We also denote the boundary of A by ∂A .

Remark: The notation \overline{A} for the closure of a subset A of E is somewhat unfortunate, since \overline{A} is often used to denote the set complement of A in E . Still, we prefer it to more cumbersome notations such as $\text{clo}(A)$, and we denote the complement of A in E by $E - A$ (or sometimes, A^c).

By definition, it is clear that a subset A of E is closed iff $A = \overline{A}$. The set \mathbb{Q} of rationals is dense in \mathbb{R} . It is easily shown that $\overline{A} = \overset{\circ}{A} \cup \partial A$ and $\overset{\circ}{A} \cap \partial A = \emptyset$. Another useful characterization of \overline{A} is given by the following proposition.

Proposition 26.2. *Given a topological space (E, \mathcal{O}) , given any subset A of E , the closure \overline{A} of A is the set of all points $x \in E$ such that for every open set U containing x , then $U \cap A \neq \emptyset$.*

Proof. If $A = \emptyset$, since \emptyset is closed, the proposition holds trivially. Thus, assume that $A \neq \emptyset$. First, assume that $x \in \overline{A}$. Let U be any open set such that $x \in U$. If $U \cap A = \emptyset$, since U is open, then $E - U$ is a closed set containing A , and since \overline{A} is the intersection of all closed sets containing A , we must have $x \in E - U$, which is impossible. Conversely, assume that $x \in E$ is a point such that for every open set U containing x , then $U \cap A \neq \emptyset$. Let F be any closed subset containing A . If $x \notin F$, since F is closed, then $U = E - F$ is an open set such that $x \in U$, and $U \cap A = \emptyset$, a contradiction. Thus, we have $x \in F$ for every closed set containing A , that is, $x \in \overline{A}$. \square

Often, it is necessary to consider a subset A of a topological space E , and to view the subset A as a topological space. The following proposition shows how to define a topology on a subset.

Proposition 26.3. *Given a topological space (E, \mathcal{O}) , given any subset A of E , let*

$$\mathcal{U} = \{U \cap A \mid U \in \mathcal{O}\}$$

be the family of all subsets of A obtained as the intersection of any open set in \mathcal{O} with A . The following properties hold.

- (1) The space (A, \mathcal{U}) is a topological space.
- (2) If E is a metric space with metric d , then the restriction $d_A: A \times A \rightarrow \mathbb{R}_+$ of the metric d to A defines a metric space. Furthermore, the topology induced by the metric d_A agrees with the topology defined by \mathcal{U} , as above.

Proof. Left as an exercise. □

Proposition 26.3 suggests the following definition.

Definition 26.7. Given a topological space (E, \mathcal{O}) , given any subset A of E , the *subspace topology on A induced by \mathcal{O}* is the family \mathcal{U} of open sets defined such that

$$\mathcal{U} = \{U \cap A \mid U \in \mathcal{O}\}$$

is the family of all subsets of A obtained as the intersection of any open set in \mathcal{O} with A . We say that (A, \mathcal{U}) has the *subspace topology*. If (E, d) is a metric space, the restriction $d_A: A \times A \rightarrow \mathbb{R}_+$ of the metric d to A is called the *subspace metric*.

For example, if $E = \mathbb{R}^n$ and d is the Euclidean metric, we obtain the subspace topology on the closed n -cube

$$\{(x_1, \dots, x_n) \in E \mid a_i \leq x_i \leq b_i, 1 \leq i \leq n\}.$$



One should realize that every open set $U \in \mathcal{O}$ which is entirely contained in A is also in the family \mathcal{U} , but \mathcal{U} may contain open sets that are not in \mathcal{O} . For example, if $E = \mathbb{R}$ with $|x - y|$, and $A = [a, b]$, then sets of the form $[a, c[$, with $a < c < b$ belong to \mathcal{U} , but they are not open sets for \mathbb{R} under $|x - y|$. However, there is agreement in the following situation.

Proposition 26.4. *Given a topological space (E, \mathcal{O}) , given any subset A of E , if \mathcal{U} is the subspace topology, then the following properties hold.*

- (1) If A is an open set $A \in \mathcal{O}$, then every open set $U \in \mathcal{U}$ is an open set $U \in \mathcal{O}$.
- (2) If A is a closed set in E , then every closed set w.r.t. the subspace topology is a closed set w.r.t. \mathcal{O} .

Proof. Left as an exercise. □

The concept of product topology is also useful. We have the following proposition.

Proposition 26.5. *Given n topological spaces (E_i, \mathcal{O}_i) , let \mathcal{B} be the family of subsets of $E_1 \times \dots \times E_n$ defined as follows:*

$$\mathcal{B} = \{U_1 \times \dots \times U_n \mid U_i \in \mathcal{O}_i, 1 \leq i \leq n\},$$

and let \mathcal{P} be the family consisting of arbitrary unions of sets in \mathcal{B} , including \emptyset . Then, \mathcal{P} is a topology on $E_1 \times \dots \times E_n$.

Proof. Left as an exercise. □

Definition 26.8. Given n topological spaces (E_i, \mathcal{O}_i) , the *product topology* on $E_1 \times \cdots \times E_n$ is the family \mathcal{P} of subsets of $E_1 \times \cdots \times E_n$ defined as follows: if

$$\mathcal{B} = \{U_1 \times \cdots \times U_n \mid U_i \in \mathcal{O}_i, 1 \leq i \leq n\},$$

then \mathcal{P} is the family consisting of arbitrary unions of sets in \mathcal{B} , including \emptyset .

If each $(E_i, \|\cdot\|_i)$ is a normed vector space, there are three natural norms that can be defined on $E_1 \times \cdots \times E_n$:

$$\begin{aligned}\|(x_1, \dots, x_n)\|_1 &= \|x_1\|_1 + \cdots + \|x_n\|_n, \\ \|(x_1, \dots, x_n)\|_2 &= \left(\|x_1\|_1^2 + \cdots + \|x_n\|_n^2\right)^{\frac{1}{2}}, \\ \|(x_1, \dots, x_n)\|_\infty &= \max\{\|x_1\|_1, \dots, \|x_n\|_n\}.\end{aligned}$$

It is easy to show that they all define the same topology, which is the product topology. It can also be verified that when $E_i = \mathbb{R}$, with the standard topology induced by $|x - y|$, the topology product on \mathbb{R}^n is the standard topology induced by the Euclidean norm.

Definition 26.9. Two metrics d_1 and d_2 on a space E are *equivalent* if they induce the same topology \mathcal{O} on E (i.e., they define the same family \mathcal{O} of open sets). Similarly, two norms $\|\cdot\|_1$ and $\|\cdot\|_2$ on a space E are *equivalent* if they induce the same topology \mathcal{O} on E .

Remark: Given a topological space (E, \mathcal{O}) , it is often useful, as in Proposition 26.5, to define the topology \mathcal{O} in terms of a subfamily \mathcal{B} of subsets of E . We say that a family \mathcal{B} of subsets of E is a *basis for the topology* \mathcal{O} , if \mathcal{B} is a subset of \mathcal{O} , and if every open set U in \mathcal{O} can be obtained as some union (possibly infinite) of sets in \mathcal{B} (agreeing that the empty union is the empty set).

It is immediately verified that if a family $\mathcal{B} = (U_i)_{i \in I}$ is a basis for the topology of (E, \mathcal{O}) , then $E = \bigcup_{i \in I} U_i$, and the intersection of any two sets $U_i, U_j \in \mathcal{B}$ is the union of some sets in the family \mathcal{B} (again, agreeing that the empty union is the empty set). Conversely, a family \mathcal{B} with these properties is the basis of the topology obtained by forming arbitrary unions of sets in \mathcal{B} .

A *subbasis* for \mathcal{O} is a family \mathcal{S} of subsets of E , such that the family \mathcal{B} of all finite intersections of sets in \mathcal{S} (including E itself, in case of the empty intersection) is a basis of \mathcal{O} .

The following proposition gives useful criteria for determining whether a family of open subsets is a basis of a topological space.

Proposition 26.6. *Given a topological space (E, \mathcal{O}) and a family \mathcal{B} of open subsets in \mathcal{O} the following properties hold:*

- (1) The family \mathcal{B} is a basis for the topology \mathcal{O} iff for every open set $U \in \mathcal{O}$ and every $x \in U$, there is some $B \in \mathcal{B}$ such that $x \in B$ and $B \subseteq U$.
- (2) The family \mathcal{B} is a basis for the topology \mathcal{O} iff
- (a) For every $x \in E$, there is some $B \in \mathcal{B}$ such that $x \in B$.
 - (b) For any two open subsets, $B_1, B_2 \in \mathcal{B}$, for every $x \in E$, if $x \in B_1 \cap B_2$, then there is some $B_3 \in \mathcal{B}$ such that $x \in B_3$ and $B_3 \subseteq B_1 \cap B_2$.

We now consider the fundamental property of continuity.

26.3 Continuous Functions, Limits

Definition 26.10. Let (E, \mathcal{O}_E) and (F, \mathcal{O}_F) be topological spaces, and let $f: E \rightarrow F$ be a function. For every $a \in E$, we say that f is *continuous at a* , if for every open set $V \in \mathcal{O}_F$ containing $f(a)$, there is some open set $U \in \mathcal{O}_E$ containing a , such that, $f(U) \subseteq V$. We say that f is *continuous* if it is continuous at every $a \in E$.

Define a *neighborhood* of $a \in E$ as any subset N of E containing some open set $O \in \mathcal{O}$ such that $a \in O$. Now, if f is continuous at a and N is any neighborhood of $f(a)$, there is some open set $V \subseteq N$ containing $f(a)$, and since f is continuous at a , there is some open set U containing a , such that $f(U) \subseteq V$. Since $V \subseteq N$, the open set U is a subset of $f^{-1}(N)$ containing a , and $f^{-1}(N)$ is a neighborhood of a . Conversely, if $f^{-1}(N)$ is a neighborhood of a whenever N is any neighborhood of $f(a)$, it is immediate that f is continuous at a . It is easy to see that Definition 26.10 is equivalent to the following statements.

Proposition 26.7. Let (E, \mathcal{O}_E) and (F, \mathcal{O}_F) be topological spaces, and let $f: E \rightarrow F$ be a function. For every $a \in E$, the function f is continuous at $a \in E$ iff for every neighborhood N of $f(a) \in F$, then $f^{-1}(N)$ is a neighborhood of a . The function f is continuous on E iff $f^{-1}(V)$ is an open set in \mathcal{O}_E for every open set $V \in \mathcal{O}_F$.

If E and F are metric spaces defined by metrics d_1 and d_2 , we can show easily that f is continuous at a iff

for every $\epsilon > 0$, there is some $\eta > 0$, such that, for every $x \in E$,

$$\text{if } d_1(a, x) \leq \eta, \text{ then } d_2(f(a), f(x)) \leq \epsilon.$$

Similarly, if E and F are normed vector spaces defined by norms $\| \cdot \|_1$ and $\| \cdot \|_2$, we can show easily that f is continuous at a iff

for every $\epsilon > 0$, there is some $\eta > 0$, such that, for every $x \in E$,

$$\text{if } \|x - a\|_1 \leq \eta, \text{ then } \|f(x) - f(a)\|_2 \leq \epsilon.$$

It is worth noting that continuity is a topological notion, in the sense that equivalent metrics (or equivalent norms) define exactly the same notion of continuity.

If (E, \mathcal{O}_E) and (F, \mathcal{O}_F) are topological spaces, and $f: E \rightarrow F$ is a function, for every nonempty subset $A \subseteq E$ of E , we say that f is *continuous on A* if the restriction of f to A is continuous with respect to (A, \mathcal{U}) and (F, \mathcal{O}_F) , where \mathcal{U} is the subspace topology induced by \mathcal{O}_E on A .

Given a product $E_1 \times \cdots \times E_n$ of topological spaces, as usual, we let $\pi_i: E_1 \times \cdots \times E_n \rightarrow E_i$ be the projection function such that, $\pi_i(x_1, \dots, x_n) = x_i$. It is immediately verified that each π_i is continuous.

Given a topological space (E, \mathcal{O}) , we say that a point $a \in E$ is *isolated* if $\{a\}$ is an open set in \mathcal{O} . Then, if (E, \mathcal{O}_E) and (F, \mathcal{O}_F) are topological spaces, any function $f: E \rightarrow F$ is continuous at every isolated point $a \in E$. In the discrete topology, every point is isolated.

In a nontrivial normed vector space $(E, \|\cdot\|)$ (with $E \neq \{0\}$), no point is isolated. To show this, we show that every open ball $B_0(u, \rho)$ contains some vectors different from u . Indeed, since E is nontrivial, there is some $v \in E$ such that $v \neq 0$, and thus $\lambda = \|v\| > 0$ (by (N1)). Let

$$w = u + \frac{\rho}{\lambda + 1}v.$$

Since $v \neq 0$ and $\rho > 0$, we have $w \neq u$. Then,

$$\|w - u\| = \left\| \frac{\rho}{\lambda + 1}v \right\| = \frac{\rho\lambda}{\lambda + 1} < \rho,$$

which shows that $\|w - u\| < \rho$, for $w \neq u$.

The following proposition is easily shown.

Proposition 26.8. *Given topological spaces (E, \mathcal{O}_E) , (F, \mathcal{O}_F) , and (G, \mathcal{O}_G) , and two functions $f: E \rightarrow F$ and $g: F \rightarrow G$, if f is continuous at $a \in E$ and g is continuous at $f(a) \in F$, then $g \circ f: E \rightarrow G$ is continuous at $a \in E$. Given n topological spaces (F_i, \mathcal{O}_i) , for every function $f: E \rightarrow F_1 \times \cdots \times F_n$, then f is continuous at $a \in E$ iff every $f_i: E \rightarrow F_i$ is continuous at a , where $f_i = \pi_i \circ f$.*

One can also show that in a metric space (E, d) , the norm $d: E \times E \rightarrow \mathbb{R}$ is continuous, where $E \times E$ has the product topology, and that for a normed vector space $(E, \|\cdot\|)$, the norm $\|\cdot\|: E \rightarrow \mathbb{R}$ is continuous.

Given a function $f: E_1 \times \cdots \times E_n \rightarrow F$, we can fix $n - 1$ of the arguments, say $a_1, \dots, a_{i-1}, a_{i+1}, \dots, a_n$, and view f as a function of the remaining argument,

$$x_i \mapsto f(a_1, \dots, a_{i-1}, x_i, a_{i+1}, \dots, a_n),$$

where $x_i \in E_i$. If f is continuous, it is clear that each f_i is continuous.



One should be careful that the converse is false! For example, consider the function $f: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, defined such that,

$$f(x, y) = \frac{xy}{x^2 + y^2} \quad \text{if } (x, y) \neq (0, 0), \quad \text{and} \quad f(0, 0) = 0.$$

The function f is continuous on $\mathbb{R} \times \mathbb{R} - \{(0, 0)\}$, but on the line $y = mx$, with $m \neq 0$, we have $f(x, y) = \frac{m}{1+m^2} \neq 0$, and thus, on this line, $f(x, y)$ does not approach 0 when (x, y) approaches $(0, 0)$.

The following proposition is useful for showing that real-valued functions are continuous.

Proposition 26.9. *If E is a topological space, and $(\mathbb{R}, |x - y|)$ the reals under the standard topology, for any two functions $f: E \rightarrow \mathbb{R}$ and $g: E \rightarrow \mathbb{R}$, for any $a \in E$, for any $\lambda \in \mathbb{R}$, if f and g are continuous at a , then $f + g$, λf , $f \cdot g$, are continuous at a , and f/g is continuous at a if $g(a) \neq 0$.*

Proof. Left as an exercise. □

Using Proposition 26.9, we can show easily that every real polynomial function is continuous.

The notion of isomorphism of topological spaces is defined as follows.

Definition 26.11. Let (E, \mathcal{O}_E) and (F, \mathcal{O}_F) be topological spaces, and let $f: E \rightarrow F$ be a function. We say that f is a *homeomorphism between E and F* if f is bijective, and both $f: E \rightarrow F$ and $f^{-1}: F \rightarrow E$ are continuous.



One should be careful that a bijective continuous function $f: E \rightarrow F$ is not necessarily an homeomorphism. For example, if $E = \mathbb{R}$ with the discrete topology, and $F = \mathbb{R}$ with the standard topology, the identity is not a homeomorphism. Another interesting example involving a parametric curve is given below. Let $L: \mathbb{R} \rightarrow \mathbb{R}^2$ be the function, defined such that,

$$\begin{aligned} L_1(t) &= \frac{t(1+t^2)}{1+t^4}, \\ L_2(t) &= \frac{t(1-t^2)}{1+t^4}. \end{aligned}$$

If we think of $(x(t), y(t)) = (L_1(t), L_2(t))$ as a geometric point in \mathbb{R}^2 , the set of points $(x(t), y(t))$ obtained by letting t vary in \mathbb{R} from $-\infty$ to $+\infty$, defines a curve having the shape of a “figure eight”, with self-intersection at the origin, called the “lemniscate of Bernoulli”. The map L is continuous, and in fact bijective, but its inverse L^{-1} is not continuous. Indeed, when we approach the origin on the branch of the curve in the upper left quadrant (i.e., points such that, $x \leq 0$, $y \geq 0$), then t goes to $-\infty$, and when we approach the origin on the

branch of the curve in the lower right quadrant (i.e., points such that, $x \geq 0$, $y \leq 0$), then t goes to $+\infty$.

We also review the concept of limit of a sequence. Given any set E , a *sequence* is any function $x: \mathbb{N} \rightarrow E$, usually denoted by $(x_n)_{n \in \mathbb{N}}$, or $(x_n)_{n \geq 0}$, or even by (x_n) .

Definition 26.12. Given a topological space (E, \mathcal{O}) , we say that a *sequence* $(x_n)_{n \in \mathbb{N}}$ *converges to some* $a \in E$ if for every open set U containing a , there is some $n_0 \geq 0$, such that, $x_n \in U$, for all $n \geq n_0$. We also say that a is a *limit of* $(x_n)_{n \in \mathbb{N}}$.

When E is a metric space with metric d , it is easy to show that this is equivalent to the fact that,

for every $\epsilon > 0$, there is some $n_0 \geq 0$, such that, $d(x_n, a) \leq \epsilon$, for all $n \geq n_0$.

When E is a normed vector space with norm $\| \cdot \|$, it is easy to show that this is equivalent to the fact that,

for every $\epsilon > 0$, there is some $n_0 \geq 0$, such that, $\|x_n - a\| \leq \epsilon$, for all $n \geq n_0$.

The following proposition shows the importance of the Hausdorff separation axiom.

Proposition 26.10. *Given a topological space (E, \mathcal{O}) , if the Hausdorff separation axiom holds, then every sequence has at most one limit.*

Proof. Left as an exercise. □

It is worth noting that the notion of limit is topological, in the sense that a sequence converge to a limit b iff it converges to the same limit b in any equivalent metric (and similarly for equivalent norms).

We still need one more concept of limit for functions.

Definition 26.13. Let (E, \mathcal{O}_E) and (F, \mathcal{O}_F) be topological spaces, let A be some nonempty subset of E , and let $f: A \rightarrow F$ be a function. For any $a \in \overline{A}$ and any $b \in F$, we say that $f(x)$ *approaches* b *as* x *approaches* a *with values in* A if for every open set $V \in \mathcal{O}_F$ containing b , there is some open set $U \in \mathcal{O}_E$ containing a , such that, $f(U \cap A) \subseteq V$. This is denoted by

$$\lim_{x \rightarrow a, x \in A} f(x) = b.$$

First, note that by Proposition 26.2, since $a \in \overline{A}$, for every open set U containing a , we have $U \cap A \neq \emptyset$, and the definition is nontrivial. Also, even if $a \in A$, the value $f(a)$ of f at a plays no role in this definition. When E and F are metric space with metrics d_1 and d_2 , it can be shown easily that the definition can be stated as follows:

For every $\epsilon > 0$, there is some $\eta > 0$, such that, for every $x \in A$,

$$\text{if } d_1(x, a) \leq \eta, \text{ then } d_2(f(x), b) \leq \epsilon.$$

When E and F are normed vector spaces with norms $\|\cdot\|_1$ and $\|\cdot\|_2$, it can be shown easily that the definition can be stated as follows:

For every $\epsilon > 0$, there is some $\eta > 0$, such that, for every $x \in A$,

$$\text{if } \|x - a\|_1 \leq \eta, \text{ then } \|f(x) - b\|_2 \leq \epsilon.$$

We have the following result relating continuity at a point and the previous notion.

Proposition 26.11. *Let (E, \mathcal{O}_E) and (F, \mathcal{O}_F) be two topological spaces, and let $f: E \rightarrow F$ be a function. For any $a \in E$, the function f is continuous at a iff $f(x)$ approaches $f(a)$ when x approaches a (with values in E).*

Proof. Left as a trivial exercise. □

Another important proposition relating the notion of convergence of a sequence to continuity, is stated without proof.

Proposition 26.12. *Let (E, \mathcal{O}_E) and (F, \mathcal{O}_F) be two topological spaces, and let $f: E \rightarrow F$ be a function.*

- (1) *If f is continuous, then for every sequence $(x_n)_{n \in \mathbb{N}}$ in E , if (x_n) converges to a , then $(f(x_n))$ converges to $f(a)$.*
- (2) *If E is a metric space, and $(f(x_n))$ converges to $f(a)$ whenever (x_n) converges to a , for every sequence $(x_n)_{n \in \mathbb{N}}$ in E , then f is continuous.*

A special case of Definition 26.13 will be used when E and F are (nontrivial) normed vector spaces with norms $\|\cdot\|_1$ and $\|\cdot\|_2$. Let U be any nonempty open subset of E . We showed earlier that E has no isolated points and that every set $\{v\}$ is closed, for every $v \in E$. Since E is nontrivial, for every $v \in U$, there is a nontrivial open ball contained in U (an open ball not reduced to its center). Then, for every $v \in U$, $A = U - \{v\}$ is open and nonempty, and clearly, $v \in \overline{A}$. For any $v \in U$, if $f(x)$ approaches b when x approaches v with values in $A = U - \{v\}$, we say that $f(x)$ approaches b when x approaches v with values $\neq v$ in U . This is denoted by

$$\lim_{x \rightarrow v, x \in U, x \neq v} f(x) = b.$$

Remark: Variations of the above case show up in the following case: $E = \mathbb{R}$, and F is some arbitrary topological space. Let A be some nonempty subset of \mathbb{R} , and let $f: A \rightarrow F$ be some function. For any $a \in A$, we say that f is continuous on the right at a if

$$\lim_{x \rightarrow a, x \in A \cap [a, +\infty[} f(x) = f(a).$$

We can define continuity on the left at a in a similar fashion.

Let us consider another variation. Let A be some nonempty subset of \mathbb{R} , and let $f: A \rightarrow F$ be some function. For any $a \in A$, we say that f has a *discontinuity of the first kind at a* if

$$\lim_{x \rightarrow a, x \in A \cap]-\infty, a[} f(x) = f(a_-)$$

and

$$\lim_{x \rightarrow a, x \in A \cap]a, +\infty[} f(x) = f(a_+)$$

both exist, and either $f(a_-) \neq f(a)$, or $f(a_+) \neq f(a)$.

Note that it is possible that $f(a_-) = f(a_+)$, but f is still discontinuous at a if this common value differs from $f(a)$. Functions defined on a nonempty subset of \mathbb{R} , and that are continuous, except for some points of discontinuity of the first kind, play an important role in analysis.

We now turn to connectivity properties of topological spaces.

26.4 Connected Sets

Connectivity properties of topological spaces play a very important role in understanding the topology of surfaces. This section gathers the facts needed to have a good understanding of the classification theorem for compact surfaces (with boundary). The main references are Ahlfors and Sario [1] and Massey [77, 78]. For general background on topology, geometry, and algebraic topology, we also highly recommend Bredon [18] and Fulton [41].

Definition 26.14. A topological space, (E, \mathcal{O}) , is *connected* if the only subsets of E that are both open and closed are the empty set and E itself. Equivalently, (E, \mathcal{O}) is connected if E cannot be written as the union, $E = U \cup V$, of two disjoint nonempty open sets, U, V , if E cannot be written as the union, $E = U \cup V$, of two disjoint nonempty closed sets. A subset, $S \subseteq E$, is *connected* if it is connected in the subspace topology on S induced by (E, \mathcal{O}) . A connected open set is called a *region* and a closed set is a *closed region* if its interior is a connected (open) set.

Intuitively, if a space is not connected, it is possible to define a continuous function which is constant on disjoint “connected components” and which takes possibly distinct values on disjoint components. This can be stated in terms of the concept of a locally constant function. Given two topological spaces, X, Y , a function, $f: X \rightarrow Y$, is *locally constant* if for every $x \in X$, there is an open set, $U \subseteq X$, such that $x \in U$ and f is constant on U .

We claim that a locally constant function is continuous. In fact, we will prove that $f^{-1}(V)$ is open for every subset, $V \subseteq Y$ (not just for an open set V). It is enough to show that $f^{-1}(y)$ is open for every $y \in Y$, since for every subset $V \subseteq Y$,

$$f^{-1}(V) = \bigcup_{y \in V} f^{-1}(y),$$

and open sets are closed under arbitrary unions. However, either $f^{-1}(y) = \emptyset$ if $y \in Y - f(X)$ or f is constant on $U = f^{-1}(y)$ if $y \in f(X)$ (with value y), and since f is locally constant, for every $x \in U$, there is some open set, $W \subseteq X$, such that $x \in W$ and f is constant on W , which implies that $f(w) = y$ for all $w \in W$ and thus, that $W \subseteq U$, showing that U is a union of open sets and thus, is open. The following proposition shows that a space is connected iff every locally constant function is constant:

Proposition 26.13. *A topological space is connected iff every locally constant function is constant.*

Proof. First, assume that X is connected. Let $f: X \rightarrow Y$ be a locally constant function to some space Y and assume that f is not constant. Pick any $y \in f(Y)$. Since f is not constant, $U_1 = f^{-1}(y) \neq X$, and of course, $U_1 \neq \emptyset$. We proved just before Proposition 26.13 that $f^{-1}(V)$ is open for every subset $V \subseteq Y$, and thus $U_1 = f^{-1}(y) = f^{-1}(\{y\})$ and $U_2 = f^{-1}(Y - \{y\})$ are both open, nonempty, and clearly $X = U_1 \cup U_2$ and U_1 and U_2 are disjoint. This contradicts the fact that X is connected and f must be constant.

Assume that every locally constant function, $f: X \rightarrow Y$, to a Hausdorff space, Y , is constant. If X is not connected, we can write $X = U_1 \cup U_2$, where both U_1, U_2 are open, disjoint, and nonempty. We can define the function, $f: X \rightarrow \mathbb{R}$, such that $f(x) = 1$ on U_1 and $f(x) = 0$ on U_2 . Since U_1 and U_2 are open, the function f is locally constant, and yet not constant, a contradiction. \square

The following standard proposition characterizing the connected subsets of \mathbb{R} can be found in most topology texts (for example, Munkres [84], Schwartz [93]). For the sake of completeness, we give a proof.

Proposition 26.14. *A subset of the real line, \mathbb{R} , is connected iff it is an interval, i.e., of the form $[a, b]$, $]a, b]$, where $a = -\infty$ is possible, $[a, b[$, where $b = +\infty$ is possible, or $]a, b[$, where $a = -\infty$ or $b = +\infty$ is possible.*

Proof. Assume that A is a connected nonempty subset of \mathbb{R} . The cases where $A = \emptyset$ or A consists of a single point are trivial. We show that whenever $a, b \in A$, $a < b$, then the entire interval $[a, b]$ is a subset of A . Indeed, if this was not the case, there would be some $c \in]a, b[$ such that $c \notin A$, and then we could write $A = (]-\infty, c[\cap A) \cup (]c, +\infty[\cap A)$, where $] - \infty, c[\cap A$ and $]c, +\infty[\cap A$ are nonempty and disjoint open subsets of A , contradicting the fact that A is connected. It follows easily that A must be an interval.

Conversely, we show that an interval, I , must be connected. Let A be any nonempty subset of I which is both open and closed in I . We show that $I = A$. Fix any $x \in A$ and consider the set, R_x , of all y such that $[x, y] \subseteq A$. If the set R_x is unbounded, then $R_x = [x, +\infty[$. Otherwise, if this set is bounded, let b be its least upper bound. We claim that b is the right boundary of the interval I . Because A is closed in I , unless I is open on the right and b is its right boundary, we must have $b \in A$. In the first case, $A \cap [x, b[= I \cap [x, b[= [x, b[$. In the second case, because A is also open in I , unless b is the

right boundary of the interval I (closed on the right), there is some open set $]b - \eta, b + \eta[$ contained in A , which implies that $[x, b + \eta/2] \subseteq A$, contradicting the fact that b is the least upper bound of the set R_x . Thus, b must be the right boundary of the interval I (closed on the right). A similar argument applies to the set, L_y , of all x such that $[x, y] \subseteq A$ and either L_y is unbounded, or its greatest lower bound a is the left boundary of I (open or closed on the left). In all cases, we showed that $A = I$, and the interval must be connected. \square

A characterization on the connected subsets of \mathbb{R}^n is harder and requires the notion of arcwise connectedness. One of the most important properties of connected sets is that they are preserved by continuous maps.

Proposition 26.15. *Given any continuous map, $f: E \rightarrow F$, if $A \subseteq E$ is connected, then $f(A)$ is connected.*

Proof. If $f(A)$ is not connected, then there exist some nonempty open sets, U, V , in F such that $f(A) \cap U$ and $f(A) \cap V$ are nonempty and disjoint, and

$$f(A) = (f(A) \cap U) \cup (f(A) \cap V).$$

Then, $f^{-1}(U)$ and $f^{-1}(V)$ are nonempty and open since f is continuous and

$$A = (A \cap f^{-1}(U)) \cup (A \cap f^{-1}(V)),$$

with $A \cap f^{-1}(U)$ and $A \cap f^{-1}(V)$ nonempty, disjoint, and open in A , contradicting the fact that A is connected. \square

An important corollary of Proposition 26.15 is that for every continuous function, $f: E \rightarrow \mathbb{R}$, where E is a connected space, $f(E)$ is an interval. Indeed, this follows from Proposition 26.14. Thus, if f takes the values a and b where $a < b$, then f takes all values $c \in [a, b]$. This is a very important property.

Even if a topological space is not connected, it turns out that it is the disjoint union of maximal connected subsets and these connected components are closed in E . In order to obtain this result, we need a few lemmas.

Lemma 26.16. *Given a topological space, E , for any family, $(A_i)_{i \in I}$, of (nonempty) connected subsets of E , if $A_i \cap A_j \neq \emptyset$ for all $i, j \in I$, then the union, $A = \bigcup_{i \in I} A_i$, of the family, $(A_i)_{i \in I}$, is also connected.*

Proof. Assume that $\bigcup_{i \in I} A_i$ is not connected. Then, there exists two nonempty open subsets, U and V , of E such that $A \cap U$ and $A \cap V$ are disjoint and nonempty and such that

$$A = (A \cap U) \cup (A \cap V).$$

Now, for every $i \in I$, we can write

$$A_i = (A_i \cap U) \cup (A_i \cap V),$$

where $A_i \cap U$ and $A_i \cap V$ are disjoint, since $A_i \subseteq A$ and $A \cap U$ and $A \cap V$ are disjoint. Since A_i is connected, either $A_i \cap U = \emptyset$ or $A_i \cap V = \emptyset$. This implies that either $A_i \subseteq A \cap U$ or $A_i \subseteq A \cap V$. However, by assumption, $A_i \cap A_j \neq \emptyset$, for all $i, j \in I$, and thus, either both $A_i \subseteq A \cap U$ and $A_j \subseteq A \cap U$, or both $A_i \subseteq A \cap V$ and $A_j \subseteq A \cap V$, since $A \cap U$ and $A \cap V$ are disjoint. Thus, we conclude that either $A_i \subseteq A \cap U$ for all $i \in I$, or $A_i \subseteq A \cap V$ for all $i \in I$. But this proves that either

$$A = \bigcup_{i \in I} A_i \subseteq A \cap U,$$

or

$$A = \bigcup_{i \in I} A_i \subseteq A \cap V,$$

contradicting the fact that both $A \cap U$ and $A \cap V$ are disjoint and nonempty. Thus, A must be connected. \square

In particular, the above lemma applies when the connected sets in a family $(A_i)_{i \in I}$ have a point in common.

Lemma 26.17. *If A is a connected subset of a topological space, E , then for every subset, B , such that $A \subseteq B \subseteq \overline{A}$, where \overline{A} is the closure of A in E , the set B is connected.*

Proof. If B is not connected, then there are two nonempty open subsets, U, V , of E such that $B \cap U$ and $B \cap V$ are disjoint and nonempty, and

$$B = (B \cap U) \cup (B \cap V).$$

Since $A \subseteq B$, the above implies that

$$A = (A \cap U) \cup (A \cap V),$$

and since A is connected, either $A \cap U = \emptyset$, or $A \cap V = \emptyset$. Without loss of generality, assume that $A \cap V = \emptyset$, which implies that $A \subseteq A \cap U \subseteq B \cap U$. However, $B \cap U$ is closed in the subspace topology for B and since $B \subseteq \overline{A}$ and \overline{A} is closed in E , the closure of A in B w.r.t. the subspace topology of B is clearly $B \cap \overline{A} = B$, which implies that $B \subseteq B \cap U$ (since the closure is the smallest closed set containing the given set). Thus, $B \cap V = \emptyset$, a contradiction. \square

In particular, Lemma 26.17 shows that if A is a connected subset, then its closure, \overline{A} , is also connected. We are now ready to introduce the connected components of a space.

Definition 26.15. Given a topological space, (E, \mathcal{O}) , we say that two points, $a, b \in E$, are *connected* if there is some connected subset, A , of E such that $a \in A$ and $b \in A$.

It is immediately verified that the relation “ a and b are connected in E ” is an equivalence relation. Only transitivity is not obvious, but it follows immediately as a special case of Lemma 26.16. Thus, the above equivalence relation defines a partition of E into nonempty disjoint *connected components*. The following proposition is easily proved using Lemma 26.16 and Lemma 26.17:

Proposition 26.18. *Given any topological space, E , for any $a \in E$, the connected component containing a is the largest connected set containing a . The connected components of E are closed.*

The notion of a locally connected space is also useful.

Definition 26.16. A topological space, (E, \mathcal{O}) , is *locally connected* if for every $a \in E$, for every neighborhood, V , of a , there is a connected neighborhood, U , of a such that $U \subseteq V$.

As we shall see in a moment, it would be equivalent to require that E has a basis of connected open sets.



There are connected spaces that are not locally connected and there are locally connected spaces that are not connected. The two properties are independent.

Proposition 26.19. *A topological space, E , is locally connected iff for every open subset, A , of E , the connected components of A are open.*

Proof. Assume that E is locally connected. Let A be any open subset of E and let C be one of the connected components of A . For any $a \in C \subseteq A$, there is some connected neighborhood, U , of a such that $U \subseteq A$ and since C is a connected component of A containing a , we must have $U \subseteq C$. This shows that for every $a \in C$, there is some open subset containing a contained in C , so C is open.

Conversely, assume that for every open subset, A , of E , the connected components of A are open. Then, for every $a \in E$ and every neighborhood, U , of a , since U contains some open set A containing a , the interior, $\overset{\circ}{U}$, of U is an open set containing a and its connected components are open. In particular, the connected component C containing a is a connected open set containing a and contained in U . \square

Proposition 26.19 shows that in a locally connected space, the connected open sets form a basis for the topology. It is easily seen that \mathbb{R}^n is locally connected. Another very important property of surfaces and more generally, manifolds, is to be arcwise connected. The intuition is that any two points can be joined by a continuous arc of curve. This is formalized as follows.

Definition 26.17. Given a topological space, (E, \mathcal{O}) , an *arc (or path)* is a continuous map, $\gamma: [a, b] \rightarrow E$, where $[a, b]$ is a closed interval of the real line, \mathbb{R} . The point $\gamma(a)$ is the *initial point* of the arc and the point $\gamma(b)$ is the *terminal point* of the arc. We say that γ is an *arc joining* $\gamma(a)$ and $\gamma(b)$. An arc is a *closed curve* if $\gamma(a) = \gamma(b)$. The set $\gamma([a, b])$ is the *trace* of the arc γ .

Typically, $a = 0$ and $b = 1$. In the sequel, this will be assumed.



One should not confuse an arc, $\gamma: [a, b] \rightarrow E$, with its trace. For example, γ could be constant, and thus, its trace reduced to a single point.

An arc is a *Jordan arc* if γ is a homeomorphism onto its trace. An arc, $\gamma: [a, b] \rightarrow E$, is a *Jordan curve* if $\gamma(a) = \gamma(b)$ and γ is injective on $[a, b[$. Since $[a, b]$ is connected, by Proposition 26.15, the trace $\gamma([a, b])$ of an arc is a connected subset of E .

Given two arcs $\gamma: [0, 1] \rightarrow E$ and $\delta: [0, 1] \rightarrow E$ such that $\gamma(1) = \delta(0)$, we can form a new arc defined as follows:

Definition 26.18. Given two arcs, $\gamma: [0, 1] \rightarrow E$ and $\delta: [0, 1] \rightarrow E$, such that $\gamma(1) = \delta(0)$, we can form their *composition (or product)*, $\gamma\delta$, defined such that

$$\gamma\delta(t) = \begin{cases} \gamma(2t) & \text{if } 0 \leq t \leq 1/2; \\ \delta(2t - 1) & \text{if } 1/2 \leq t \leq 1. \end{cases}$$

The *inverse*, γ^{-1} , of the arc, γ , is the arc defined such that $\gamma^{-1}(t) = \gamma(1 - t)$, for all $t \in [0, 1]$.

It is trivially verified that Definition 26.18 yields continuous arcs.

Definition 26.19. A topological space, E , is *arcwise connected* if for any two points, $a, b \in E$, there is an arc, $\gamma: [0, 1] \rightarrow E$, joining a and b , i.e., such that $\gamma(0) = a$ and $\gamma(1) = b$. A topological space, E , is *locally arcwise connected* if for every $a \in E$, for every neighborhood, V , of a , there is an arcwise connected neighborhood, U , of a such that $U \subseteq V$.

The space \mathbb{R}^n is locally arcwise connected, since for any open ball, any two points in this ball are joined by a line segment. Manifolds and surfaces are also locally arcwise connected. Proposition 26.15 also applies to arcwise connectedness (this is a simple exercise). The following theorem is crucial to the theory of manifolds and surfaces:

Theorem 26.20. *If a topological space, E , is arcwise connected, then it is connected. If a topological space, E , is connected and locally arcwise connected, then E is arcwise connected.*

Proof. First, assume that E is arcwise connected. Pick any point, a , in E . Since E is arcwise connected, for every $b \in E$, there is a path, $\gamma_b: [0, 1] \rightarrow E$, from a to b and so,

$$E = \bigcup_{b \in E} \gamma_b([0, 1])$$

a union of connected subsets all containing a . By Lemma 26.16, E is connected.

Now assume that E is connected and locally arcwise connected. For any point $a \in E$, let F_a be the set of all points, b , such that there is an arc, $\gamma_b: [0, 1] \rightarrow E$, from a to b . Clearly, F_a contains a . We show that F_a is both open and closed. For any $b \in F_a$, since E is locally arcwise connected, there is an arcwise connected neighborhood U containing b (because E is

a neighborhood of b). Thus, b can be joined to every point $c \in U$ by an arc, and since by the definition of F_a , there is an arc from a to b , the composition of these two arcs yields an arc from a to c , which shows that $c \in F_a$. But then $U \subseteq F_a$ and thus, F_a is open. Now assume that b is in the complement of F_a . As in the previous case, there is some arcwise connected neighborhood U containing b . Thus, every point $c \in U$ can be joined to b by an arc. If there was an arc joining a to c , we would get an arc from a to b , contradicting the fact that b is in the complement of F_a . Thus, every point $c \in U$ is in the complement of F_a , which shows that U is contained in the complement of F_a , and thus, that the complement of F_a is open. Consequently, we have shown that F_a is both open and closed and since it is nonempty, we must have $E = F_a$, which shows that E is arcwise connected. \square

If E is locally arcwise connected, the above argument shows that the connected components of E are arcwise connected.



It is not true that a connected space is arcwise connected. For example, the space consisting of the graph of the function

$$f(x) = \sin(1/x),$$

where $x > 0$, together with the portion of the y -axis, for which $-1 \leq y \leq 1$, is connected, but not arcwise connected.

A trivial modification of the proof of Theorem 26.20 shows that in a normed vector space, E , a connected open set is arcwise connected by polygonal lines (i.e., arcs consisting of line segments). This is because in every open ball, any two points are connected by a line segment. Furthermore, if E is finite dimensional, these polygonal lines can be forced to be parallel to basis vectors.

We now consider compactness.

26.5 Compact Sets

The property of compactness is very important in topology and analysis. We provide a quick review geared towards the study of surfaces and for details, we refer the reader to Munkres [84], Schwartz [93]. In this section, we will need to assume that the topological spaces are Hausdorff spaces. This is not a luxury, as many of the results are false otherwise.

There are various equivalent ways of defining compactness. For our purposes, the most convenient way involves the notion of open cover.

Definition 26.20. Given a topological space, E , for any subset, A , of E , an *open cover*, $(U_i)_{i \in I}$, of A is a family of open subsets of E such that $A \subseteq \bigcup_{i \in I} U_i$. An *open subcover* of an open cover, $(U_i)_{i \in I}$, of A is any subfamily, $(U_j)_{j \in J}$, which is an open cover of A , with $J \subseteq I$. An open cover, $(U_i)_{i \in I}$, of A is *finite* if I is finite. The topological space, E , is *compact* if it

is Hausdorff and for every open cover, $(U_i)_{i \in I}$, of E , there is a finite open subcover, $(U_j)_{j \in J}$, of E . Given any subset, A , of E , we say that A is *compact* if it is compact with respect to the subspace topology. We say that A is *relatively compact* if its closure \bar{A} is compact.

It is immediately verified that a subset, A , of E is compact in the subspace topology relative to A iff for every open cover, $(U_i)_{i \in I}$, of A by open subsets of E , there is a finite open subcover, $(U_j)_{j \in J}$, of A . The property that every open cover contains a finite open subcover is often called the *Heine-Borel-Lebesgue* property. By considering complements, a Hausdorff space is compact iff for every family, $(F_i)_{i \in I}$, of closed sets, if $\bigcap_{i \in I} F_i = \emptyset$, then $\bigcap_{j \in J} F_j = \emptyset$ for some finite subset, J , of I .



Definition 26.20 requires that a compact space be Hausdorff. There are books in which a compact space is not necessarily required to be Hausdorff. Following Schwartz, we prefer calling such a space *quasi-compact*.

Another equivalent and useful characterization can be given in terms of families having the finite intersection property. A family, $(F_i)_{i \in I}$, of sets has the *finite intersection property* if $\bigcap_{j \in J} F_j \neq \emptyset$ for every finite subset, J , of I . We have the following proposition:

Proposition 26.21. *A topological Hausdorff space, E , is compact iff for every family, $(F_i)_{i \in I}$, of closed sets having the finite intersection property, then $\bigcap_{i \in I} F_i \neq \emptyset$.*

Proof. If E is compact and $(F_i)_{i \in I}$ is a family of closed sets having the finite intersection property, then $\bigcap_{i \in I} F_i$ cannot be empty, since otherwise we would have $\bigcap_{j \in J} F_j = \emptyset$ for some finite subset, J , of I , a contradiction. The converse is equally obvious. \square

Another useful consequence of compactness is as follows. For any family, $(F_i)_{i \in I}$, of closed sets such that $F_{i+1} \subseteq F_i$ for all $i \in I$, if $\bigcap_{i \in I} F_i = \emptyset$, then $F_i = \emptyset$ for some $i \in I$. Indeed, there must be some finite subset, J , of I such that $\bigcap_{j \in J} F_j = \emptyset$ and since $F_{i+1} \subseteq F_i$ for all $i \in I$, we must have $F_j = \emptyset$ for the smallest F_j in $(F_j)_{j \in J}$. Using this fact, we note that \mathbb{R} is *not* compact. Indeed, the family of closed sets, $([n, +\infty[)_{n \geq 0}$, is decreasing and has an empty intersection.

Given a metric space, if we define a *bounded subset* to be a subset that can be enclosed in some closed ball (of finite radius), then any nonbounded subset of a metric space is not compact. However, a closed interval $[a, b]$ of the real line is compact.

Proposition 26.22. *Every closed interval, $[a, b]$, of the real line is compact.*

Proof. We proceed by contradiction. Let $(U_i)_{i \in I}$ be any open cover of $[a, b]$ and assume that there is no finite open subcover. Let $c = (a + b)/2$. If both $[a, c]$ and $[c, b]$ had some finite open subcover, so would $[a, b]$, and thus, either $[a, c]$ does not have any finite subcover, or $[c, b]$ does not have any finite open subcover. Let $[a_1, b_1]$ be such a bad subinterval. The same argument applies and we split $[a_1, b_1]$ into two equal subintervals, one of which must be bad. Thus, having defined $[a_n, b_n]$ of length $(b - a)/2^n$ as an interval having no finite open

subcover, splitting $[a_n, b_n]$ into two equal intervals, we know that at least one of the two has no finite open subcover and we denote such a bad interval by $[a_{n+1}, b_{n+1}]$. The sequence (a_n) is nondecreasing and bounded from above by b , and thus, by a fundamental property of the real line, it converges to its least upper bound, α . Similarly, the sequence (b_n) is nonincreasing and bounded from below by a and thus, it converges to its greatest lower bound, β . Since $[a_n, b_n]$ has length $(b - a)/2^n$, we must have $\alpha = \beta$. However, the common limit $\alpha = \beta$ of the sequences (a_n) and (b_n) must belong to some open set, U_i , of the open cover and since U_i is open, it must contain some interval $[c, d]$ containing α . Then, because α is the common limit of the sequences (a_n) and (b_n) , there is some N such that the intervals $[a_n, b_n]$ are all contained in the interval $[c, d]$ for all $n \geq N$, which contradicts the fact that none of the intervals $[a_n, b_n]$ has a finite open subcover. Thus, $[a, b]$ is indeed compact. \square

The argument of Proposition 26.22 can be adapted to show that in \mathbb{R}^m , every closed set, $[a_1, b_1] \times \cdots \times [a_m, b_m]$, is compact. At every stage, we need to divide into 2^m subpieces instead of 2.

The following two propositions give very important properties of the compact sets, and they only hold for Hausdorff spaces:

Proposition 26.23. *Given a topological Hausdorff space, E , for every compact subset, A , and every point, b , not in A , there exist disjoint open sets, U and V , such that $A \subseteq U$ and $b \in V$. As a consequence, every compact subset is closed.*

Proof. Since E is Hausdorff, for every $a \in A$, there are some disjoint open sets, U_a and V_a , containing a and b respectively. Thus, the family, $(U_a)_{a \in A}$, forms an open cover of A . Since A is compact there is a finite open subcover, $(U_j)_{j \in J}$, of A , where $J \subseteq A$, and then $\bigcup_{j \in J} U_j$ is an open set containing A disjoint from the open set $\bigcap_{j \in J} V_j$ containing b . This shows that every point, b , in the complement of A belongs to some open set in this complement and thus, that the complement is open, i.e., that A is closed. \square

Actually, the proof of Proposition 26.23 can be used to show the following useful property:

Proposition 26.24. *Given a topological Hausdorff space, E , for every pair of compact disjoint subsets, A and B , there exist disjoint open sets, U and V , such that $A \subseteq U$ and $B \subseteq V$.*

Proof. We repeat the argument of Proposition 26.23 with B playing the role of b and use Proposition 26.23 to find disjoint open sets, U_a , containing $a \in A$ and, V_a , containing B . \square

The following proposition shows that in a compact topological space, every closed set is compact:

Proposition 26.25. *Given a compact topological space, E , every closed set is compact.*

Proof. Since A is closed, $E - A$ is open and from any open cover, $(U_i)_{i \in I}$, of A , we can form an open cover of E by adding $E - A$ to $(U_i)_{i \in I}$ and, since E is compact, a finite subcover, $(U_j)_{j \in J} \cup \{E - A\}$, of E can be extracted such that $(U_j)_{j \in J}$ is a finite subcover of A . \square

Remark: Proposition 26.25 also holds for quasi-compact spaces, i.e., the Hausdorff separation property is not needed.

Putting Proposition 26.24 and Proposition 26.25 together, we note that if X is compact, then for every pair of disjoint closed, sets A and B , there exist disjoint open sets, U and V , such that $A \subseteq U$ and $B \subseteq V$. We say that X is a *normal* space.

Proposition 26.26. *Given a compact topological space, E , for every $a \in E$, for every neighborhood, V , of a , there exists a compact neighborhood, U , of a such that $U \subseteq V$*

Proof. Since V is a neighborhood of a , there is some open subset, O , of V containing a . Then the complement, $K = E - O$, of O is closed and since E is compact, by Proposition 26.25, K is compact. Now, if we consider the family of all closed sets of the form, $K \cap F$, where F is any closed neighborhood of a , since $a \notin K$, this family has an empty intersection and thus, there is a finite number of closed neighborhood, F_1, \dots, F_n , of a , such that $K \cap F_1 \cap \dots \cap F_n = \emptyset$. Then, $U = F_1 \cap \dots \cap F_n$ is a compact neighborhood of a contained in $O \subseteq V$. \square

It can be shown that in a normed vector space of finite dimension, a subset is compact iff it is closed and bounded. For \mathbb{R}^n , the proof is simple.



In a normed vector space of infinite dimension, there are closed and bounded sets that are not compact!

More could be said about compactness in metric spaces but we will only need the notion of Lebesgue number, which will be discussed a little later. Another crucial property of compactness is that it is preserved under continuity.

Proposition 26.27. *Let E be a topological space and let F be a topological Hausdorff space. For every compact subset, A , of E , for every continuous map, $f: E \rightarrow F$, the subspace $f(A)$ is compact.*

Proof. Let $(U_i)_{i \in I}$ be an open cover of $f(A)$. We claim that $(f^{-1}(U_i))_{i \in I}$ is an open cover of A , which is easily checked. Since A is compact, there is a finite open subcover, $(f^{-1}(U_j))_{j \in J}$, of A , and thus, $(U_j)_{j \in J}$ is an open subcover of $f(A)$. \square

As a corollary of Proposition 26.27, if E is compact, F is Hausdorff, and $f: E \rightarrow F$ is continuous and bijective, then f is a homeomorphism. Indeed, it is enough to show that f^{-1} is continuous, which is equivalent to showing that f maps closed sets to closed sets. However, closed sets are compact and Proposition 26.27 shows that compact sets are mapped to compact sets, which, by Proposition 26.23, are closed.

It can also be shown that if E is a compact nonempty space and $f: E \rightarrow \mathbb{R}$ is a continuous function, then there are points $a, b \in E$ such that $f(a)$ is the minimum of $f(E)$ and $f(b)$ is the maximum of $f(E)$. Indeed, $f(E)$ is a compact subset of \mathbb{R} and thus, a closed and bounded set which contains its greatest lower bound and its least upper bound.

Another useful notion is that of local compactness. Indeed, manifolds and surfaces are locally compact.

Definition 26.21. A topological space, E , is *locally compact* if it is Hausdorff and for every $a \in E$, there is some compact neighborhood, K , of a .

From Proposition 26.26, every compact space is locally compact but the converse is false. It can be shown that a normed vector space of finite dimension is locally compact.

Proposition 26.28. *Given a locally compact topological space, E , for every $a \in E$, for every neighborhood, N , of a , there exists a compact neighborhood, U , of a , such that $U \subseteq N$.*

Proof. For any $a \in E$, there is some compact neighborhood, V , of a . By Proposition 26.26, every neighborhood of a relative to V contains some compact neighborhood U of a relative to V . But every neighborhood of a relative to V is a neighborhood of a relative to E and every neighborhood N of a in E yields a neighborhood, $V \cap N$, of a in V and thus, for every neighborhood, N , of a , there exists a compact neighborhood, U , of a such that $U \subseteq N$. \square

It is much harder to deal with noncompact surfaces (or manifolds) than it is to deal with compact surfaces (or manifolds). However, surfaces (and manifolds) are locally compact and it turns out that there are various ways of embedding a locally compact Hausdorff space into a compact Hausdorff space. The most economical construction consists in adding just one point. This construction, known as the *Alexandroff compactification*, is technically useful, and we now describe it and sketch the proof that it achieves its goal.

To help the reader's intuition, let us consider the case of the plane, \mathbb{R}^2 . If we view the plane, \mathbb{R}^2 , as embedded in 3-space, \mathbb{R}^3 , say as the xOy plane of equation $z = 0$, we can consider the sphere, Σ , of radius 1 centered on the z -axis at the point $(0, 0, 1)$ and tangent to the xOy plane at the origin (sphere of equation $x^2 + y^2 + (z - 1)^2 = 1$). If N denotes the north pole on the sphere, i.e., the point of coordinates $(0, 0, 2)$, then any line, D , passing through the north pole and not tangent to the sphere (i.e., not parallel to the xOy plane) intersects the xOy plane in a unique point, M , and the sphere in a unique point, P , other than the north pole, N . This, way, we obtain a bijection between the xOy plane and the punctured sphere Σ , i.e., the sphere with the north pole N deleted. This bijection is called a *stereographic projection*. The Alexandroff compactification of the plane puts the north pole back on the sphere, which amounts to adding a single point at infinity ∞ to the plane. Intuitively, as we travel away from the origin O towards infinity (in any direction!), we tend towards an ideal point at infinity ∞ . Imagine that we "bend" the plane so that it gets wrapped around the sphere, according to stereographic projection. A simpler example takes a line and gets a circle as its compactification. The Alexandroff compactification is a generalization of these simple constructions.

Definition 26.22. Let (E, \mathcal{O}) be a locally compact space. Let ω be any point not in E , and let $E_\omega = E \cup \{\omega\}$. Define the family, \mathcal{O}_ω , as follows:

$$\mathcal{O}_\omega = \mathcal{O} \cup \{(E - K) \cup \{\omega\} \mid K \text{ compact in } E\}.$$

The pair, $(E_\omega, \mathcal{O}_\omega)$, is called the *Alexandroff compactification (or one point compactification)* of (E, \mathcal{O}) .

The following theorem shows that $(E_\omega, \mathcal{O}_\omega)$ is indeed a topological space, and that it is compact.

Theorem 26.29. *Let E be a locally compact topological space. The Alexandroff compactification, E_ω , of E is a compact space such that E is a subspace of E_ω and if E is not compact, then $\bar{E} = E_\omega$.*

Proof. The verification that \mathcal{O}_ω is a family of open sets is not difficult but a bit tedious. Details can be found in Munkres [84] or Schwartz [93]. Let us show that E_ω is compact. For every open cover, $(U_i)_{i \in I}$, of E_ω , since ω must be covered, there is some U_{i_0} of the form

$$U_{i_0} = (E - K_0) \cup \{\omega\}$$

where K_0 is compact in E . Consider the family, $(V_i)_{i \in I}$, defined as follows:

$$\begin{aligned} V_i &= U_i & \text{if } U_i \in \mathcal{O}, \\ V_i &= E - K & \text{if } U_i = (E - K) \cup \{\omega\}, \end{aligned}$$

where K is compact in E . Then, because each K is compact and thus closed in E (since E is Hausdorff), $E - K$ is open, and every V_i is an open subset of E . Furthermore, the family, $(V_i)_{i \in I - \{i_0\}}$, is an open cover of K_0 . Since K_0 is compact, there is a finite open subcover, $(V_j)_{j \in J}$, of K_0 , and thus, $(U_j)_{j \in J \cup \{i_0\}}$ is a finite open cover of E_ω .

Let us show that E_ω is Hausdorff. Given any two points, $a, b \in E_\omega$, if both $a, b \in E$, since E is Hausdorff and every open set in \mathcal{O} is an open set in \mathcal{O}_ω , there exist disjoint open sets, U, V (in \mathcal{O}), such that $a \in U$ and $b \in V$. If $b = \omega$, since E is locally compact, there is some compact set, K , containing an open set, U , containing a and then, U and $V = (E - K) \cup \{\omega\}$ are disjoint open sets (in \mathcal{O}_ω) such that $a \in U$ and $b \in V$.

The space E is a subspace of E_ω because for every open set, U , in \mathcal{O}_ω , either $U \in \mathcal{O}$ and $E \cap U = U$ is open in E , or $U = (E - K) \cup \{\omega\}$, where K is compact in E , and thus, $U \cap E = E - K$, which is open in E , since K is compact in E and thus, closed (since E is Hausdorff). Finally, if E is not compact, for every compact subset, K , of E , $E - K$ is nonempty and thus, for every open set, $U = (E - K) \cup \{\omega\}$, containing ω , we have $U \cap E \neq \emptyset$, which shows that $\omega \in \bar{E}$ and thus, that $\bar{E} = E_\omega$. \square

Finally, in studying surfaces and manifolds, an important property is the existence of a countable basis for the topology. Indeed, this property guarantees the existence of triangulations of surfaces, a crucial property.

Definition 26.23. A topological space E is called *second-countable* if there is a countable basis for its topology, i.e., if there is a countable family, $(U_i)_{i \geq 0}$, of open sets such that every open set of E is a union of open sets U_i .

It is easily seen that \mathbb{R}^n is second-countable and more generally, that every normed vector space of finite dimension is second-countable. It can also be shown that if E is a locally compact space that has a countable basis, then E_ω also has a countable basis (and in fact, is metrizable). We have the following properties.

Proposition 26.30. *Given a second-countable topological space E , every open cover $(U_i)_{i \in I}$, of E contains some countable subcover.*

Proof. Let $(O_n)_{n \geq 0}$ be a countable basis for the topology. Then, all sets O_n contained in some U_i can be arranged into a countable subsequence, $(\Omega_m)_{m \geq 0}$, of $(O_n)_{n \geq 0}$ and for every Ω_m , there is some U_{i_m} such that $\Omega_m \subseteq U_{i_m}$. Furthermore, every U_i is some union of sets Ω_j , and thus, every $a \in E$ belongs to some Ω_j , which shows that $(\Omega_m)_{m \geq 0}$ is a countable open subcover of $(U_i)_{i \in I}$. \square

As an immediate corollary of Proposition 26.30, a locally connected second-countable space has countably many connected components.

In second-countable Hausdorff spaces, compactness can be characterized in terms of accumulation points (this is also true for metric spaces).

Definition 26.24. Given a topological Hausdorff space, E , given any sequence, (x_n) , of points in E , a point, $l \in E$, is an *accumulation point (or cluster point)* of the sequence (x_n) if every open set, U , containing l contains x_n for infinitely many n .

Clearly, if l is a limit of the sequence, (x_n) , then it is an accumulation point, since every open set, U , containing a contains all x_n except for finitely many n .

Proposition 26.31. *A second-countable topological Hausdorff space, E , is compact iff every sequence, (x_n) , has some accumulation point.*

Proof. Assume that every sequence, (x_n) , has some accumulation point. Let $(U_i)_{i \in I}$ be some open cover of E . By Proposition 26.30, there is a countable open subcover, $(O_n)_{n \geq 0}$, for E . Now, if E is not covered by any finite subcover of $(O_n)_{n \geq 0}$, we can define a sequence, (x_m) , by induction as follows:

Let x_0 be arbitrary and for every $m \geq 1$, let x_m be some point in E not in $O_1 \cup \cdots \cup O_m$, which exists, since $O_1 \cup \cdots \cup O_m$ is not an open cover of E . We claim that the sequence, (x_m) , does not have any accumulation point. Indeed, for every $l \in E$, since $(O_n)_{n \geq 0}$ is an open cover of E , there is some O_m such that $l \in O_m$, and by construction, every x_n with $n \geq m + 1$ does not belong to O_m , which means that $x_n \in O_m$ for only finitely many n and l is not an accumulation point.

Conversely, assume that E is compact, and let (x_n) be any sequence. If $l \in E$ is not an accumulation point of the sequence, then there is some open set, U_l , such that $l \in U_l$ and $x_n \in U_l$ for only finitely many n . Thus, if (x_n) does not have any accumulation point, the family, $(U_l)_{l \in E}$, is an open cover of E and since E is compact, it has some finite open subcover, $(U_l)_{l \in J}$, where J is a finite subset of E . But every U_l with $l \in J$ is such that $x_n \in U_l$ for only finitely many n , and since J is finite, $x_n \in \bigcup_{l \in J} U_l$ for only finitely many n , which contradicts the fact that $(U_l)_{l \in J}$ is an open cover of E , and thus contains all the x_n . Thus, (x_n) has some accumulation point. \square

Remark: It should be noted that the proof showing that if E is compact, then every sequence has some accumulation point, holds for any arbitrary compact space (the proof does not use a countable basis for the topology). The converse also holds for metric spaces. We will prove this converse since it is a major property of metric spaces.

Given a metric space in which every sequence has some accumulation point, we first prove the existence of a *Lebesgue number*.

Lemma 26.32. *Given a metric space, E , if every sequence, (x_n) , has an accumulation point, for every open cover, $(U_i)_{i \in I}$, of E , there is some $\delta > 0$ (a Lebesgue number for $(U_i)_{i \in I}$) such that, for every open ball, $B_0(a, \epsilon)$, of radius $\epsilon \leq \delta$, there is some open subset, U_i , such that $B_0(a, \epsilon) \subseteq U_i$.*

Proof. If there was no δ with the above property, then, for every natural number, n , there would be some open ball, $B_0(a_n, 1/n)$, which is not contained in any open set, U_i , of the open cover, $(U_i)_{i \in I}$. However, the sequence, (a_n) , has some accumulation point, a , and since $(U_i)_{i \in I}$ is an open cover of E , there is some U_i such that $a \in U_i$. Since U_i is open, there is some open ball of center a and radius ϵ contained in U_i . Now, since a is an accumulation point of the sequence, (a_n) , every open set containing a contains a_n for infinitely many n and thus, there is some n large enough so that

$$1/n \leq \epsilon/2 \quad \text{and} \quad a_n \in B_0(a, \epsilon/2),$$

which implies that

$$B_0(a_n, 1/n) \subseteq B_0(a, \epsilon) \subseteq U_i,$$

a contradiction. \square

By a previous remark, since the proof of Proposition 26.31 implies that in a compact topological space, every sequence has some accumulation point, by Lemma 26.32, in a compact metric space, every open cover has a Lebesgue number. This fact can be used to prove another important property of compact metric spaces, the uniform continuity theorem.

Definition 26.25. Given two metric spaces, (E, d_E) and (F, d_F) , a function, $f: E \rightarrow F$, is *uniformly continuous* if for every $\epsilon > 0$, there is some $\eta > 0$, such that, for all $a, b \in E$,

$$\text{if } d_E(a, b) \leq \eta \text{ then } d_F(f(a), f(b)) \leq \epsilon.$$

The *uniform continuity theorem* can be stated as follows:

Theorem 26.33. *Given two metric spaces, (E, d_E) and (F, d_F) , if E is compact and $f: E \rightarrow F$ is a continuous function, then it is uniformly continuous.*

Proof. Consider any $\epsilon > 0$ and let $(B_0(y, \epsilon/2))_{y \in F}$ be the open cover of F consisting of open balls of radius $\epsilon/2$. Since f is continuous, the family,

$$(f^{-1}(B_0(y, \epsilon/2)))_{y \in F},$$

is an open cover of E . Since, E is compact, by Lemma 26.32, there is a Lebesgue number, δ , such that for every open ball, $B_0(a, \eta)$, of radius $\eta \leq \delta$, then $B_0(a, \eta) \subseteq f^{-1}(B_0(y, \epsilon/2))$, for some $y \in F$. In particular, for any $a, b \in E$ such that $d_E(a, b) \leq \eta = \delta/2$, we have $a, b \in B_0(a, \delta)$ and thus, $a, b \in f^{-1}(B_0(y, \epsilon/2))$, which implies that $f(a), f(b) \in B_0(y, \epsilon/2)$. But then, $d_F(f(a), f(b)) \leq \epsilon$, as desired. \square

We now prove another lemma needed to obtain the characterization of compactness in metric spaces in terms of accumulation points.

Lemma 26.34. *Given a metric space, E , if every sequence, (x_n) , has an accumulation point, then for every $\epsilon > 0$, there is a finite open cover, $B_0(a_0, \epsilon) \cup \cdots \cup B_0(a_n, \epsilon)$, of E by open balls of radius ϵ .*

Proof. Let a_0 be any point in E . If $B_0(a_0, \epsilon) = E$, then the lemma is proved. Otherwise, assume that a sequence, (a_0, a_1, \dots, a_n) , has been defined, such that $B_0(a_0, \epsilon) \cup \cdots \cup B_0(a_n, \epsilon)$ does not cover E . Then, there is some a_{n+1} not in $B_0(a_0, \epsilon) \cup \cdots \cup B_0(a_n, \epsilon)$ and either

$$B_0(a_0, \epsilon) \cup \cdots \cup B_0(a_{n+1}, \epsilon) = E,$$

in which case the lemma is proved, or we obtain a sequence, $(a_0, a_1, \dots, a_{n+1})$, such that $B_0(a_0, \epsilon) \cup \cdots \cup B_0(a_{n+1}, \epsilon)$ does not cover E . If this process goes on forever, we obtain an infinite sequence, (a_n) , such that $d(a_m, a_n) > \epsilon$ for all $m \neq n$. Since every sequence in E has some accumulation point, the sequence, (a_n) , has some accumulation point, a . Then, for infinitely many n , we must have $d(a_n, a) \leq \epsilon/3$ and thus, for at least two distinct natural numbers, p, q , we must have $d(a_p, a) \leq \epsilon/3$ and $d(a_q, a) \leq \epsilon/3$, which implies $d(a_p, a_q) \leq 2\epsilon/3$, contradicting the fact that $d(a_m, a_n) > \epsilon$ for all $m \neq n$. Thus, there must be some n such that

$$B_0(a_0, \epsilon) \cup \cdots \cup B_0(a_n, \epsilon) = E.$$

\square

A metric space satisfying the condition of Lemma 26.34 is sometimes called *precompact* (or *totally bounded*). We now obtain the *Weierstrass–Bolzano* property.

Theorem 26.35. *A metric space, E , is compact iff every sequence, (x_n) , has an accumulation point.*

Proof. We already observed that the proof of Proposition 26.31 shows that for any compact space (not necessarily metric), every sequence, (x_n) , has an accumulation point. Conversely, let E be a metric space, and assume that every sequence, (x_n) , has an accumulation point. Given any open cover, $(U_i)_{i \in I}$, for E , we must find a finite open subcover of E . By Lemma 26.32, there is some $\delta > 0$ (a Lebesgue number for $(U_i)_{i \in I}$) such that, for every open ball, $B_0(a, \epsilon)$, of radius $\epsilon \leq \delta$, there is some open subset, U_j , such that $B_0(a, \epsilon) \subseteq U_j$. By Lemma 26.34, for every $\delta > 0$, there is a finite open cover, $B_0(a_0, \delta) \cup \cdots \cup B_0(a_n, \delta)$, of E by open balls of radius δ . But from the previous statement, every open ball, $B_0(a_i, \delta)$, is contained in some open set, U_{j_i} , and thus, $\{U_{j_1}, \dots, U_{j_n}\}$ is an open cover of E . \square

Another very useful characterization of compact metric spaces is obtained in terms of Cauchy sequences. Such a characterization is quite useful in fractal geometry (and elsewhere). First, recall the definition of a Cauchy sequence and of a complete metric space.

Definition 26.26. Given a metric space, (E, d) , a sequence, $(x_n)_{n \in \mathbb{N}}$, in E is a *Cauchy sequence* if the following condition holds: for every $\epsilon > 0$, there is some $p \geq 0$, such that, for all $m, n \geq p$, then $d(x_m, x_n) \leq \epsilon$.

If every Cauchy sequence in (E, d) converges we say that (E, d) is a *complete metric space*.

First, let us show the following proposition:

Proposition 26.36. *Given a metric space, E , if a Cauchy sequence, (x_n) , has some accumulation point, a , then a is the limit of the sequence, (x_n) .*

Proof. Since (x_n) is a Cauchy sequence, for every $\epsilon > 0$, there is some $p \geq 0$, such that, for all $m, n \geq p$, then $d(x_m, x_n) \leq \epsilon/2$. Since a is an accumulation point for (x_n) , for infinitely many n , we have $d(x_n, a) \leq \epsilon/2$, and thus, for at least some $n \geq p$, we have $d(x_n, a) \leq \epsilon/2$. Then, for all $m \geq p$,

$$d(x_m, a) \leq d(x_m, x_n) + d(x_n, a) \leq \epsilon,$$

which shows that a is the limit of the sequence (x_n) . \square

Recall that a metric space is *precompact* (or *totally bounded*) if for every $\epsilon > 0$, there is a finite open cover, $B_0(a_0, \epsilon) \cup \cdots \cup B_0(a_n, \epsilon)$, of E by open balls of radius ϵ . We can now prove the following theorem.

Theorem 26.37. *A metric space, E , is compact iff it is precompact and complete.*

Proof. Let E be compact. For every $\epsilon > 0$, the family of all open balls of radius ϵ is an open cover for E and since E is compact, there is a finite subcover, $B_0(a_0, \epsilon) \cup \cdots \cup B_0(a_n, \epsilon)$, of E by open balls of radius ϵ . Thus, E is precompact. Since E is compact, by Theorem 26.35, every sequence, (x_n) , has some accumulation point. Thus, every Cauchy sequence, (x_n) , has some accumulation point, a , and, by Proposition 26.36, a is the limit of (x_n) . Thus, E is complete.

Now, assume that E is precompact and complete. We prove that every sequence, (x_n) , has an accumulation point. By the other direction of Theorem 26.35, this shows that E is compact. Given any sequence, (x_n) , we construct a Cauchy subsequence, (y_n) , of (x_n) as follows: Since E is precompact, letting $\epsilon = 1$, there exists a finite cover, \mathcal{U}_1 , of E by open balls of radius 1. Thus, some open ball, B_o^1 , in the cover, \mathcal{U}_1 , contains infinitely many elements from the sequence (x_n) . Let y_0 be any element of (x_n) in B_o^1 . By induction, assume that a sequence of open balls, $(B_o^i)_{1 \leq i \leq m}$, has been defined, such that every ball, B_o^i , has radius $\frac{1}{2^i}$, contains infinitely many elements from the sequence (x_n) and contains some y_i from (x_n) such that

$$d(y_i, y_{i+1}) \leq \frac{1}{2^i},$$

for all i , $0 \leq i \leq m-1$. Then, letting $\epsilon = \frac{1}{2^{m+1}}$, because E is precompact, there is some finite cover, \mathcal{U}_{m+1} , of E by open balls of radius ϵ and thus, of the open ball B_o^m . Thus, some open ball, B_o^{m+1} , in the cover, \mathcal{U}_{m+1} , contains infinitely many elements from the sequence, (x_n) , and we let y_{m+1} be any element of (x_n) in B_o^{m+1} . Thus, we have defined by induction a sequence, (y_n) , which is a subsequence of, (x_n) , and such that

$$d(y_i, y_{i+1}) \leq \frac{1}{2^i},$$

for all i . However, for all $m, n \geq 1$, we have

$$d(y_m, y_n) \leq d(y_m, y_{m+1}) + \cdots + d(y_{n-1}, y_n) \leq \sum_{i=m}^n \frac{1}{2^i} \leq \frac{1}{2^{m-1}},$$

and thus, (y_n) is a Cauchy sequence. Since E is complete, the sequence, (y_n) , has a limit, and since it is a subsequence of (x_n) , the sequence, (x_n) , has some accumulation point. \square

If (E, d) is a nonempty complete metric space, every map, $f: E \rightarrow E$, for which there is some k such that $0 \leq k < 1$ and

$$d(f(x), f(y)) \leq kd(x, y)$$

for all $x, y \in E$, has the very important property that it has a unique fixed point, that is, there is a unique, $a \in E$, such that $f(a) = a$. A map as above is called a *contraction mapping*. Furthermore, the fixed point of a contraction mapping can be computed as the limit of a fast converging sequence.

The fixed point property of contraction mappings is used to show some important theorems of analysis, such as the implicit function theorem and the existence of solutions to certain differential equations. It can also be used to show the existence of fractal sets defined in terms of iterated function systems. Since the proof is quite simple, we prove the fixed point property of contraction mappings. First, observe that a contraction mapping is (uniformly) continuous.

Proposition 26.38. *If (E, d) is a nonempty complete metric space, every contraction mapping, $f: E \rightarrow E$, has a unique fixed point. Furthermore, for every $x_0 \in E$, defining the sequence, (x_n) , such that $x_{n+1} = f(x_n)$, the sequence, (x_n) , converges to the unique fixed point of f .*

Proof. First, we prove that f has at most one fixed point. Indeed, if $f(a) = a$ and $f(b) = b$, since

$$d(a, b) = d(f(a), f(b)) \leq kd(a, b)$$

and $0 \leq k < 1$, we must have $d(a, b) = 0$, that is, $a = b$.

Next, we prove that (x_n) is a Cauchy sequence. Observe that

$$\begin{aligned} d(x_2, x_1) &\leq kd(x_1, x_0), \\ d(x_3, x_2) &\leq kd(x_2, x_1) \leq k^2d(x_1, x_0), \\ &\vdots \\ d(x_{n+1}, x_n) &\leq kd(x_n, x_{n-1}) \leq \cdots \leq k^nd(x_1, x_0). \end{aligned}$$

Thus, we have

$$\begin{aligned} d(x_{n+p}, x_n) &\leq d(x_{n+p}, x_{n+p-1}) + d(x_{n+p-1}, x_{n+p-2}) + \cdots + d(x_{n+1}, x_n) \\ &\leq (k^{p-1} + k^{p-2} + \cdots + k + 1)k^nd(x_1, x_0) \\ &\leq \frac{k^n}{1-k} d(x_1, x_0). \end{aligned}$$

We conclude that $d(x_{n+p}, x_n)$ converges to 0 when n goes to infinity, which shows that (x_n) is a Cauchy sequence. Since E is complete, the sequence (x_n) has a limit, a . Since f is continuous, the sequence $(f(x_n))$ converges to $f(a)$. But $x_{n+1} = f(x_n)$ converges to a and so $f(a) = a$, the unique fixed point of f . \square

Note that no matter how the starting point x_0 of the sequence (x_n) is chosen, (x_n) converges to the unique fixed point of f . Also, the convergence is fast, since

$$d(x_n, a) \leq \frac{k^n}{1-k} d(x_1, x_0).$$

The Hausdorff distance between compact subsets of a metric space provides a very nice illustration of some of the theorems on complete and compact metric spaces just presented.

Definition 26.27. Given a metric space, (X, d) , for any subset, $A \subseteq X$, for any, $\epsilon \geq 0$, define the ϵ -hull of A as the set

$$V_\epsilon(A) = \{x \in X, \exists a \in A \mid d(a, x) \leq \epsilon\}.$$

Given any two nonempty bounded subsets, A, B of X , define $D(A, B)$, the Hausdorff distance between A and B , by

$$D(A, B) = \inf\{\epsilon \geq 0 \mid A \subseteq V_\epsilon(B) \text{ and } B \subseteq V_\epsilon(A)\}.$$

Note that since we are considering nonempty bounded subsets, $D(A, B)$ is well defined (i.e., not infinite). However, D is not necessarily a distance function. It is a distance function if we restrict our attention to nonempty compact subsets of X (actually, it is also a metric on closed and bounded subsets). We let $\mathcal{K}(X)$ denote the set of all nonempty compact subsets of X . The remarkable fact is that D is a distance on $\mathcal{K}(X)$ and that if X is complete or compact, then so is $\mathcal{K}(X)$. The following theorem is taken from Edgar [33].

Theorem 26.39. *If (X, d) is a metric space, then the Hausdorff distance, D , on the set, $\mathcal{K}(X)$, of nonempty compact subsets of X is a distance. If (X, d) is complete, then $(\mathcal{K}(X), D)$ is complete and if (X, d) is compact, then $(\mathcal{K}(X), D)$ is compact.*

Proof. Since (nonempty) compact sets are bounded, $D(A, B)$ is well defined. Clearly, D is symmetric. Assume that $D(A, B) = 0$. Then, for every $\epsilon > 0$, $A \subseteq V_\epsilon(B)$, which means that for every $a \in A$, there is some $b \in B$ such that $d(a, b) \leq \epsilon$, and thus, that $A \subseteq \overline{B}$. Since B is closed, $\overline{B} = B$, and we have $A \subseteq B$. Similarly, $B \subseteq A$, and thus, $A = B$. Clearly, if $A = B$, we have $D(A, B) = 0$. It remains to prove the triangle inequality. If $B \subseteq V_{\epsilon_1}(A)$ and $C \subseteq V_{\epsilon_2}(B)$, then

$$V_{\epsilon_2}(B) \subseteq V_{\epsilon_2}(V_{\epsilon_1}(A)),$$

and since

$$V_{\epsilon_2}(V_{\epsilon_1}(A)) \subseteq V_{\epsilon_1 + \epsilon_2}(A),$$

we get

$$C \subseteq V_{\epsilon_2}(B) \subseteq V_{\epsilon_1 + \epsilon_2}(A).$$

Similarly, we can prove that

$$A \subseteq V_{\epsilon_1 + \epsilon_2}(C),$$

and thus, the triangle inequality follows.

Next, we need to prove that if (X, d) is complete, then $(\mathcal{K}(X), D)$ is also complete. First, we show that if (A_n) is a sequence of nonempty compact sets converging to a nonempty compact set A in the Hausdorff metric, then

$$A = \{x \in X \mid \text{there is a sequence, } (x_n), \text{ with } x_n \in A_n \text{ converging to } x\}.$$

Indeed, if (x_n) is a sequence with $x_n \in A_n$ converging to x and (A_n) converges to A then, for every $\epsilon > 0$, there is some x_n such that $d(x_n, x) \leq \epsilon/2$ and there is some $a_n \in A$ such that

$d(a_n, x_n) \leq \epsilon/2$ and thus, $d(a_n, x) \leq \epsilon$, which shows that $x \in \overline{A}$. Since A is compact, it is closed, and $x \in A$. Conversely, since (A_n) converges to A , for every $x \in A$, for every $n \geq 1$, there is some $x_n \in A_n$ such that $d(x_n, x) \leq 1/n$ and the sequence (x_n) converges to x .

Now, let (A_n) be a Cauchy sequence in $\mathcal{K}(X)$. It can be proven that (A_n) converges to the set

$$A = \{x \in X \mid \text{there is a sequence, } (x_n), \text{ with } x_n \in A_n \text{ converging to } x\},$$

and that A is nonempty and compact. To prove that A is compact, one proves that it is totally bounded and complete. Details are given in Edgar [33].

Finally, we need to prove that if (X, d) is compact, then $(\mathcal{K}(X), D)$ is compact. Since we already know that $(\mathcal{K}(X), D)$ is complete if (X, d) is, it is enough to prove that $(\mathcal{K}(X), D)$ is totally bounded if (X, d) is, which is not hard. \square

In view of Theorem 26.39 and Theorem 26.38, it is possible to define some nonempty compact subsets of X in terms of fixed points of contraction maps. This can be done in terms of iterated function systems, yielding a large class of fractals. However, we will omit this topic and instead refer the reader to Edgar [33].

Finally, returning to second-countable spaces, we give another characterization of accumulation points.

Proposition 26.40. *Given a second-countable topological Hausdorff space, E , a point, l , is an accumulation point of the sequence, (x_n) , iff l is the limit of some subsequence, (x_{n_k}) , of (x_n) .*

Proof. Clearly, if l is the limit of some subsequence (x_{n_k}) of (x_n) , it is an accumulation point of (x_n) .

Conversely, let $(U_k)_{k \geq 0}$ be the sequence of open sets containing l , where each U_k belongs to a countable basis of E , and let $V_k = U_1 \cap \cdots \cap U_k$. For every $k \geq 1$, we can find some $n_k > n_{k-1}$ such that $x_{n_k} \in V_k$, since l is an accumulation point of (x_n) . Now, since every open set containing l contains some U_{k_0} and since $x_{n_k} \in U_{k_0}$ for all $k \geq 0$, the sequence (x_{n_k}) has limit l . \square

Remark: Proposition 26.40 also holds for metric spaces.

In Chapter 27 we show how certain fractals can be defined by iterated function systems, using Theorem 26.39 and Theorem 26.38.

Before considering differentials, we need to look at the continuity of linear maps.

26.6 Continuous Linear and Multilinear Maps

If E and F are normed vector spaces, we first characterize when a linear map $f: E \rightarrow F$ is continuous.

Proposition 26.41. *Given two normed vector spaces E and F , for any linear map $f: E \rightarrow F$, the following conditions are equivalent:*

- (1) *The function f is continuous at 0.*
- (2) *There is a constant $k \geq 0$ such that,*

$$\|f(u)\| \leq k, \text{ for every } u \in E \text{ such that } \|u\| \leq 1.$$

- (3) *There is a constant $k \geq 0$ such that,*

$$\|f(u)\| \leq k\|u\|, \text{ for every } u \in E.$$

- (4) *The function f is continuous at every point of E .*

Proof. Assume (1). Then, for every $\epsilon > 0$, there is some $\eta > 0$ such that, for every $u \in E$, if $\|u\| \leq \eta$, then $\|f(u)\| \leq \epsilon$. Pick $\epsilon = 1$, so that there is some $\eta > 0$ such that, if $\|u\| \leq \eta$, then $\|f(u)\| \leq 1$. If $\|u\| \leq 1$, then $\|\eta u\| \leq \eta\|u\| \leq \eta$, and so, $\|f(\eta u)\| \leq 1$, that is, $\eta\|f(u)\| \leq 1$, which implies $\|f(u)\| \leq \eta^{-1}$. Thus, (2) holds with $k = \eta^{-1}$.

Assume that (2) holds. If $u = 0$, then by linearity, $f(0) = 0$, and thus $\|f(0)\| \leq k\|0\|$ holds trivially for all $k \geq 0$. If $u \neq 0$, then $\|u\| > 0$, and since

$$\left\| \frac{u}{\|u\|} \right\| = 1,$$

we have

$$\left\| f\left(\frac{u}{\|u\|}\right) \right\| \leq k,$$

which implies that

$$\|f(u)\| \leq k\|u\|.$$

Thus, (3) holds.

If (3) holds, then for all $u, v \in E$, we have

$$\|f(v) - f(u)\| = \|f(v - u)\| \leq k\|v - u\|.$$

If $k = 0$, then f is the zero function, and continuity is obvious. Otherwise, if $k > 0$, for every $\epsilon > 0$, if $\|v - u\| \leq \frac{\epsilon}{k}$, then $\|f(v - u)\| \leq \epsilon$, which shows continuity at every $u \in E$. Finally, it is obvious that (4) implies (1). \square

Among other things, Proposition 26.41 shows that a linear map is continuous iff the image of the unit (closed) ball is bounded. If E and F are normed vector spaces, the set of all continuous linear maps $f: E \rightarrow F$ is denoted by $\mathcal{L}(E; F)$.

Using Proposition 26.41, we can define a norm on $\mathcal{L}(E; F)$ which makes it into a normed vector space. This definition has already been given in Chapter 7 (Definition 7.7) but for the reader's convenience, we repeat it here.

Definition 26.28. Given two normed vector spaces E and F , for every continuous linear map $f: E \rightarrow F$, we define the *norm* $\|f\|$ of f as

$$\|f\| = \min \{k \geq 0 \mid \|f(x)\| \leq k\|x\|, \text{ for all } x \in E\} = \max \{\|f(x)\| \mid \|x\| \leq 1\}.$$

From Definition 26.28, for every continuous linear map $f \in \mathcal{L}(E; F)$, we have

$$\|f(x)\| \leq \|f\|\|x\|,$$

for every $x \in E$. It is easy to verify that $\mathcal{L}(E; F)$ is a normed vector space under the norm of Definition 26.28. Furthermore, if E, F, G , are normed vector spaces, and $f: E \rightarrow F$ and $g: F \rightarrow G$ are continuous linear maps, we have

$$\|g \circ f\| \leq \|g\|\|f\|.$$

We can now show that when $E = \mathbb{R}^n$ or $E = \mathbb{C}^n$, with any of the norms $\|\cdot\|_1$, $\|\cdot\|_2$, or $\|\cdot\|_\infty$, then every linear map $f: E \rightarrow F$ is continuous.

Proposition 26.42. *If $E = \mathbb{R}^n$ or $E = \mathbb{C}^n$, with any of the norms $\|\cdot\|_1$, $\|\cdot\|_2$, or $\|\cdot\|_\infty$, and F is any normed vector space, then every linear map $f: E \rightarrow F$ is continuous.*

Proof. Let (e_1, \dots, e_n) be the standard basis of \mathbb{R}^n (a similar proof applies to \mathbb{C}^n). In view of Proposition 7.2, it is enough to prove the proposition for the norm

$$\|x\|_\infty = \max\{|x_i| \mid 1 \leq i \leq n\}.$$

We have,

$$\|f(v) - f(u)\| = \|f(v - u)\| = \left\| f\left(\sum_{1 \leq i \leq n} (v_i - u_i)e_i\right) \right\| = \left\| \sum_{1 \leq i \leq n} (v_i - u_i)f(e_i) \right\|,$$

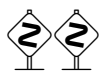
and so,

$$\|f(v) - f(u)\| \leq \left(\sum_{1 \leq i \leq n} \|f(e_i)\| \right) \max_{1 \leq i \leq n} |v_i - u_i| = \left(\sum_{1 \leq i \leq n} \|f(e_i)\| \right) \|v - u\|_\infty.$$

By the argument used in Proposition 26.41 to prove that (3) implies (4), f is continuous. \square

Actually, we proved in Theorem 7.3 that if E is a vector space of finite dimension, then any two norms are equivalent, so that they define the same topology. This fact together with Proposition 26.42 prove the following:

Theorem 26.43. *If E is a vector space of finite dimension (over \mathbb{R} or \mathbb{C}), then all norms are equivalent (define the same topology). Furthermore, for any normed vector space F , every linear map $f: E \rightarrow F$ is continuous.*



If E is a normed vector space of infinite dimension, a linear map $f: E \rightarrow F$ may not be continuous. As an example, let E be the infinite vector space of all polynomials over \mathbb{R} . Let

$$\|P(X)\| = \max_{0 \leq x \leq 1} |P(x)|.$$

We leave as an exercise to show that this is indeed a norm. Let $F = \mathbb{R}$, and let $f: E \rightarrow F$ be the map defined such that, $f(P(X)) = P(3)$. It is clear that f is linear. Consider the sequence of polynomials

$$P_n(X) = \left(\frac{X}{2}\right)^n.$$

It is clear that $\|P_n\| = \left(\frac{1}{2}\right)^n$, and thus, the sequence P_n has the null polynomial as a limit. However, we have

$$f(P_n(X)) = P_n(3) = \left(\frac{3}{2}\right)^n,$$

and the sequence $f(P_n(X))$ diverges to $+\infty$. Consequently, in view of Proposition 26.12 (1), f is not continuous.

We now consider the continuity of multilinear maps. We treat explicitly bilinear maps, the general case being a straightforward extension.

Proposition 26.44. *Given normed vector spaces E , F and G , for any bilinear map $f: E \times E \rightarrow G$, the following conditions are equivalent:*

(1) *The function f is continuous at $\langle 0, 0 \rangle$.*

(2) *There is a constant $k \geq 0$ such that,*

$$\|f(u, v)\| \leq k, \text{ for all } u, v \in E \text{ such that } \|u\|, \|v\| \leq 1.$$

(3) *There is a constant $k \geq 0$ such that,*

$$\|f(u, v)\| \leq k\|u\|\|v\|, \text{ for all } u, v \in E.$$

(4) *The function f is continuous at every point of $E \times F$.*

Proof. It is similar to that of Proposition 26.41, with a small subtlety in proving that (3) implies (4), namely that two different η 's that are not independent are needed. \square

If E , F , and G , are normed vector spaces, we denote the set of all continuous bilinear maps $f: E \times F \rightarrow G$ by $\mathcal{L}_2(E, F; G)$. Using Proposition 26.44, we can define a norm on $\mathcal{L}_2(E, F; G)$ which makes it into a normed vector space.

Definition 26.29. Given normed vector spaces E , F , and G , for every continuous bilinear map $f: E \times F \rightarrow G$, we define the *norm* $\|f\|$ of f as

$$\begin{aligned}\|f\| &= \min \{k \geq 0 \mid \|f(x, y)\| \leq k\|x\|\|y\|, \text{ for all } x, y \in E\} \\ &= \max \{\|f(x, y)\| \mid \|x\|, \|y\| \leq 1\}.\end{aligned}$$

From Definition 26.28, for every continuous bilinear map $f \in \mathcal{L}_2(E, F; G)$, we have

$$\|f(x, y)\| \leq \|f\|\|x\|\|y\|,$$

for all $x, y \in E$. It is easy to verify that $\mathcal{L}_2(E, F; G)$ is a normed vector space under the norm of Definition 26.29.

Given a bilinear map $f: E \times F \rightarrow G$, for every $u \in E$, we obtain a linear map denoted $fu: F \rightarrow G$, defined such that, $fu(v) = f(u, v)$. Furthermore, since

$$\|f(x, y)\| \leq \|f\|\|x\|\|y\|,$$

it is clear that fu is continuous. We can then consider the map $\varphi: E \rightarrow \mathcal{L}(F; G)$, defined such that, $\varphi(u) = fu$, for any $u \in E$, or equivalently, such that,

$$\varphi(u)(v) = f(u, v).$$

Actually, it is easy to show that φ is linear and continuous, and that $\|\varphi\| = \|f\|$. Thus, $f \mapsto \varphi$ defines a map from $\mathcal{L}_2(E, F; G)$ to $\mathcal{L}(E; \mathcal{L}(F; G))$. We can also go back from $\mathcal{L}(E; \mathcal{L}(F; G))$ to $\mathcal{L}_2(E, F; G)$. We summarize all this in the following proposition.

Proposition 26.45. *Let E, F, G be three normed vector spaces. The map $f \mapsto \varphi$, from $\mathcal{L}_2(E, F; G)$ to $\mathcal{L}(E; \mathcal{L}(F; G))$, defined such that, for every $f \in \mathcal{L}_2(E, F; G)$,*

$$\varphi(u)(v) = f(u, v),$$

is an isomorphism of vector spaces, and furthermore, $\|\varphi\| = \|f\|$.

As a corollary of Proposition 26.45, we get the following proposition which will be useful when we define second-order derivatives.

Proposition 26.46. *Let E, F be normed vector spaces. The map app from $\mathcal{L}(E; F) \times E$ to F , defined such that, for every $f \in \mathcal{L}(E; F)$, for every $u \in E$,*

$$\text{app}(f, u) = f(u),$$

is a continuous bilinear map.

Remark: If E and F are nontrivial, it can be shown that $\|\text{app}\| = 1$. It can also be shown that composition

$$\circ: \mathcal{L}(E; F) \times \mathcal{L}(F; G) \rightarrow \mathcal{L}(E; G),$$

is bilinear and continuous.

The above propositions and definition generalize to arbitrary n -multilinear maps, with $n \geq 2$. Proposition 26.44 extends in the obvious way to any n -multilinear map $f: E_1 \times \cdots \times E_n \rightarrow F$, but condition (3) becomes:

There is a constant $k \geq 0$ such that,

$$\|f(u_1, \dots, u_n)\| \leq k\|u_1\| \cdots \|u_n\|, \text{ for all } u_1 \in E_1, \dots, u_n \in E_n.$$

Definition 26.29 also extends easily to

$$\begin{aligned} \|f\| &= \min \{k \geq 0 \mid \|f(x_1, \dots, x_n)\| \leq k\|x_1\| \cdots \|x_n\|, \text{ for all } x_i \in E_i, 1 \leq i \leq n\} \\ &= \max \{\|f(x_1, \dots, x_n)\| \mid \|x_1\|, \dots, \|x_n\| \leq 1\}. \end{aligned}$$

Proposition 26.45 is also easily extended, and we get an isomorphism between continuous n -multilinear maps in $\mathcal{L}_n(E_1, \dots, E_n; F)$, and continuous linear maps in

$$\mathcal{L}(E_1; \mathcal{L}(E_2; \dots; \mathcal{L}(E_n; F)))$$

An obvious extension of Proposition 26.46 also holds.

Definition 26.30. A normed vector space $(E, \|\cdot\|)$ over \mathbb{R} (or \mathbb{C}) which is a complete metric space for the distance $\|v - u\|$, is called a *Banach space*.

It can be shown that every normed vector space of finite dimension is a Banach space (is complete). It can also be shown that if E and F are normed vector spaces, and F is a Banach space, then $\mathcal{L}(E; F)$ is a Banach space. If E, F and G are normed vector spaces, and G is a Banach space, then $\mathcal{L}_2(E, F; G)$ is a Banach space.

Finally, we consider normed affine spaces.

26.7 Normed Affine Spaces

For geometric applications, we will need to consider affine spaces (E, \vec{E}) where the associated space of translations \vec{E} is a vector space equipped with a norm.

Definition 26.31. Given an affine space (E, \vec{E}) , where the space of translations \vec{E} is a vector space over \mathbb{R} or \mathbb{C} , we say that (E, \vec{E}) is a *normed affine space* if \vec{E} is a normed vector space with norm $\| \cdot \|$.

Given a normed affine space, there is a natural metric on E itself, defined such that

$$d(a, b) = \|\vec{ab}\|.$$

Observe that this metric is invariant under translation, that is,

$$d(a + u, b + u) = d(a, b).$$

Also, for every fixed $a \in E$ and $\lambda > 0$, if we consider the map $h: E \rightarrow E$, defined such that,

$$h(x) = a + \lambda \vec{ax},$$

then $d(h(x), h(y)) = \lambda d(x, y)$.

Note that the map $(a, b) \mapsto \vec{ab}$ from $E \times E$ to \vec{E} is continuous, and similarly for the map $a \mapsto a + u$ from $E \times \vec{E}$ to E . In fact, the map $u \mapsto a + u$ is a homeomorphism from \vec{E} to E_a .

Of course, \mathbb{R}^n is a normed affine space under the Euclidean metric, and it is also complete.

If an affine space E is a finite direct sum $(E_1, a_1) \oplus \cdots \oplus (E_m, a_m)$, and each E_i is also a normed affine space with norm $\| \cdot \|_i$, we make $(E_1, a_1) \oplus \cdots \oplus (E_m, a_m)$ into a normed affine space, by giving it the norm

$$\|(x_1, \dots, x_n)\| = \max(\|x_1\|_1, \dots, \|x_n\|_n).$$

Similarly, the finite product $E_1 \times \cdots \times E_m$ is made into a normed affine space, under the same norm.

We are now ready to define the derivative (or differential) of a map between two normed affine spaces. This will lead to tangent spaces to curves and surfaces (in normed affine spaces).

26.8 Futher Readings

A thorough treatment of general topology can be found in Munkres [84, 85], Dixmier [29], Lang [70], Schwartz [93, 92], Bredon [18], and the classic, Seifert and Threlfall [95].

Chapter 27

A Detour On Fractals

27.1 Iterated Function Systems and Fractals

A pleasant application of the Hausdorff distance and of the fixed point theorem for contracting mappings is a method for defining a class of “self-similar” fractals. For this, we can use iterated function systems.

Definition 27.1. Given a metric space, (X, d) , an *iterated function system*, for short, an *ifs*, is a finite sequence of functions, (f_1, \dots, f_n) , where each $f_i: X \rightarrow X$ is a contracting mapping. A nonempty compact subset, K , of X is an *invariant set (or attractor)* for the ifs, (f_1, \dots, f_n) , if

$$K = f_1(K) \cup \dots \cup f_n(K).$$

The major result about ifs's is the following:

Theorem 27.1. *If (X, d) is a nonempty complete metric space, then every iterated function system, (f_1, \dots, f_n) , has a unique invariant set, A , which is a nonempty compact subset of X . Furthermore, for every nonempty compact subset, A_0 , of X , this invariant set, A , is the limit of the sequence, (A_m) , where $A_{m+1} = f_1(A_m) \cup \dots \cup f_n(A_m)$.*

Proof. Since X is complete, by Theorem 26.39, the space $(\mathcal{K}(X), D)$ is a complete metric space. The theorem will follow from Theorem 26.38 if we can show that the map, $F: \mathcal{K}(X) \rightarrow \mathcal{K}(X)$, defined such that

$$F(K) = f_1(K) \cup \dots \cup f_n(K),$$

for every nonempty compact set, K , is a contracting mapping. Let A, B be any two nonempty compact subsets of X and consider any $\eta \geq D(A, B)$. Since each $f_i: X \rightarrow X$ is a contracting mapping, there is some λ_i , with $0 \leq \lambda_i < 1$, such that

$$d(f_i(a), f_i(b)) \leq \lambda_i d(a, b),$$

for all $a, b \in X$. Let $\lambda = \max\{\lambda_1, \dots, \lambda_n\}$. We claim that

$$D(F(A), F(B)) \leq \lambda D(A, B).$$

For any $x \in F(A) = f_1(A) \cup \dots \cup f_n(A)$, there is some $a_i \in A_i$ such that $x = f_i(a_i)$ and since $\eta \geq D(A, B)$, there is some $b_i \in B$ such that

$$d(a_i, b_i) \leq \eta,$$

and thus,

$$d(x, f_i(b_i)) = d(f_i(a_i), f_i(b_i)) \leq \lambda_i d(a_i, b_i) \leq \lambda \eta.$$

This show that

$$F(A) \subseteq V_{\lambda\eta}(F(B)).$$

Similarly, we can prove that

$$F(B) \subseteq V_{\lambda\eta}(F(A)),$$

and since this holds for all $\eta \geq D(A, B)$, we proved that

$$D(F(A), F(B)) \leq \lambda D(A, B)$$

where $\lambda = \max\{\lambda_1, \dots, \lambda_n\}$. Since $0 \leq \lambda_i < 1$, we have $0 \leq \lambda < 1$ and F is indeed a contracting mapping. \square

Theorem 27.1 justifies the existence of many familiar “self-similar” fractals. One of the best known fractals is the *Sierpinski gasket*.

Example 27.1. Consider an equilateral triangle with vertices a, b, c , and let f_1, f_2, f_3 be the dilatations of centers a, b, c and ratio $1/2$. The Sierpinski gasket is the invariant set of the ifs (f_1, f_2, f_3) . The dilations f_1, f_2, f_3 can be defined explicitly as follows, assuming that $a = (-1/2, 0)$, $b = (1/2, 0)$, and $c = (0, \sqrt{3}/2)$. The contractions f_1, f_2, f_3 are specified by

$$x' = \frac{1}{2}x - \frac{1}{4},$$

$$y' = \frac{1}{2}y,$$

$$x' = \frac{1}{2}x + \frac{1}{4},$$

$$y' = \frac{1}{2}y,$$

and

$$x' = \frac{1}{2}x,$$

$$y' = \frac{1}{2}y + \frac{\sqrt{3}}{4}.$$

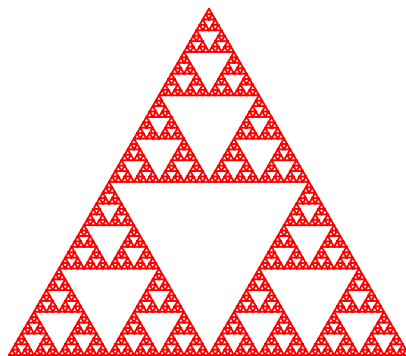


Figure 27.1: The Sierpinski gasket

We wrote a *Mathematica* program that iterates any finite number of affine maps on any input figure consisting of combinations of points, line segments, and polygons (with their interior points). Starting with the edges of the triangle a, b, c , after 6 iterations, we get the picture shown in Figure 27.1.

It is amusing that the same fractal is obtained no matter what the initial nonempty compact figure is. It is interesting to see what happens if we start with a solid triangle (with its interior points). The result after 6 iterations is shown in Figure 27.2. The convergence towards the Sierpinski gasket is very fast. Incidentally, there are many other ways of defining the Sierpinski gasket.

A nice variation on the theme of the Sierpinski gasket is the *Sierpinski dragon*.

Example 27.2. The Sierpinski dragon is specified by the following three contractions:

$$\begin{aligned} x' &= -\frac{1}{4}x - \frac{\sqrt{3}}{4}y + \frac{3}{4}, \\ y' &= \frac{\sqrt{3}}{4}x - \frac{1}{4}y + \frac{\sqrt{3}}{4}, \\ x' &= -\frac{1}{4}x + \frac{\sqrt{3}}{4}y - \frac{3}{4}, \\ y' &= -\frac{\sqrt{3}}{4}x - \frac{1}{4}y + \frac{\sqrt{3}}{4}, \\ x' &= \frac{1}{2}x, \\ y' &= \frac{1}{2}y + \frac{\sqrt{3}}{2}. \end{aligned}$$

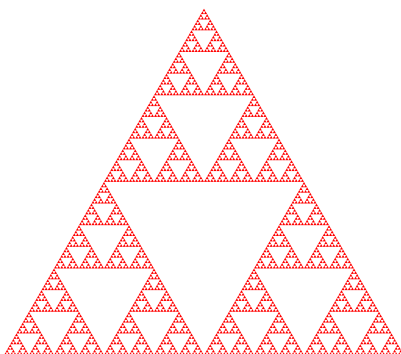


Figure 27.2: The Sierpinski gasket, version 2

The result of 7 iterations starting from the line segment $(-1, 0), (1, 0)$, is shown in Figure 27.3. This curve converges to the boundary of the Sierpinski gasket.

A different kind of fractal is the *Heighway dragon*.

Example 27.3. The Heighway dragon is specified by the following two contractions:

$$\begin{aligned} x' &= \frac{1}{2}x - \frac{1}{2}y, \\ y' &= \frac{1}{2}x + \frac{1}{2}y, \\ x' &= -\frac{1}{2}x - \frac{1}{2}y, \\ y' &= \frac{1}{2}x - \frac{1}{2}y + 1. \end{aligned}$$

It can be shown that for any number of iterations, the polygon does not cross itself. This means that no edge is traversed twice and that if a point is traversed twice, then this point is the endpoint of some edge. The result of 13 iterations, starting with the line segment $((0, 0), (0, 1))$, is shown in Figure 27.4.

The Heighway dragon turns out to fill a closed and bounded set. It can also be shown that the plane can be tiled with copies of the Heighway dragon.

Another well known example is the *Koch curve*.

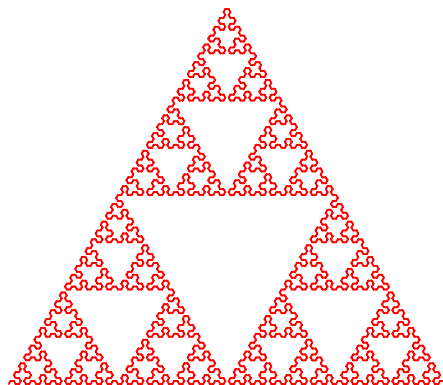


Figure 27.3: The Sierpinski dragon

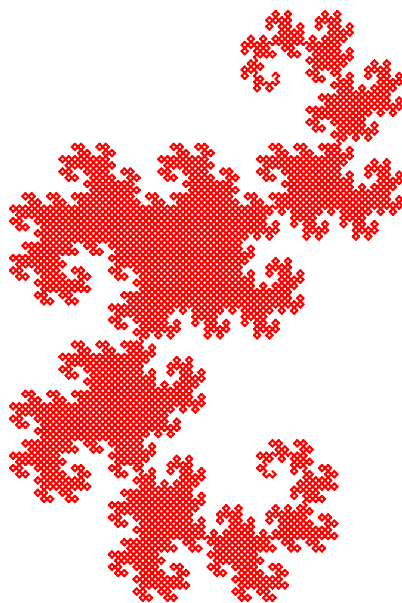


Figure 27.4: The Heighway dragon

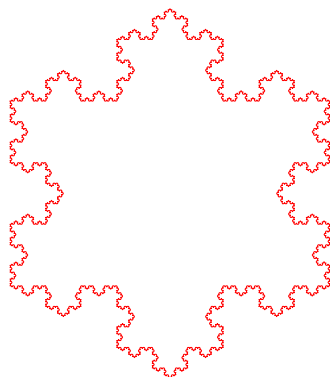


Figure 27.6: The snowflake curve

Example 27.5. The snowflake curve obtained after 5 iterations is shown in Figure 27.6.

The snowflake curve is an example of a closed curve of infinite length bounding a finite area.

We conclude with another famous example, a variant of the *Hilbert curve*.

Example 27.6. This version of the Hilbert curve is defined by the following four contractions:

$$\begin{aligned}
 x' &= \frac{1}{2}x - \frac{1}{2}, \\
 y' &= \frac{1}{2}y + 1, \\
 x' &= \frac{1}{2}x + \frac{1}{2}, \\
 y' &= \frac{1}{2}y + 1, \\
 x' &= -\frac{1}{2}y + 1, \\
 y' &= \frac{1}{2}x + \frac{1}{2}, \\
 x' &= \frac{1}{2}y - 1, \\
 y' &= -\frac{1}{2}x + \frac{1}{2}.
 \end{aligned}$$

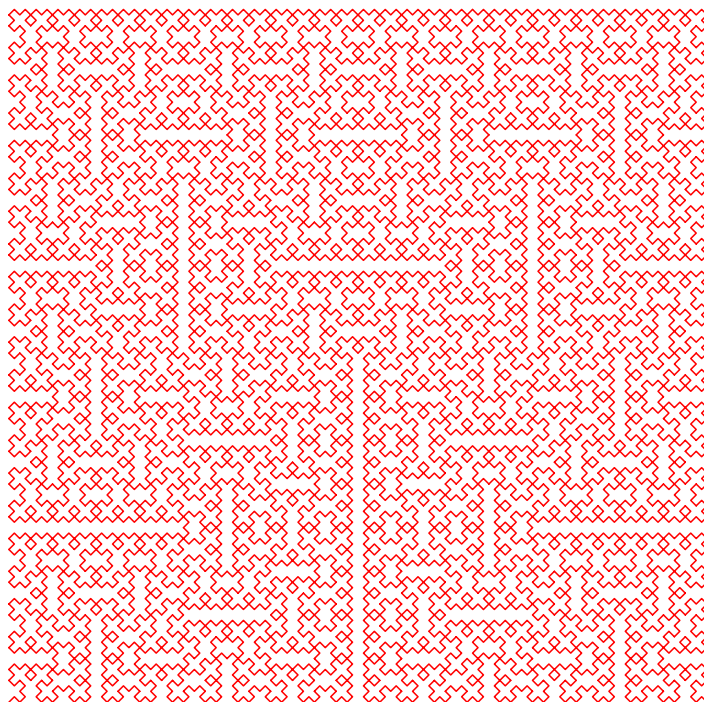


Figure 27.7: A Hilbert curve

This continuous curve is a space-filling curve, in the sense that its image is the entire unit square. The result of 6 iterations, starting with the two line segments $((-1, 0), (0, 1))$ and $((0, 1), (1, 0))$, is shown in Figure 27.7.

For more on iterated function systems and fractals, we recommend Edgar [33].

Chapter 28

Differential Calculus

28.1 Directional Derivatives, Total Derivatives

This chapter contains a review of basic notions of differential calculus. First, we review the definition of the derivative of a function $f: \mathbb{R} \rightarrow \mathbb{R}$. Next, we define directional derivatives and the total derivative of a function $f: E \rightarrow F$ between normed affine spaces. Basic properties of derivatives are shown, including the chain rule. We show how derivatives are represented by Jacobian matrices. The mean value theorem is stated, as well as the implicit function theorem and the inverse function theorem. Diffeomorphisms and local diffeomorphisms are defined. Tangent spaces are defined. Higher-order derivatives are defined, as well as the Hessian. Schwarz's lemma (about the commutativity of partials) is stated. Several versions of Taylor's formula are stated, and a famous formula due to Faà di Bruno's is given.

We first review the notion of the derivative of a real-valued function whose domain is an open subset of \mathbb{R} .

Let $f: A \rightarrow \mathbb{R}$, where A is a nonempty open subset of \mathbb{R} , and consider any $a \in A$. The main idea behind the concept of the derivative of f at a , denoted by $f'(a)$, is that locally around a (that is, in some small open set $U \subseteq A$ containing a), the function f is approximated linearly by the map

$$x \mapsto f(a) + f'(a)(x - a).$$

Part of the difficulty in extending this idea to more complex spaces is to give an adequate notion of linear approximation. Of course, we will use linear maps! Let us now review the formal definition of the derivative of a real-valued function.

Definition 28.1. Let A be any nonempty open subset of \mathbb{R} , and let $a \in A$. For any function $f: A \rightarrow \mathbb{R}$, the *derivative of f at $a \in A$* is the limit (if it exists)

$$\lim_{h \rightarrow 0, h \in U} \frac{f(a+h) - f(a)}{h},$$

where $U = \{h \in \mathbb{R} \mid a + h \in A, h \neq 0\}$. This limit is denoted by $f'(a)$, or $Df(a)$, or $\frac{df}{dx}(a)$. If $f'(a)$ exists for every $a \in A$, we say that f is *differentiable on A* . In this case, the map $a \mapsto f'(a)$ is denoted by f' , or Df , or $\frac{df}{dx}$.

Note that since A is assumed to be open, $A - \{a\}$ is also open, and since the function $h \mapsto a + h$ is continuous and U is the inverse image of $A - \{a\}$ under this function, U is indeed open and the definition makes sense.

We can also define $f'(a)$ as follows: there is some function ϵ , such that,

$$f(a + h) = f(a) + f'(a) \cdot h + \epsilon(h)h,$$

whenever $a + h \in A$, where $\epsilon(h)$ is defined for all h such that $a + h \in A$, and

$$\lim_{h \rightarrow 0, h \in U} \epsilon(h) = 0.$$

Remark: We can also define the notion of *derivative of f at a on the left*, and *derivative of f at a on the right*. For example, we say that the *derivative of f at a on the left* is the limit $f'(a_-)$ (if it exists)

$$\lim_{h \rightarrow 0, h \in U} \frac{f(a + h) - f(a)}{h},$$

where $U = \{h \in \mathbb{R} \mid a + h \in A, h < 0\}$.

If a function f as in Definition 28.1 has a derivative $f'(a)$ at a , then it is continuous at a . If f is differentiable on A , then f is continuous on A . The composition of differentiable functions is differentiable.

Remark: A function f has a derivative $f'(a)$ at a iff the derivative of f on the left at a and the derivative of f on the right at a exist, and if they are equal. Also, if the derivative of f on the left at a exists, then f is continuous on the left at a (and similarly on the right).

We would like to extend the notion of derivative to functions $f: A \rightarrow F$, where E and F are normed affine spaces, and A is some nonempty open subset of E . The first difficulty is to make sense of the quotient

$$\frac{f(a + h) - f(a)}{h}.$$

If E and F are normed affine spaces, it will be notationally convenient to assume that the vector space associated with E is denoted by \vec{E} , and that the vector space associated with F is denoted as \vec{F} .

Since F is a normed affine space, making sense of $f(a + h) - f(a)$ is easy: we can define this as $\overrightarrow{f(a)f(a + h)}$, the unique vector translating $f(a)$ to $f(a + h)$. We should note however, that this quantity is a vector and not a point. Nevertheless, in defining derivatives, it is

notationally more pleasant to denote $\overrightarrow{f(a)f(a+h)}$ by $f(a+h) - f(a)$. Thus, in the rest of this chapter, the vector \overrightarrow{ab} will be denoted by $b - a$. But now, how do we define the quotient by a vector? Well, we don't!

A first possibility is to consider the *directional derivative* with respect to a vector $u \neq 0$ in \overrightarrow{E} . We can consider the vector $f(a+tu) - f(a)$, where $t \in \mathbb{R}$ (or $t \in \mathbb{C}$). Now,

$$\frac{f(a+tu) - f(a)}{t}$$

makes sense. The idea is that in E , the points of the form $a+tu$ for t in some small interval $[-\epsilon, +\epsilon]$ in \mathbb{R} (or \mathbb{C}) form a line segment $[r, s]$ in A containing a , and that the image of this line segment defines a small curve segment on $f(A)$. This curve segment is defined by the map $t \mapsto f(a+tu)$, from $[r, s]$ to F , and the directional derivative $D_u f(a)$ defines the direction of the tangent line at a to this curve. This leads us to the following definition.

Definition 28.2. Let E and F be two normed affine spaces, let A be a nonempty open subset of E , and let $f: A \rightarrow F$ be any function. For any $a \in A$, for any $u \neq 0$ in \overrightarrow{E} , the *directional derivative of f at a w.r.t. the vector u* , denoted by $D_u f(a)$, is the limit (if it exists)

$$\lim_{t \rightarrow 0, t \in U} \frac{f(a+tu) - f(a)}{t},$$

where $U = \{t \in \mathbb{R} \mid a+tu \in A, t \neq 0\}$ (or $U = \{t \in \mathbb{C} \mid a+tu \in A, t \neq 0\}$).

Since the map $t \mapsto a+tu$ is continuous, and since $A - \{a\}$ is open, the inverse image U of $A - \{a\}$ under the above map is open, and the definition of the limit in Definition 28.2 makes sense.

Remark: Since the notion of limit is purely topological, the existence and value of a directional derivative is independent of the choice of norms in E and F , as long as they are equivalent norms.

The directional derivative is sometimes called the *Gâteaux derivative*.

In the special case where $E = \mathbb{R}$ and $F = \mathbb{R}$, and we let $u = 1$ (i.e., the real number 1, viewed as a vector), it is immediately verified that $D_1 f(a) = f'(a)$, in the sense of Definition 28.1. When $E = \mathbb{R}$ (or $E = \mathbb{C}$) and F is any normed vector space, the derivative $D_1 f(a)$, also denoted by $f'(a)$, provides a suitable generalization of the notion of derivative.

However, when E has dimension ≥ 2 , directional derivatives present a serious problem, which is that their definition is not sufficiently uniform. Indeed, there is no reason to believe that the directional derivatives w.r.t. all nonnull vectors u share something in common. As a consequence, a function can have all directional derivatives at a , and yet not be continuous at a . Two functions may have all directional derivatives in some open sets, and yet their composition may not. Thus, we introduce a more uniform notion.

Definition 28.3. Let E and F be two normed affine spaces, let A be a nonempty open subset of E , and let $f: A \rightarrow F$ be any function. For any $a \in A$, we say that f is *differentiable at* $a \in A$ if there is a linear continuous map $L: \vec{E} \rightarrow \vec{F}$ and a function ϵ , such that

$$f(a+h) = f(a) + L(h) + \epsilon(h)\|h\|$$

for every $a+h \in A$, where $\epsilon(h)$ is defined for every h such that $a+h \in A$ and

$$\lim_{h \rightarrow 0, h \in U} \epsilon(h) = 0,$$

where $U = \{h \in \vec{E} \mid a+h \in A, h \neq 0\}$. The linear map L is denoted by $Df(a)$, or Df_a , or $df(a)$, or df_a , or $f'(a)$, and it is called the *Fréchet derivative*, or *derivative*, or *total derivative*, or *total differential*, or *differential*, of f at a .

Since the map $h \mapsto a+h$ from \vec{E} to E is continuous, and since A is open in E , the inverse image U of $A - \{a\}$ under the above map is open in \vec{E} , and it makes sense to say that

$$\lim_{h \rightarrow 0, h \in U} \epsilon(h) = 0.$$

Note that for every $h \in U$, since $h \neq 0$, $\epsilon(h)$ is uniquely determined since

$$\epsilon(h) = \frac{f(a+h) - f(a) - L(h)}{\|h\|},$$

and that the value $\epsilon(0)$ plays absolutely no role in this definition. The condition for f to be differentiable at a amounts to the fact that

$$\lim_{h \rightarrow 0} \frac{\|f(a+h) - f(a) - L(h)\|}{\|h\|} = 0$$

as $h \neq 0$ approaches 0, when $a+h \in A$. However, it does no harm to assume that $\epsilon(0) = 0$, and we will assume this from now on.

Again, we note that the derivative $Df(a)$ of f at a provides an affine approximation of f , locally around a .

Remark: Since the notion of limit is purely topological, the existence and value of a derivative is independent of the choice of norms in E and F , as long as they are equivalent norms.

Note that the continuous linear map L is unique, if it exists. In fact, the next proposition implies this as a corollary. The following proposition shows that our new definition is consistent with the definition of the directional derivative.

Proposition 28.1. *Let E and F be two normed affine spaces, let A be a nonempty open subset of E , and let $f: A \rightarrow F$ be any function. For any $a \in A$, if $Df(a)$ is defined, then f is continuous at a and f has a directional derivative $D_u f(a)$ for every $u \neq 0$ in \vec{E} , and furthermore,*

$$D_u f(a) = Df(a)(u).$$

Proof. If $h \neq 0$ approaches 0, since L is continuous, $\epsilon(h)\|h\|$ approaches 0, and thus, f is continuous at a . For any $u \neq 0$ in \vec{E} , for $|t| \in \mathbb{R}$ small enough (where $t \in \mathbb{R}$ or $t \in \mathbb{C}$), we have $a + tu \in A$, and letting $h = tu$, we have

$$f(a + tu) = f(a) + tL(u) + \epsilon(tu)|t|\|u\|,$$

and for $t \neq 0$,

$$\frac{f(a + tu) - f(a)}{t} = L(u) + \frac{|t|}{t}\epsilon(tu)\|u\|,$$

and the limit when $t \neq 0$ approaches 0 is indeed $D_u f(a)$. \square

The uniqueness of L follows from Proposition 28.1. Also, when E is of finite dimension, it is easily shown that every linear map is continuous, and this assumption is then redundant.

It is important to note that the derivative $Df(a)$ of f at a is a continuous linear map from the vector space \vec{E} to the vector space \vec{F} , and not a function from the affine space E to the affine space F .

As an example, consider the map, $f: M_n(\mathbb{R}) \rightarrow M_n(\mathbb{R})$, given by

$$f(A) = A^\top A - I,$$

where $M_n(\mathbb{R})$ is equipped with any matrix norm, since they are all equivalent; for example, pick the Frobenius norm, $\|A\|_F = \sqrt{\text{tr}(A^\top A)}$. We claim that

$$Df(A)(H) = A^\top H + H^\top A, \quad \text{for all } A \text{ and } H \text{ in } M_n(\mathbb{R}).$$

We have

$$\begin{aligned} f(A + H) - f(A) - (A^\top H + H^\top A) &= (A + H)^\top (A + H) - I - (A^\top A - I) - A^\top H - H^\top A \\ &= A^\top A + A^\top H + H^\top A + H^\top H - A^\top A - A^\top H - H^\top A \\ &= H^\top H. \end{aligned}$$

It follows that

$$\epsilon(H) = \frac{f(A + H) - f(A) - (A^\top H + H^\top A)}{\|H\|} = \frac{H^\top H}{\|H\|},$$

and since our norm is the Frobenius norm,

$$\|\epsilon(H)\| = \left\| \frac{H^\top H}{\|H\|} \right\| \leq \frac{\|H^\top\| \|H\|}{\|H\|} = \|H^\top\| = \|H\|,$$

so

$$\lim_{H \rightarrow 0} \epsilon(H) = 0,$$

and we conclude that

$$Df(A)(H) = A^\top H + H^\top A.$$

If $Df(a)$ exists for every $a \in A$, we get a map

$$Df: A \rightarrow \mathcal{L}(\vec{E}; \vec{F}),$$

called the *derivative of f on A* , and also denoted by df . Recall that $\mathcal{L}(\vec{E}; \vec{F})$ denotes the vector space of all continuous maps from \vec{E} to \vec{F} .

When E is of finite dimension n , for any frame $(a_0, (u_1, \dots, u_n))$ of E , where (u_1, \dots, u_n) is a basis of \vec{E} , we can define the directional derivatives with respect to the vectors in the basis (u_1, \dots, u_n) (actually, we can also do it for an infinite frame). This way, we obtain the definition of partial derivatives, as follows.

Definition 28.4. For any two normed affine spaces E and F , if E is of finite dimension n , for every frame $(a_0, (u_1, \dots, u_n))$ for E , for every $a \in E$, for every function $f: E \rightarrow F$, the directional derivatives $D_{u_j}f(a)$ (if they exist) are called the *partial derivatives of f with respect to the frame $(a_0, (u_1, \dots, u_n))$* . The partial derivative $D_{u_j}f(a)$ is also denoted by $\partial_j f(a)$, or $\frac{\partial f}{\partial x_j}(a)$.

The notation $\frac{\partial f}{\partial x_j}(a)$ for a partial derivative, although customary and going back to Leibniz, is a “logical obscenity.” Indeed, the variable x_j really has nothing to do with the formal definition. This is just another of these situations where tradition is just too hard to overthrow!

We now consider a number of standard results about derivatives.

Proposition 28.2. *Given two normed affine spaces E and F , if $f: E \rightarrow F$ is a constant function, then $Df(a) = 0$, for every $a \in E$. If $f: E \rightarrow F$ is a continuous affine map, then $Df(a) = f$, for every $a \in E$, the linear map associated with f .*

Proof. Straightforward. □

Proposition 28.3. *Given a normed affine space E and a normed vector space F , for any two functions $f, g: E \rightarrow F$, for every $a \in E$, if $Df(a)$ and $Dg(a)$ exist, then $D(f+g)(a)$ and $D(\lambda f)(a)$ exist, and*

$$\begin{aligned} D(f+g)(a) &= Df(a) + Dg(a), \\ D(\lambda f)(a) &= \lambda Df(a). \end{aligned}$$

Proof. Straightforward. □

Proposition 28.4. *Given three normed vector spaces E_1 , E_2 , and F , for any continuous bilinear map*

$f: E_1 \times E_2 \rightarrow F$, for every $(a, b) \in E_1 \times E_2$, $Df(a, b)$ exists, and for every $u \in E_1$ and $v \in E_2$,

$$Df(a, b)(u, v) = f(u, b) + f(a, v).$$

Proof. Straightforward. □

We now state the very useful *chain rule*.

Theorem 28.5. *Given three normed affine spaces E , F , and G , let A be an open set in E , and let B an open set in F . For any functions $f: A \rightarrow F$ and $g: B \rightarrow G$, such that $f(A) \subseteq B$, for any $a \in A$, if $Df(a)$ exists and $Dg(f(a))$ exists, then $D(g \circ f)(a)$ exists, and*

$$D(g \circ f)(a) = Dg(f(a)) \circ Df(a).$$

Proof. It is not difficult, but more involved than the previous two. □

Theorem 28.5 has many interesting consequences. We mention two corollaries.

Proposition 28.6. *Given three normed affine spaces E , F , and G , for any open subset A in E , for any $a \in A$, let $f: A \rightarrow F$ such that $Df(a)$ exists, and let $g: F \rightarrow G$ be a continuous affine map. Then, $D(g \circ f)(a)$ exists, and*

$$D(g \circ f)(a) = g \circ Df(a),$$

where g is the linear map associated with the affine map g .

Proposition 28.7. *Given two normed affine spaces E and F , let A be some open subset in E , let B be some open subset in F , let $f: A \rightarrow B$ be a bijection from A to B , and assume that Df exists on A and that Df^{-1} exists on B . Then, for every $a \in A$,*

$$Df^{-1}(f(a)) = (Df(a))^{-1}.$$

Proposition 28.7 has the remarkable consequence that the two vector spaces \vec{E} and \vec{F} have the same dimension. In other words, a local property, the existence of a bijection f between an open set A of E and an open set B of F , such that f is differentiable on A and f^{-1} is differentiable on B , implies a global property, that the two vector spaces \vec{E} and \vec{F} have the same dimension.

We now consider the situation where the normed affine space F is a finite direct sum $F = (F_1, b_1) \oplus \cdots \oplus (F_m, b_m)$.

Proposition 28.8. *Given normed affine spaces E and $F = (F_1, b_1) \oplus \cdots \oplus (F_m, b_m)$, given any open subset A of E , for any $a \in A$, for any function $f: A \rightarrow F$, letting $f = (f_1, \dots, f_m)$, $Df(a)$ exists iff every $Df_i(a)$ exists, and*

$$Df(a) = in_1 \circ Df_1(a) + \cdots + in_m \circ Df_m(a).$$

Proof. Observe that $f(a+h) - f(a) = (f(a+h) - b) - (f(a) - b)$, where $b = (b_1, \dots, b_m)$, and thus, as far as dealing with derivatives, $Df(a)$ is equal to $Df_b(a)$, where $f_b: E \rightarrow \vec{F}$ is defined such that $f_b(x) = f(x) - b$, for every $x \in E$. Thus, we can work with the vector space \vec{F} instead of the affine space F . The proposition is then a simple application of Theorem 28.5. □

In the special case where F is a normed affine space of finite dimension m , for any frame $(b_0, (v_1, \dots, v_m))$ of F , where (v_1, \dots, v_m) is a basis of \vec{F} , every point $x \in F$ can be expressed uniquely as

$$x = b_0 + x_1 v_1 + \dots + x_m v_m,$$

where $(x_1, \dots, x_m) \in K^m$, the coordinates of x in the frame $(b_0, (v_1, \dots, v_m))$ (where $K = \mathbb{R}$ or $K = \mathbb{C}$). Thus, letting F_i be the standard normed affine space K with its natural structure, we note that F is isomorphic to the direct sum $F = (K, 0) \oplus \dots \oplus (K, 0)$. Then, every function $f: E \rightarrow F$ is represented by m functions (f_1, \dots, f_m) , where $f_i: E \rightarrow K$ (where $K = \mathbb{R}$ or $K = \mathbb{C}$), and

$$f(x) = b_0 + f_1(x)v_1 + \dots + f_m(x)v_m,$$

for every $x \in E$. The following proposition is an immediate corollary of Proposition 28.8.

Proposition 28.9. *For any two normed affine spaces E and F , if F is of finite dimension m , for any frame $(b_0, (v_1, \dots, v_m))$ of F , where (v_1, \dots, v_m) is a basis of \vec{F} , for every $a \in E$, a function $f: E \rightarrow F$ is differentiable at a iff each f_i is differentiable at a , and*

$$Df(a)(u) = Df_1(a)(u)v_1 + \dots + Df_m(a)(u)v_m,$$

for every $u \in \vec{E}$.

We now consider the situation where E is a finite direct sum. Given a normed affine space $E = (E_1, a_1) \oplus \dots \oplus (E_n, a_n)$ and a normed affine space F , given any open subset A of E , for any $c = (c_1, \dots, c_n) \in A$, we define the continuous functions $i_j^c: E_j \rightarrow E$, such that

$$i_j^c(x) = (c_1, \dots, c_{j-1}, x, c_{j+1}, \dots, c_n).$$

For any function $f: A \rightarrow F$, we have functions $f \circ i_j^c: E_j \rightarrow F$, defined on $(i_j^c)^{-1}(A)$, which contains c_j . If $D(f \circ i_j^c)(c_j)$ exists, we call it the *partial derivative of f w.r.t. its j th argument, at c* . We also denote this derivative by $D_j f(c)$. Note that $D_j f(c) \in \mathcal{L}(\vec{E}_j; \vec{F})$.

This notion is a generalization of the notion defined in Definition 28.4. In fact, when E is of dimension n , and a frame $(a_0, (u_1, \dots, u_n))$ has been chosen, we can write $E = (E_1, a_1) \oplus \dots \oplus (E_n, a_n)$, for some obvious (E_j, a_j) (as explained just after Proposition 28.8), and then

$$D_j f(c)(\lambda u_j) = \lambda \partial_j f(c),$$

and the two notions are consistent.

The definition of i_j^c and of $D_j f(c)$ also makes sense for a finite product $E_1 \times \dots \times E_n$ of affine spaces E_i . We will use freely the notation $\partial_j f(c)$ instead of $D_j f(c)$.

The notion $\partial_j f(c)$ introduced in Definition 28.4 is really that of the vector derivative, whereas $D_j f(c)$ is the corresponding linear map. Although perhaps confusing, we identify the two notions. The following proposition holds.

Proposition 28.10. *Given a normed affine space $E = (E_1, a_1) \oplus \cdots \oplus (E_n, a_n)$, and a normed affine space F , given any open subset A of E , for any function $f: A \rightarrow F$, for every $c \in A$, if $Df(c)$ exists, then each $D_j f(c)$ exists, and*

$$Df(c)(u_1, \dots, u_n) = D_1 f(c)(u_1) + \cdots + D_n f(c)(u_n),$$

for every $u_i \in E_i$, $1 \leq i \leq n$. The same result holds for the finite product $E_1 \times \cdots \times E_n$.

Proof. Since every $c \in E$ can be written as $c = a + c - a$, where $a = (a_1, \dots, a_n)$, defining $f_a: \vec{E} \rightarrow F$ such that, $f_a(u) = f(a + u)$, for every $u \in \vec{E}$, clearly, $Df(c) = Df_a(c - a)$, and thus, we can work with the function f_a whose domain is the vector space \vec{E} . The proposition is then a simple application of Theorem 28.5. \square

28.2 Jacobian Matrices

If both E and F are of finite dimension, for any frame $(a_0, (u_1, \dots, u_n))$ of E and any frame $(b_0, (v_1, \dots, v_m))$ of F , every function $f: E \rightarrow F$ is determined by m functions $f_i: E \rightarrow \mathbb{R}$ (or $f_i: E \rightarrow \mathbb{C}$), where

$$f(x) = b_0 + f_1(x)v_1 + \cdots + f_m(x)v_m,$$

for every $x \in E$. From Proposition 28.1, we have

$$Df(a)(u_j) = D_{u_j} f(a) = \partial_j f(a),$$

and from Proposition 28.9, we have

$$Df(a)(u_j) = Df_1(a)(u_j)v_1 + \cdots + Df_i(a)(u_j)v_i + \cdots + Df_m(a)(u_j)v_m,$$

that is,

$$Df(a)(u_j) = \partial_j f_1(a)v_1 + \cdots + \partial_j f_i(a)v_i + \cdots + \partial_j f_m(a)v_m.$$

Since the j -th column of the $m \times n$ -matrix representing $Df(a)$ w.r.t. the bases (u_1, \dots, u_n) and (v_1, \dots, v_m) is equal to the components of the vector $Df(a)(u_j)$ over the basis (v_1, \dots, v_m) , the linear map $Df(a)$ is determined by the $m \times n$ -matrix $J(f)(a) = (\partial_j f_i(a))$, (or $J(f)(a) = (\frac{\partial f_i}{\partial x_j}(a))$):

$$J(f)(a) = \begin{pmatrix} \partial_1 f_1(a) & \partial_2 f_1(a) & \cdots & \partial_n f_1(a) \\ \partial_1 f_2(a) & \partial_2 f_2(a) & \cdots & \partial_n f_2(a) \\ \vdots & \vdots & \ddots & \vdots \\ \partial_1 f_m(a) & \partial_2 f_m(a) & \cdots & \partial_n f_m(a) \end{pmatrix}$$

or

$$J(f)(a) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1}(a) & \frac{\partial f_1}{\partial x_2}(a) & \cdots & \frac{\partial f_1}{\partial x_n}(a) \\ \frac{\partial f_2}{\partial x_1}(a) & \frac{\partial f_2}{\partial x_2}(a) & \cdots & \frac{\partial f_2}{\partial x_n}(a) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(a) & \frac{\partial f_m}{\partial x_2}(a) & \cdots & \frac{\partial f_m}{\partial x_n}(a) \end{pmatrix}$$

This matrix is called the *Jacobian matrix* of Df at a . When $m = n$, the determinant, $\det(J(f)(a))$, of $J(f)(a)$ is called the *Jacobian* of $Df(a)$. From a previous remark, we know that this determinant in fact only depends on $Df(a)$, and not on specific bases. However, partial derivatives give a means for computing it.

When $E = \mathbb{R}^n$ and $F = \mathbb{R}^m$, for any function $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$, it is easy to compute the partial derivatives $\frac{\partial f_i}{\partial x_j}(a)$. We simply treat the function $f_i: \mathbb{R}^n \rightarrow \mathbb{R}$ as a function of its j -th argument, leaving the others fixed, and compute the derivative as in Definition 28.1, that is, the usual derivative.

Example 28.1. For example, consider the function $f: \mathbb{R}^2 \rightarrow \mathbb{R}^2$, defined such that

$$f(r, \theta) = (r \cos(\theta), r \sin(\theta)).$$

Then, we have

$$J(f)(r, \theta) = \begin{pmatrix} \cos(\theta) & -r \sin(\theta) \\ \sin(\theta) & r \cos(\theta) \end{pmatrix}$$

and the Jacobian (determinant) has value $\det(J(f)(r, \theta)) = r$.

In the case where $E = \mathbb{R}$ (or $E = \mathbb{C}$), for any function $f: \mathbb{R} \rightarrow F$ (or $f: \mathbb{C} \rightarrow F$), the Jacobian matrix of $Df(a)$ is a column vector. In fact, this column vector is just $D_1 f(a)$. Then, for every $\lambda \in \mathbb{R}$ (or $\lambda \in \mathbb{C}$),

$$Df(a)(\lambda) = \lambda D_1 f(a).$$

This case is sufficiently important to warrant a definition.

Definition 28.5. Given a function $f: \mathbb{R} \rightarrow F$ (or $f: \mathbb{C} \rightarrow F$), where F is a normed affine space, the vector

$$Df(a)(1) = D_1 f(a)$$

is called the *vector derivative* or *velocity vector* (in the real case) at a . We usually identify $Df(a)$ with its Jacobian matrix $D_1 f(a)$, which is the column vector corresponding to $D_1 f(a)$. By abuse of notation, we also let $Df(a)$ denote the vector $Df(a)(1) = D_1 f(a)$.

When $E = \mathbb{R}$, the physical interpretation is that f defines a (parametric) curve that is the trajectory of some particle moving in \mathbb{R}^m as a function of time, and the vector $D_1f(a)$ is the *velocity* of the moving particle $f(t)$ at $t = a$.

It is often useful to consider functions $f: [a, b] \rightarrow F$ from a closed interval $[a, b] \subseteq \mathbb{R}$ to a normed affine space F , and its derivative $Df(a)$ on $[a, b]$, even though $[a, b]$ is not open. In this case, as in the case of a real-valued function, we define the right derivative $D_1f(a_+)$ at a , and the left derivative $D_1f(b_-)$ at b , and we assume their existence.

Example 28.2.

1. When $E = [0, 1]$, and $F = \mathbb{R}^3$, a function $f: [0, 1] \rightarrow \mathbb{R}^3$ defines a (parametric) curve in \mathbb{R}^3 . Letting $f = (f_1, f_2, f_3)$, its Jacobian matrix at $a \in \mathbb{R}$ is

$$J(f)(a) = \begin{pmatrix} \frac{\partial f_1}{\partial t}(a) \\ \frac{\partial f_2}{\partial t}(a) \\ \frac{\partial f_3}{\partial t}(a) \end{pmatrix}$$

2. When $E = \mathbb{R}^2$, and $F = \mathbb{R}^3$, a function $\varphi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$ defines a parametric surface. Letting $\varphi = (f, g, h)$, its Jacobian matrix at $a \in \mathbb{R}^2$ is

$$J(\varphi)(a) = \begin{pmatrix} \frac{\partial f}{\partial u}(a) & \frac{\partial f}{\partial v}(a) \\ \frac{\partial g}{\partial u}(a) & \frac{\partial g}{\partial v}(a) \\ \frac{\partial h}{\partial u}(a) & \frac{\partial h}{\partial v}(a) \end{pmatrix}$$

3. When $E = \mathbb{R}^3$, and $F = \mathbb{R}$, for a function $f: \mathbb{R}^3 \rightarrow \mathbb{R}$, the Jacobian matrix at $a \in \mathbb{R}^3$ is

$$J(f)(a) = \begin{pmatrix} \frac{\partial f}{\partial x}(a) & \frac{\partial f}{\partial y}(a) & \frac{\partial f}{\partial z}(a) \end{pmatrix}.$$

More generally, when $f: \mathbb{R}^n \rightarrow \mathbb{R}$, the Jacobian matrix at $a \in \mathbb{R}^n$ is the row vector

$$J(f)(a) = \left(\frac{\partial f}{\partial x_1}(a) \cdots \frac{\partial f}{\partial x_n}(a) \right).$$

Its transpose is a column vector called the *gradient* of f at a , denoted by $\text{grad}f(a)$ or $\nabla f(a)$. Then, given any $v \in \mathbb{R}^n$, note that

$$Df(a)(v) = \frac{\partial f}{\partial x_1}(a) v_1 + \cdots + \frac{\partial f}{\partial x_n}(a) v_n = \text{grad}f(a) \cdot v,$$

the scalar product of $\text{grad}f(a)$ and v .

When E , F , and G have finite dimensions, and $(a_0, (u_1, \dots, u_p))$ is an affine frame for E , $(b_0, (v_1, \dots, v_n))$ is an affine frame for F , and $(c_0, (w_1, \dots, w_m))$ is an affine frame for G , if A is an open subset of E , B is an open subset of F , for any functions $f: A \rightarrow F$ and $g: B \rightarrow G$, such that $f(A) \subseteq B$, for any $a \in A$, letting $b = f(a)$, and $h = g \circ f$, if $Df(a)$ exists and $Dg(b)$ exists, by Theorem 28.5, the Jacobian matrix $J(h)(a) = J(g \circ f)(a)$ w.r.t. the bases (u_1, \dots, u_p) and (w_1, \dots, w_m) is the product of the Jacobian matrices $J(g)(b)$ w.r.t. the bases (v_1, \dots, v_n) and (w_1, \dots, w_m) , and $J(f)(a)$ w.r.t. the bases (u_1, \dots, u_p) and (v_1, \dots, v_n) :

$$J(h)(a) = \begin{pmatrix} \partial_1 g_1(b) & \partial_2 g_1(b) & \dots & \partial_n g_1(b) \\ \partial_1 g_2(b) & \partial_2 g_2(b) & \dots & \partial_n g_2(b) \\ \vdots & \vdots & \ddots & \vdots \\ \partial_1 g_m(b) & \partial_2 g_m(b) & \dots & \partial_n g_m(b) \end{pmatrix} \begin{pmatrix} \partial_1 f_1(a) & \partial_2 f_1(a) & \dots & \partial_p f_1(a) \\ \partial_1 f_2(a) & \partial_2 f_2(a) & \dots & \partial_p f_2(a) \\ \vdots & \vdots & \ddots & \vdots \\ \partial_1 f_n(a) & \partial_2 f_n(a) & \dots & \partial_p f_n(a) \end{pmatrix}$$

or

$$J(h)(a) = \begin{pmatrix} \frac{\partial g_1}{\partial y_1}(b) & \frac{\partial g_1}{\partial y_2}(b) & \dots & \frac{\partial g_1}{\partial y_n}(b) \\ \frac{\partial g_2}{\partial y_1}(b) & \frac{\partial g_2}{\partial y_2}(b) & \dots & \frac{\partial g_2}{\partial y_n}(b) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial g_m}{\partial y_1}(b) & \frac{\partial g_m}{\partial y_2}(b) & \dots & \frac{\partial g_m}{\partial y_n}(b) \end{pmatrix} \begin{pmatrix} \frac{\partial f_1}{\partial x_1}(a) & \frac{\partial f_1}{\partial x_2}(a) & \dots & \frac{\partial f_1}{\partial x_p}(a) \\ \frac{\partial f_2}{\partial x_1}(a) & \frac{\partial f_2}{\partial x_2}(a) & \dots & \frac{\partial f_2}{\partial x_p}(a) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1}(a) & \frac{\partial f_n}{\partial x_2}(a) & \dots & \frac{\partial f_n}{\partial x_p}(a) \end{pmatrix}.$$

Thus, we have the familiar formula

$$\frac{\partial h_i}{\partial x_j}(a) = \sum_{k=1}^{k=n} \frac{\partial g_i}{\partial y_k}(b) \frac{\partial f_k}{\partial x_j}(a).$$

Given two normed affine spaces E and F of finite dimension, given an open subset A of E , if a function $f: A \rightarrow F$ is differentiable at $a \in A$, then its Jacobian matrix is well defined.



One should be warned that the converse is false. There are functions such that all the partial derivatives exist at some $a \in A$, but yet, the function is not differentiable at a , and not even continuous at a . For example, consider the function $f: \mathbb{R}^2 \rightarrow \mathbb{R}$, defined such that $f(0, 0) = 0$, and

$$f(x, y) = \frac{x^2 y}{x^4 + y^2} \quad \text{if } (x, y) \neq (0, 0).$$

For any $u \neq 0$, letting $u = \begin{pmatrix} h \\ k \end{pmatrix}$, we have

$$\frac{f(0 + tu) - f(0)}{t} = \frac{h^2 k}{t^2 h^4 + k^2},$$

so that

$$D_u f(0,0) = \begin{cases} \frac{h^2}{k} & \text{if } k \neq 0 \\ 0 & \text{if } k = 0. \end{cases}$$

Thus, $D_u f(0,0)$ exists for all $u \neq 0$. On the other hand, if $Df(0,0)$ existed, it would be a linear map $Df(0,0): \mathbb{R}^2 \rightarrow \mathbb{R}$ represented by a row matrix $(\alpha \ \beta)$, and we would have $D_u f(0,0) = Df(0,0)(u) = \alpha h + \beta k$, but the explicit formula for $D_u f(0,0)$ is not linear. As a matter of fact, the function f is not continuous at $(0,0)$. For example, on the parabola $y = x^2$, $f(x,y) = \frac{1}{2}$, and when we approach the origin on this parabola, the limit is $\frac{1}{2}$, when in fact, $f(0,0) = 0$.

However, there are sufficient conditions on the partial derivatives for $Df(a)$ to exist, namely, continuity of the partial derivatives.

If f is differentiable on A , then f defines a function $Df: A \rightarrow \mathcal{L}(\vec{E}; \vec{F})$. It turns out that the continuity of the partial derivatives on A is a necessary and sufficient condition for Df to exist and to be continuous on A .

If $f: [a,b] \rightarrow \mathbb{R}$ is a function which is continuous on $[a,b]$ and differentiable on $]a,b[$, then there is some c with $a < c < b$ such that

$$f(b) - f(a) = (b - a)f'(c).$$

This result is known as the *mean value theorem* and is a generalization of *Rolle's theorem*, which corresponds to the case where $f(a) = f(b)$.

Unfortunately, the mean value theorem fails for vector-valued functions. For example, the function $f: [0, 2\pi] \rightarrow \mathbb{R}^2$ given by

$$f(t) = (\cos t, \sin t)$$

is such that $f(2\pi) - f(0) = (0,0)$, yet its derivative $f'(t) = (-\sin t, \cos t)$ does not vanish in $]0, 2\pi[$.

A suitable generalization of the mean value theorem to vector-valued functions is possible if we consider an inequality (an upper bound) instead of an equality. This generalized version of the mean value theorem plays an important role in the proof of several major results of differential calculus.

If E is an affine space (over \mathbb{R} or \mathbb{C}), given any two points $a, b \in E$, the *closed segment* $[a, b]$ is the set of all points $a + \lambda(b - a)$, where $0 \leq \lambda \leq 1$, $\lambda \in \mathbb{R}$, and the *open segment* $]a, b[$ is the set of all points $a + \lambda(b - a)$, where $0 < \lambda < 1$, $\lambda \in \mathbb{R}$.

Lemma 28.11. *Let E and F be two normed affine spaces, let A be an open subset of E , and let $f: A \rightarrow F$ be a continuous function on A . Given any $a \in A$ and any $h \neq 0$ in \vec{E} , if the closed segment $[a, a + h]$ is contained in A , if $f: A \rightarrow F$ is differentiable at every point of the open segment $]a, a + h[$, and*

$$\sup_{x \in]a, a+h[} \|Df(x)\| \leq M,$$

for some $M \geq 0$, then

$$\|f(a+h) - f(a)\| \leq M\|h\|.$$

As a corollary, if $L: \vec{E} \rightarrow \vec{F}$ is a continuous linear map, then

$$\|f(a+h) - f(a) - L(h)\| \leq M\|h\|,$$

where $M = \sup_{x \in]a, a+h[} \|Df(x) - L\|$.

The above lemma is sometimes called the “mean value theorem.” Lemma 28.11 can be used to show the following important result.

Theorem 28.12. *Given two normed affine spaces E and F , where E is of finite dimension n , and where $(a_0, (u_1, \dots, u_n))$ is a frame of E , given any open subset A of E , given any function $f: A \rightarrow F$, the derivative $Df: A \rightarrow \mathcal{L}(\vec{E}; \vec{F})$ is defined and continuous on A iff every partial derivative $\partial_j f$ (or $\frac{\partial f}{\partial x_j}$) is defined and continuous on A , for all j , $1 \leq j \leq n$. As a corollary, if F is of finite dimension m , and $(b_0, (v_1, \dots, v_m))$ is a frame of F , the derivative $Df: A \rightarrow \mathcal{L}(\vec{E}; \vec{F})$ is defined and continuous on A iff every partial derivative $\partial_j f_i$ (or $\frac{\partial f_i}{\partial x_j}$) is defined and continuous on A , for all i, j , $1 \leq i \leq m$, $1 \leq j \leq n$.*

Theorem 28.12 gives a necessary and sufficient condition for the existence and continuity of the derivative of a function on an open set. It should be noted that a more general version of Theorem 28.12 holds, assuming that $E = (E_1, a_1) \oplus \dots \oplus (E_n, a_n)$, or $E = E_1 \times \dots \times E_n$, and using the more general partial derivatives $D_j f$ introduced before Proposition 28.10.

Definition 28.6. Given two normed affine spaces E and F , and an open subset A of E , we say that a function $f: A \rightarrow F$ is of class C^0 on A or a C^0 -function on A if f is continuous on A . We say that $f: A \rightarrow F$ is of class C^1 on A or a C^1 -function on A if Df exists and is continuous on A .

Since the existence of the derivative on an open set implies continuity, a C^1 -function is of course a C^0 -function. Theorem 28.12 gives a necessary and sufficient condition for a function f to be a C^1 -function (when E is of finite dimension). It is easy to show that the composition of C^1 -functions (on appropriate open sets) is a C^1 -function.

28.3 The Implicit and The Inverse Function Theorems

Given three normed affine spaces E, F , and G , given a function $f: E \times F \rightarrow G$, given any $c \in G$, it may happen that the equation

$$f(x, y) = c$$

has the property that, for some open sets $A \subseteq E$, and $B \subseteq F$, there is a function $g: A \rightarrow B$, such that

$$f(x, g(x)) = c,$$

for all $x \in A$. Such a situation is usually very rare, but if some solution $(a, b) \in E \times F$ such that $f(a, b) = c$ is known, under certain conditions, for some small open sets $A \subseteq E$ containing a and $B \subseteq F$ containing b , the existence of a unique $g: A \rightarrow B$, such that

$$f(x, g(x)) = c,$$

for all $x \in A$, can be shown. Under certain conditions, it can also be shown that g is continuous, and differentiable. Such a theorem, known as the *implicit function theorem*, can be shown. We state a version of this result below, following Schwartz [94]. The proof (see Schwartz [94]) is fairly involved, and uses a fixed-point theorem for contracting mappings in complete metric spaces. Other proofs can be found in Lang [69] and Cartan [20].

Theorem 28.13. *Let E, F , and G , be normed affine spaces, let Ω be an open subset of $E \times F$, let $f: \Omega \rightarrow G$ be a function defined on Ω , let $(a, b) \in \Omega$, let $c \in G$, and assume that $f(a, b) = c$. If the following assumptions hold*

- (1) *The function $f: \Omega \rightarrow G$ is continuous on Ω ;*
- (2) *F is a complete normed affine space (and so is G);*
- (3) *$\frac{\partial f}{\partial y}(x, y)$ exists for every $(x, y) \in \Omega$, and $\frac{\partial f}{\partial y}: \Omega \rightarrow \mathcal{L}(\vec{F}; \vec{G})$ is continuous;*
- (4) *$\frac{\partial f}{\partial y}(a, b)$ is a bijection of $\mathcal{L}(\vec{F}; \vec{G})$, and $\left(\frac{\partial f}{\partial y}(a, b)\right)^{-1} \in \mathcal{L}(\vec{G}; \vec{F})$;*

then the following properties hold:

- (a) *There exist some open subset $A \subseteq E$ containing a and some open subset $B \subseteq F$ containing b , such that $A \times B \subseteq \Omega$, and for every $x \in A$, the equation $f(x, y) = c$ has a single solution $y = g(x)$, and thus, there is a unique function $g: A \rightarrow B$ such that $f(x, g(x)) = c$, for all $x \in A$;*
- (b) *The function $g: A \rightarrow B$ is continuous.*

If we also assume that

- (5) *The derivative $Df(a, b)$ exists;*

then

- (c) *The derivative $Dg(a)$ exists, and*

$$Dg(a) = -\left(\frac{\partial f}{\partial y}(a, b)\right)^{-1} \circ \frac{\partial f}{\partial x}(a, b);$$

and if in addition

(6) $\frac{\partial f}{\partial x}: \Omega \rightarrow \mathcal{L}(\vec{E}; \vec{G})$ is also continuous (and thus, in view of (3), f is C^1 on Ω);

then

(d) The derivative $Dg: A \rightarrow \mathcal{L}(\vec{E}; \vec{F})$ is continuous, and

$$Dg(x) = -\left(\frac{\partial f}{\partial y}(x, g(x))\right)^{-1} \circ \frac{\partial f}{\partial x}(x, g(x)),$$

for all $x \in A$.

The implicit function theorem plays an important role in the calculus of variations. We now consider another very important notion, that of a (local) diffeomorphism.

Definition 28.7. Given two topological spaces E and F , and an open subset A of E , we say that a function $f: A \rightarrow F$ is a *local homeomorphism from A to F* if for every $a \in A$, there is an open set $U \subseteq A$ containing a and an open set V containing $f(a)$ such that f is a homeomorphism from U to $V = f(U)$. If B is an open subset of F , we say that $f: A \rightarrow F$ is a *(global) homeomorphism from A to B* if f is a homeomorphism from A to $B = f(A)$. If E and F are normed affine spaces, we say that $f: A \rightarrow F$ is a *local diffeomorphism from A to F* if for every $a \in A$, there is an open set $U \subseteq A$ containing a and an open set V containing $f(a)$ such that f is a bijection from U to V , f is a C^1 -function on U , and f^{-1} is a C^1 -function on $V = f(U)$. We say that $f: A \rightarrow F$ is a *(global) diffeomorphism from A to B* if f is a homeomorphism from A to $B = f(A)$, f is a C^1 -function on A , and f^{-1} is a C^1 -function on B .

Note that a local diffeomorphism is a local homeomorphism. Also, as a consequence of Proposition 28.7, if f is a diffeomorphism on A , then $Df(a)$ is a linear isomorphism for every $a \in A$. The following theorem can be shown. In fact, there is a fairly simple proof using Theorem 28.13; see Schwartz [94], Lang [69] and Cartan [20].

Theorem 28.14. *Let E and F be complete normed affine spaces, let A be an open subset of E , and let $f: A \rightarrow F$ be a C^1 -function on A . The following properties hold:*

- (1) *For every $a \in A$, if $Df(a)$ is a linear isomorphism (which means that both $Df(a)$ and $(Df(a))^{-1}$ are linear and continuous),¹ then there exist some open subset $U \subseteq A$ containing a , and some open subset V of F containing $f(a)$, such that f is a diffeomorphism from U to $V = f(U)$. Furthermore,*

$$Df^{-1}(f(a)) = (Df(a))^{-1}.$$

For every neighborhood N of a , its image $f(N)$ is a neighborhood of $f(a)$, and for every open ball $U \subseteq A$ of center a , its image $f(U)$ contains some open ball of center $f(a)$.

¹Actually, since E and F are Banach spaces, by the Open Mapping Theorem, it is sufficient to assume that $Df(a)$ is continuous and bijective; see Lang [69].

- (2) If $Df(a)$ is invertible for every $a \in A$, then $B = f(A)$ is an open subset of F , and f is a local diffeomorphism from A to B . Furthermore, if f is injective, then f is a diffeomorphism from A to B .

Part (1) of Theorem 28.14 is often referred to as the “(local) inverse function theorem.” It plays an important role in the study of manifolds and (ordinary) differential equations.

If E and F are both of finite dimension, and some frames have been chosen, the invertibility of $Df(a)$ is equivalent to the fact that the Jacobian determinant $\det(J(f)(a))$ is nonnull. The case where $Df(a)$ is just injective or just surjective is also important for defining manifolds, using implicit definitions.

Definition 28.8. Let E and F be normed affine spaces, where E and F are of finite dimension (or both E and F are complete), and let A be an open subset of E . For any $a \in A$, a C^1 -function $f: A \rightarrow F$ is an *immersion at a* if $Df(a)$ is injective. A C^1 -function $f: A \rightarrow F$ is a *submersion at a* if $Df(a)$ is surjective. A C^1 -function $f: A \rightarrow F$ is an *immersion on A* (resp. a *submersion on A*) if $Df(a)$ is injective (resp. surjective) for every $a \in A$.

The following results can be shown.

Proposition 28.15. Let A be an open subset of \mathbb{R}^n , and let $f: A \rightarrow \mathbb{R}^m$ be a function. For every $a \in A$, $f: A \rightarrow \mathbb{R}^m$ is a submersion at a iff there exists an open subset U of A containing a , an open subset $W \subseteq \mathbb{R}^{n-m}$, and a diffeomorphism $\varphi: U \rightarrow f(U) \times W$, such that,

$$f = \pi_1 \circ \varphi,$$

where $\pi_1: f(U) \times W \rightarrow f(U)$ is the first projection. Equivalently,

$$(f \circ \varphi^{-1})(y_1, \dots, y_m, \dots, y_n) = (y_1, \dots, y_m).$$

$$\begin{array}{ccc} U \subseteq A & \xrightarrow{\varphi} & f(U) \times W \\ & \searrow f & \downarrow \pi_1 \\ & & f(U) \subseteq \mathbb{R}^m \end{array}$$

Furthermore, the image of every open subset of A under f is an open subset of F . (The same result holds for \mathbb{C}^n and \mathbb{C}^m).

Proposition 28.16. Let A be an open subset of \mathbb{R}^n , and let $f: A \rightarrow \mathbb{R}^m$ be a function. For every $a \in A$, $f: A \rightarrow \mathbb{R}^m$ is an immersion at a iff there exists an open subset U of A containing a , an open subset V containing $f(a)$ such that $f(U) \subseteq V$, an open subset W containing 0 such that $W \subseteq \mathbb{R}^{m-n}$, and a diffeomorphism $\varphi: V \rightarrow U \times W$, such that,

$$\varphi \circ f = in_1,$$

where $in_1: U \rightarrow U \times W$ is the injection map such that $in_1(u) = (u, 0)$, or equivalently,

$$(\varphi \circ f)(x_1, \dots, x_n) = (x_1, \dots, x_n, 0, \dots, 0).$$

$$\begin{array}{ccc} U \subseteq A & \xrightarrow{f} & f(U) \subseteq V \\ & \searrow in_1 & \downarrow \varphi \\ & & U \times W \end{array}$$

(The same result holds for \mathbb{C}^n and \mathbb{C}^m).

28.4 Tangent Spaces and Differentials

In this section, we discuss briefly a geometric interpretation of the notion of derivative. We consider sets of points defined by a differentiable function. This is a special case of the notion of a (differential) manifold.

Given two normed affine spaces E and F , let A be an open subset of E , and let $f: A \rightarrow F$ be a function.

Definition 28.9. Given $f: A \rightarrow F$ as above, its *graph* $\Gamma(f)$ is the set of all points

$$\Gamma(f) = \{(x, y) \in E \times F \mid x \in A, y = f(x)\}.$$

If Df is defined on A , we say that $\Gamma(f)$ is a *differential submanifold* of $E \times F$ of equation $y = f(x)$.

It should be noted that this is a very particular kind of differential manifold.

Example 28.3. If $E = \mathbb{R}$ and $F = \mathbb{R}^2$, letting $f = (g, h)$, where $g: \mathbb{R} \rightarrow \mathbb{R}$ and $h: \mathbb{R} \rightarrow \mathbb{R}$, $\Gamma(f)$ is a curve in \mathbb{R}^3 , of equations $y = g(x)$, $z = h(x)$. When $E = \mathbb{R}^2$ and $F = \mathbb{R}$, $\Gamma(f)$ is a surface in \mathbb{R}^3 , of equation $z = f(x, y)$.

We now define the notion of affine tangent space in a very general way. Next, we will see what it means for manifolds $\Gamma(f)$, as in Definition 28.9.

Definition 28.10. Given a normed affine space E , given any nonempty subset M of E , given any point $a \in M$, we say that a vector $u \in \overrightarrow{E}$ is *tangent at a to M* if there exist a sequence $(a_n)_{n \in \mathbb{N}}$ of points in M converging to a , and a sequence $(\lambda_n)_{n \in \mathbb{N}}$, with $\lambda_i \in \mathbb{R}$ and $\lambda_n \geq 0$, such that the sequence $(\lambda_n(a_n - a))_{n \in \mathbb{N}}$ converges to u .

The set of all vectors tangent at a to M is called the *family of tangent vectors at a to M* and the set of all points of E of the form $a + u$ where u belongs to the family of tangent vectors at a to M is called the *affine tangent family at a to M* .

Clearly, 0 is always tangent, and if u is tangent, then so is every λu , for $\lambda \in \mathbb{R}$, $\lambda \geq 0$. If $u \neq 0$, then the sequence $(\lambda_n)_{n \in \mathbb{N}}$ must tend towards $+\infty$. We have the following proposition.

Proposition 28.17. *Let E and F be two normed affine spaces, let A be an open subset of E , let $a \in A$, and let $f: A \rightarrow F$ be a function. If $Df(a)$ exists, then the family of tangent vectors at $(a, f(a))$ to Γ is a subspace $T_a(\Gamma)$ of $\vec{E} \times \vec{F}$, defined by the condition (equation)*

$$(u, v) \in T_a(\Gamma) \quad \text{iff} \quad v = Df(a)(u),$$

and the affine tangent family at $(a, f(a))$ to Γ is an affine variety $T_a(\Gamma)$ of $E \times F$, defined by the condition (equation)

$$(x, y) \in T_a(\Gamma) \quad \text{iff} \quad y = f(a) + Df(a)(x - a),$$

where Γ is the graph of f .

The proof is actually rather simple. We have $T_a(\Gamma) = a + T_a(\Gamma)$, and since $T_a(\Gamma)$ is a subspace of $\vec{E} \times \vec{F}$, the set $T_a(\Gamma)$ is an affine variety. Thus, the affine tangent space at a point $(a, f(a))$ is a familiar object, a line, a plane, etc.

As an illustration, when $E = \mathbb{R}^2$ and $F = \mathbb{R}$, the affine tangent plane at the point (a, b, c) to the surface of equation $z = f(x, y)$, is defined by the equation

$$z = c + \frac{\partial f}{\partial x}(a, b)(x - a) + \frac{\partial f}{\partial y}(a, b)(y - b).$$

If $E = \mathbb{R}$ and $F = \mathbb{R}^2$, the tangent line at (a, b, c) , to the curve of equations $y = g(x)$, $z = h(x)$, is defined by the equations

$$\begin{aligned} y &= b + Dg(a)(x - a), \\ z &= c + Dh(a)(x - a). \end{aligned}$$

Thus, derivatives and partial derivatives have the desired intended geometric interpretation as tangent spaces. Of course, in order to deal with this topic properly, we really would have to go deeper into the study of (differential) manifolds.

We now briefly consider second-order and higher-order derivatives.

28.5 Second-Order and Higher-Order Derivatives

Given two normed affine spaces E and F , and some open subset A of E , if $Df(a)$ is defined for every $a \in A$, then we have a mapping $Df: A \rightarrow \mathcal{L}(\vec{E}; \vec{F})$. Since $\mathcal{L}(\vec{E}; \vec{F})$ is a normed vector space, if Df exists on an open subset U of A containing a , we can consider taking the derivative of Df at some $a \in A$. If $D(Df)(a)$ exists for every $a \in A$, we get a mapping

$D^2f: A \rightarrow \mathcal{L}(\vec{E}; \mathcal{L}(\vec{E}; \vec{F}))$, where $D^2f(a) = D(Df)(a)$, for every $a \in A$. If $D^2f(a)$ exists, then for every $u \in \vec{E}$,

$$D^2f(a)(u) = D(Df)(a)(u) = D_u(Df)(a) \in \mathcal{L}(\vec{E}; \vec{F}).$$

Recall from Proposition 26.46, that the map app from $\mathcal{L}(\vec{E}; \vec{F}) \times \vec{E}$ to \vec{F} , defined such that for every $L \in \mathcal{L}(\vec{E}; \vec{F})$, for every $v \in \vec{E}$,

$$\text{app}(L, v) = L(v),$$

is a continuous bilinear map. Thus, in particular, given a fixed $v \in \vec{E}$, the linear map $\text{app}_v: \mathcal{L}(\vec{E}; \vec{F}) \rightarrow \vec{F}$, defined such that $\text{app}_v(L) = L(v)$, is a continuous map.

Also recall from Proposition 28.6, that if $h: A \rightarrow G$ is a function such that $Dh(a)$ exists, and $k: G \rightarrow H$ is a continuous linear map, then, $D(k \circ h)(a)$ exists, and

$$k(Dh(a)(u)) = D(k \circ h)(a)(u),$$

that is,

$$k(D_u h(a)) = D_u(k \circ h)(a),$$

Applying these two facts to $h = Df$, and to $k = \text{app}_v$, we have

$$D_u(Df)(a)(v) = D_u(\text{app}_v \circ Df)(a).$$

But $(\text{app}_v \circ Df)(x) = Df(x)(v) = D_v f(x)$, for every $x \in A$, that is, $\text{app}_v \circ Df = D_v f$ on A . So, we have

$$D_u(Df)(a)(v) = D_u(D_v f)(a),$$

and since $D^2f(a)(u) = D_u(Df)(a)$, we get

$$D^2f(a)(u)(v) = D_u(D_v f)(a).$$

Thus, when $D^2f(a)$ exists, $D_u(D_v f)(a)$ exists, and

$$D^2f(a)(u)(v) = D_u(D_v f)(a),$$

for all $u, v \in \vec{E}$. We also denote $D_u(D_v f)(a)$ by $D_{u,v}^2 f(a)$, or $D_u D_v f(a)$.

Recall from Proposition 26.45, that the map from $\mathcal{L}_2(\vec{E}, \vec{E}; \vec{F})$ to $\mathcal{L}(\vec{E}; \mathcal{L}(\vec{E}; \vec{F}))$ defined such that $g \mapsto \varphi$ iff for every $g \in \mathcal{L}_2(\vec{E}, \vec{E}; \vec{F})$,

$$\varphi(u)(v) = g(u, v),$$

is an isomorphism of vector spaces. Thus, we will consider $D^2f(a) \in \mathcal{L}(\vec{E}; \mathcal{L}(\vec{E}; \vec{F}))$ as a continuous bilinear map in $\mathcal{L}_2(\vec{E}, \vec{E}; \vec{F})$, and we will write $D^2f(a)(u, v)$, instead of $D^2f(a)(u)(v)$.

Then, the above discussion can be summarized by saying that when $D^2f(a)$ is defined, we have

$$D^2f(a)(u, v) = D_u D_v f(a).$$

When E has finite dimension and $(a_0, (e_1, \dots, e_n))$ is a frame for E , we denote $D_{e_j} D_{e_i} f(a)$ by $\frac{\partial^2 f}{\partial x_i \partial x_j}(a)$, when $i \neq j$, and we denote $D_{e_i} D_{e_i} f(a)$ by $\frac{\partial^2 f}{\partial x_i^2}(a)$.

The following important lemma attributed to Schwarz can be shown, using Lemma 28.11. Given a bilinear map $f: \vec{E} \times \vec{E} \rightarrow \vec{F}$, recall that f is *symmetric*, if

$$f(u, v) = f(v, u),$$

for all $u, v \in \vec{E}$.

Lemma 28.18. (*Schwarz's lemma*) *Given two normed affine spaces E and F , given any open subset A of E , given any $f: A \rightarrow F$, for every $a \in A$, if $D^2f(a)$ exists, then $D^2f(a) \in \mathcal{L}_2(\vec{E}, \vec{E}; \vec{F})$ is a continuous symmetric bilinear map. As a corollary, if E is of finite dimension n , and $(a_0, (e_1, \dots, e_n))$ is a frame for E , we have*

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(a) = \frac{\partial^2 f}{\partial x_j \partial x_i}(a).$$

Remark: There is a variation of the above lemma which does not assume the existence of $D^2f(a)$, but instead assumes that $D_u D_v f$ and $D_v D_u f$ exist on an open subset containing a and are continuous at a , and concludes that $D_u D_v f(a) = D_v D_u f(a)$. This is just a different result which does not imply Lemma 28.18, and is not a consequence of Lemma 28.18.



When $E = \mathbb{R}^2$, the only existence of $\frac{\partial^2 f}{\partial x \partial y}(a)$ and $\frac{\partial^2 f}{\partial y \partial x}(a)$ is not sufficient to insure the existence of $D^2f(a)$.

When E is of finite dimension n and $(a_0, (e_1, \dots, e_n))$ is a frame for E , if $D^2f(a)$ exists, for every $u = u_1 e_1 + \dots + u_n e_n$ and $v = v_1 e_1 + \dots + v_n e_n$ in \vec{E} , since $D^2f(a)$ is a symmetric bilinear form, we have

$$D^2f(a)(u, v) = \sum_{i=1, j=1}^n u_i v_j \frac{\partial^2 f}{\partial x_i \partial x_j}(a),$$

which can be written in matrix form as:

$$D^2f(a)(u, v) = U^T \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2}(a) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(a) & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(a) \\ \frac{\partial^2 f}{\partial x_1 \partial x_2}(a) & \frac{\partial^2 f}{\partial x_2^2}(a) & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n}(a) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n}(a) & \frac{\partial^2 f}{\partial x_2 \partial x_n}(a) & \cdots & \frac{\partial^2 f}{\partial x_n^2}(a) \end{pmatrix} V$$

where U is the column matrix representing u , and V is the column matrix representing v , over the frame $(a_0, (e_1, \dots, e_n))$.

The above symmetric matrix is called the *Hessian of f at a* . If F itself is of finite dimension, and $(b_0, (v_1, \dots, v_m))$ is a frame for F , then $f = (f_1, \dots, f_m)$, and each component $D^2 f(a)_i(u, v)$ of $D^2 f(a)(u, v)$ ($1 \leq i \leq m$), can be written as

$$D^2 f(a)_i(u, v) = U^\top \begin{pmatrix} \frac{\partial^2 f_i}{\partial x_1^2}(a) & \frac{\partial^2 f_i}{\partial x_1 \partial x_2}(a) & \dots & \frac{\partial^2 f_i}{\partial x_1 \partial x_n}(a) \\ \frac{\partial^2 f_i}{\partial x_1 \partial x_2}(a) & \frac{\partial^2 f_i}{\partial x_2^2}(a) & \dots & \frac{\partial^2 f_i}{\partial x_2 \partial x_n}(a) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f_i}{\partial x_1 \partial x_n}(a) & \frac{\partial^2 f_i}{\partial x_2 \partial x_n}(a) & \dots & \frac{\partial^2 f_i}{\partial x_n^2}(a) \end{pmatrix} V$$

Thus, we could describe the vector $D^2 f(a)(u, v)$ in terms of an $mn \times mn$ -matrix consisting of m diagonal blocks, which are the above Hessians, and the row matrix (U^\top, \dots, U^\top) (m times) and the column matrix consisting of m copies of V .

We now indicate briefly how higher-order derivatives are defined. Let $m \geq 2$. Given a function $f: A \rightarrow F$ as before, for any $a \in A$, if the derivatives $D^i f$ exist on A for all i , $1 \leq i \leq m-1$, by induction, $D^{m-1} f$ can be considered to be a continuous function $D^{m-1} f: A \rightarrow \mathcal{L}_{m-1}(\overrightarrow{E^{m-1}}; \overrightarrow{F})$, and we define

$$D^m f(a) = D(D^{m-1} f)(a).$$

Then, $D^m f(a)$ can be identified with a continuous m -multilinear map in $\mathcal{L}_m(\overrightarrow{E^m}; \overrightarrow{F})$. We can then show (as we did before), that if $D^m f(a)$ is defined, then

$$D^m f(a)(u_1, \dots, u_m) = D_{u_1} \dots D_{u_m} f(a).$$

When E is of finite dimension n and $(a_0, (e_1, \dots, e_n))$ is a frame for E , if $D^m f(a)$ exists, for every $j_1, \dots, j_m \in \{1, \dots, n\}$, we denote $D_{e_{j_m}} \dots D_{e_{j_1}} f(a)$ by

$$\frac{\partial^m f}{\partial x_{j_1} \dots \partial x_{j_m}}(a).$$

Given a m -multilinear map $f \in \mathcal{L}_m(\overrightarrow{E^m}; \overrightarrow{F})$, recall that f is *symmetric* if

$$f(u_{\pi(1)}, \dots, u_{\pi(m)}) = f(u_1, \dots, u_m),$$

for all $u_1, \dots, u_m \in \overrightarrow{E}$, and all permutations π on $\{1, \dots, m\}$. Then, the following generalization of Schwarz's lemma holds.

Lemma 28.19. *Given two normed affine spaces E and F , given any open subset A of E , given any $f: A \rightarrow F$, for every $a \in A$, for every $m \geq 1$, if $D^m f(a)$ exists, then $D^m f(a) \in \mathcal{L}_m(\overrightarrow{E^m}; \overrightarrow{F})$ is a continuous symmetric m -multilinear map. As a corollary, if E is of finite dimension n , and $(a_0, (e_1, \dots, e_n))$ is a frame for E , we have*

$$\frac{\partial^m f}{\partial x_{j_1} \dots \partial x_{j_m}}(a) = \frac{\partial^m f}{\partial x_{\pi(j_1)} \dots \partial x_{\pi(j_m)}}(a),$$

for every $j_1, \dots, j_m \in \{1, \dots, n\}$, and for every permutation π on $\{1, \dots, m\}$.

If E is of finite dimension n , and $(a_0, (e_1, \dots, e_n))$ is a frame for E , $D^m f(a)$ is a symmetric m -multilinear map, and we have

$$D^m f(a)(u_1, \dots, u_m) = \sum_j u_{1,j_1} \dots u_{m,j_m} \frac{\partial^m f}{\partial x_{j_1} \dots \partial x_{j_m}}(a),$$

where j ranges over all functions $j: \{1, \dots, m\} \rightarrow \{1, \dots, n\}$, for any m vectors

$$u_j = u_{j,1}e_1 + \dots + u_{j,n}e_n.$$

The concept of C^1 -function is generalized to the concept of C^m -function, and Theorem 28.12 can also be generalized.

Definition 28.11. Given two normed affine spaces E and F , and an open subset A of E , for any $m \geq 1$, we say that a function $f: A \rightarrow F$ is of class C^m on A or a C^m -function on A if $D^k f$ exists and is continuous on A for every k , $1 \leq k \leq m$. We say that $f: A \rightarrow F$ is of class C^∞ on A or a C^∞ -function on A if $D^k f$ exists and is continuous on A for every $k \geq 1$. A C^∞ -function (on A) is also called a *smooth function* (on A). A C^m -diffeomorphism $f: A \rightarrow B$ between A and B (where A is an open subset of E and B is an open subset of F) is a bijection between A and $B = f(A)$, such that both $f: A \rightarrow B$ and its inverse $f^{-1}: B \rightarrow A$ are C^m -functions.

Equivalently, f is a C^m -function on A if f is a C^1 -function on A and Df is a C^{m-1} -function on A .

We have the following theorem giving a necessary and sufficient condition for f to a C^m -function on A . A generalization to the case where $E = (E_1, a_1) \oplus \dots \oplus (E_n, a_n)$ also holds.

Theorem 28.20. *Given two normed affine spaces E and F , where E is of finite dimension n , and where $(a_0, (u_1, \dots, u_n))$ is a frame of E , given any open subset A of E , given any function $f: A \rightarrow F$, for any $m \geq 1$, the derivative $D^m f$ is a C^m -function on A iff every partial derivative $D_{u_{j_k}} \dots D_{u_{j_1}} f$ (or $\frac{\partial^k f}{\partial x_{j_1} \dots \partial x_{j_k}}(a)$) is defined and continuous on A , for all*

k , $1 \leq k \leq m$, and all $j_1, \dots, j_k \in \{1, \dots, n\}$. As a corollary, if F is of finite dimension p , and $(b_0, (v_1, \dots, v_p))$ is a frame of F , the derivative $D^m f$ is defined and continuous on A iff every partial derivative $D_{u_{j_k}} \dots D_{u_{j_1}} f_i$ (or $\frac{\partial^k f_i}{\partial x_{j_1} \dots \partial x_{j_k}}(a)$) is defined and continuous on A , for all k , $1 \leq k \leq m$, for all i , $1 \leq i \leq p$, and all $j_1, \dots, j_k \in \{1, \dots, n\}$.

When $E = \mathbb{R}$ (or $E = \mathbb{C}$), for any $a \in E$, $D^m f(a)(1, \dots, 1)$ is a vector in \vec{F} , called the m th-order vector derivative. As in the case $m = 1$, we will usually identify the multilinear map $D^m f(a)$ with the vector $D^m f(a)(1, \dots, 1)$. Some notational conventions can also be introduced to simplify the notation of higher-order derivatives, and we discuss such conventions very briefly.

Recall that when E is of finite dimension n , and $(a_0, (e_1, \dots, e_n))$ is a frame for E , $D^m f(a)$ is a symmetric m -multilinear map, and we have

$$D^m f(a)(u_1, \dots, u_m) = \sum_j u_{1,j_1} \dots u_{m,j_m} \frac{\partial^m f}{\partial x_{j_1} \dots \partial x_{j_m}}(a),$$

where j ranges over all functions $j: \{1, \dots, m\} \rightarrow \{1, \dots, n\}$, for any m vectors

$$u_j = u_{j,1}e_1 + \dots + u_{j,n}e_n.$$

We can then group the various occurrences of ∂x_{j_k} corresponding to the same variable x_{j_k} , and this leads to the notation

$$\left(\frac{\partial}{\partial x_1}\right)^{\alpha_1} \left(\frac{\partial}{\partial x_2}\right)^{\alpha_2} \dots \left(\frac{\partial}{\partial x_n}\right)^{\alpha_n} f(a),$$

where $\alpha_1 + \alpha_2 + \dots + \alpha_n = m$.

If we denote $(\alpha_1, \dots, \alpha_n)$ simply by α , then we denote

$$\left(\frac{\partial}{\partial x_1}\right)^{\alpha_1} \left(\frac{\partial}{\partial x_2}\right)^{\alpha_2} \dots \left(\frac{\partial}{\partial x_n}\right)^{\alpha_n} f$$

by

$$\partial^\alpha f, \quad \text{or} \quad \left(\frac{\partial}{\partial x}\right)^\alpha f.$$

If $\alpha = (\alpha_1, \dots, \alpha_n)$, we let $|\alpha| = \alpha_1 + \alpha_2 + \dots + \alpha_n$, $\alpha! = \alpha_1! \dots \alpha_n!$, and if $h = (h_1, \dots, h_n)$, we denote $h_1^{\alpha_1} \dots h_n^{\alpha_n}$ by h^α .

In the next section, we survey various versions of Taylor's formula.

28.6 Taylor's formula, Faà di Bruno's formula

We discuss, without proofs, several versions of Taylor's formula. The hypotheses required in each version become increasingly stronger. The first version can be viewed as a generalization of the notion of derivative. Given an m -linear map $f: \overrightarrow{E^m} \rightarrow \overrightarrow{F}$, for any vector $h \in \overrightarrow{E}$, we abbreviate

$$f(\underbrace{h, \dots, h}_m)$$

by $f(h^m)$. The version of Taylor's formula given next is sometimes referred to as the *formula of Taylor–Young*.

Theorem 28.21. (*Taylor–Young*) *Given two normed affine spaces E and F , for any open subset $A \subseteq E$, for any function $f: A \rightarrow F$, for any $a \in A$, if $D^k f$ exists in A for all k , $1 \leq k \leq m-1$, and if $D^m f(a)$ exists, then we have:*

$$f(a+h) = f(a) + \frac{1}{1!}D^1 f(a)(h) + \cdots + \frac{1}{m!}D^m f(a)(h^m) + \|h\|^m \epsilon(h),$$

for any h such that $a+h \in A$, and where $\lim_{h \rightarrow 0, h \neq 0} \epsilon(h) = 0$.

The above version of Taylor's formula has applications to the study of relative maxima (or minima) of real-valued functions. It is also used to study the local properties of curves and surfaces.

The next version of Taylor's formula can be viewed as a generalization of Lemma 28.11. It is sometimes called the *Taylor formula with Lagrange remainder* or *generalized mean value theorem*.

Theorem 28.22. (*Generalized mean value theorem*) *Let E and F be two normed affine spaces, let A be an open subset of E , and let $f: A \rightarrow F$ be a function on A . Given any $a \in A$ and any $h \neq 0$ in \overrightarrow{E} , if the closed segment $[a, a+h]$ is contained in A , $D^k f$ exists in A for all k , $1 \leq k \leq m$, $D^{m+1} f(x)$ exists at every point x of the open segment $]a, a+h[$, and*

$$\max_{x \in]a, a+h[} \|D^{m+1} f(x)\| \leq M,$$

for some $M \geq 0$, then

$$\left\| f(a+h) - f(a) - \left(\frac{1}{1!}D^1 f(a)(h) + \cdots + \frac{1}{m!}D^m f(a)(h^m) \right) \right\| \leq M \frac{\|h\|^{m+1}}{(m+1)!}.$$

As a corollary, if $L: \overrightarrow{E^{m+1}} \rightarrow \overrightarrow{F}$ is a continuous $(m+1)$ -linear map, then

$$\left\| f(a+h) - f(a) - \left(\frac{1}{1!}D^1 f(a)(h) + \cdots + \frac{1}{m!}D^m f(a)(h^m) + \frac{L(h^{m+1})}{(m+1)!} \right) \right\| \leq M \frac{\|h\|^{m+1}}{(m+1)!},$$

where $M = \max_{x \in]a, a+h[} \|D^{m+1} f(x) - L\|$.

The above theorem is sometimes stated under the slightly stronger assumption that f is a C^m -function on A . If $f: A \rightarrow \mathbb{R}$ is a real-valued function, Theorem 28.22 can be refined a little bit. This version is often called the *formula of Taylor–MacLaurin*.

Theorem 28.23. (*Taylor–MacLaurin*) Let E be a normed affine space, let A be an open subset of E , and let $f: A \rightarrow \mathbb{R}$ be a real-valued function on A . Given any $a \in A$ and any $h \neq 0$ in \vec{E} , if the closed segment $[a, a + h]$ is contained in A , if $D^k f$ exists in A for all k , $1 \leq k \leq m$, and $D^{m+1}f(x)$ exists at every point x of the open segment $]a, a + h[$, then there is some $\theta \in \mathbb{R}$, with $0 < \theta < 1$, such that

$$f(a + h) = f(a) + \frac{1}{1!}D^1 f(a)(h) + \cdots + \frac{1}{m!}D^m f(a)(h^m) + \frac{1}{(m+1)!}D^{m+1}f(a + \theta h)(h^{m+1}).$$

We also mention for “mathematical culture,” a version with integral remainder, in the case of a real-valued function. This is usually called *Taylor’s formula with integral remainder*.

Theorem 28.24. (*Taylor’s formula with integral remainder*) Let E be a normed affine space, let A be an open subset of E , and let $f: A \rightarrow \mathbb{R}$ be a real-valued function on A . Given any $a \in A$ and any $h \neq 0$ in \vec{E} , if the closed segment $[a, a + h]$ is contained in A , and if f is a C^{m+1} -function on A , then we have

$$f(a + h) = f(a) + \frac{1}{1!}D^1 f(a)(h) + \cdots + \frac{1}{m!}D^m f(a)(h^m) + \int_0^1 \frac{(1-t)^m}{m!} [D^{m+1}f(a + th)(h^{m+1})] dt.$$

The advantage of the above formula is that it gives an explicit remainder. We now examine briefly the situation where E is of finite dimension n , and $(a_0, (e_1, \dots, e_n))$ is a frame for E . In this case, we get a more explicit expression for the expression

$$\sum_{i=0}^{k=m} \frac{1}{k!} D^k f(a)(h^k)$$

involved in all versions of Taylor’s formula, where by convention, $D^0 f(a)(h^0) = f(a)$. If $h = h_1 e_1 + \cdots + h_n e_n$, then we have

$$\sum_{k=0}^{k=m} \frac{1}{k!} D^k f(a)(h^k) = \sum_{k_1 + \cdots + k_n \leq m} \frac{h_1^{k_1} \cdots h_n^{k_n}}{k_1! \cdots k_n!} \left(\frac{\partial}{\partial x_1} \right)^{k_1} \cdots \left(\frac{\partial}{\partial x_n} \right)^{k_n} f(a),$$

which, using the abbreviated notation introduced at the end of Section 28.5, can also be written as

$$\sum_{k=0}^{k=m} \frac{1}{k!} D^k f(a)(h^k) = \sum_{|\alpha| \leq m} \frac{h^\alpha}{\alpha!} \partial^\alpha f(a).$$

The advantage of the above notation is that it is the same as the notation used when $n = 1$, i.e., when $E = \mathbb{R}$ (or $E = \mathbb{C}$). Indeed, in this case, the Taylor–MacLaurin formula reads as:

$$f(a + h) = f(a) + \frac{h}{1!}D^1f(a) + \cdots + \frac{h^m}{m!}D^mf(a) + \frac{h^{m+1}}{(m+1)!}D^{m+1}f(a + \theta h),$$

for some $\theta \in \mathbb{R}$, with $0 < \theta < 1$, where $D^k f(a)$ is the value of the k -th derivative of f at a (and thus, as we have already said several times, this is the k th-order vector derivative, which is just a scalar, since $F = \mathbb{R}$).

In the above formula, the assumptions are that $f: [a, a + h] \rightarrow \mathbb{R}$ is a C^m -function on $[a, a + h]$, and that $D^{m+1}f(x)$ exists for every $x \in]a, a + h[$.

Taylor's formula is useful to study the local properties of curves and surfaces. In the case of a curve, we consider a function $f: [r, s] \rightarrow F$ from a closed interval $[r, s]$ of \mathbb{R} to some affine space F , the derivatives $D^k f(a)(h^k)$ correspond to vectors $h^k D^k f(a)$, where $D^k f(a)$ is the k th vector derivative of f at a (which is really $D^k f(a)(1, \dots, 1)$), and for any $a \in]r, s[$, Theorem 28.21 yields the following formula:

$$f(a + h) = f(a) + \frac{h}{1!}D^1f(a) + \cdots + \frac{h^m}{m!}D^mf(a) + h^m\epsilon(h),$$

for any h such that $a + h \in]r, s[$, and where $\lim_{h \rightarrow 0, h \neq 0} \epsilon(h) = 0$.

In the case of functions $f: \mathbb{R}^n \rightarrow \mathbb{R}$, it is convenient to have formulae for the Taylor–Young formula and the Taylor–MacLaurin formula in terms of the gradient and the Hessian. Recall that the *gradient* $\nabla f(a)$ of f at $a \in \mathbb{R}^n$ is the column vector

$$\nabla f(a) \begin{pmatrix} \frac{\partial f}{\partial x_1}(a) \\ \frac{\partial f}{\partial x_2}(a) \\ \vdots \\ \frac{\partial f}{\partial x_n}(a) \end{pmatrix},$$

and that

$$f'(a)(u) = Df(a)(u) = \nabla f(a) \cdot u,$$

for any $u \in \mathbb{R}^n$ (where \cdot means inner product). The *Hessian matrix* $\nabla^2 f(a)$ of f at $a \in \mathbb{R}^n$

is the $n \times n$ symmetric matrix

$$\nabla^2 f(a) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2}(a) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(a) & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(a) \\ \frac{\partial^2 f}{\partial x_1 \partial x_2}(a) & \frac{\partial^2 f}{\partial x_2^2}(a) & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n}(a) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n}(a) & \frac{\partial^2 f}{\partial x_2 \partial x_n}(a) & \cdots & \frac{\partial^2 f}{\partial x_n^2}(a) \end{pmatrix},$$

and we have

$$D^2 f(a)(u, v) = u^\top \nabla^2 f(a) v = u \cdot \nabla^2 f(a) v = \nabla^2 f(a) u \cdot v,$$

for all $u, v \in \mathbb{R}^n$. Then, we have the following three formulations of the formula of Taylor–Young of order 2:

$$\begin{aligned} f(a+h) &= f(a) + Df(a)(h) + \frac{1}{2} D^2 f(a)(h, h) + \|h\|^2 \epsilon(h) \\ f(a+h) &= f(a) + \nabla f(a) \cdot h + \frac{1}{2} (h \cdot \nabla^2 f(a) h) + (h \cdot h) \epsilon(h) \\ f(a+h) &= f(a) + (\nabla f(a))^\top h + \frac{1}{2} (h^\top \nabla^2 f(a) h) + (h^\top h) \epsilon(h). \end{aligned}$$

with $\lim_{h \rightarrow 0} \epsilon(h) = 0$.

One should keep in mind that only the first formula is intrinsic (i.e., does not depend on the choice of a basis), whereas the other two depend on the basis and the inner product chosen on \mathbb{R}^n . As an exercise, the reader should write similar formulae for the Taylor–MacLaurin formula of order 2.

Another application of Taylor’s formula is the derivation of a formula which gives the m -th derivative of the composition of two functions, usually known as “Faà di Bruno’s formula.” This formula is useful when dealing with geometric continuity of splines curves and surfaces.

Proposition 28.25. *Given any normed affine space E , for any function $f: \mathbb{R} \rightarrow \mathbb{R}$ and any function $g: \mathbb{R} \rightarrow E$, for any $a \in \mathbb{R}$, letting $b = f(a)$, $f^{(i)}(a) = D^i f(a)$, and $g^{(i)}(b) = D^i g(b)$, for any $m \geq 1$, if $f^{(i)}(a)$ and $g^{(i)}(b)$ exist for all i , $1 \leq i \leq m$, then $(g \circ f)^{(m)}(a) = D^m(g \circ f)(a)$ exists and is given by the following formula:*

$$(g \circ f)^{(m)}(a) = \sum_{0 \leq j \leq m} \sum_{\substack{i_1 + i_2 + \cdots + i_m = j \\ i_1 + 2i_2 + \cdots + mi_m = m \\ i_1, i_2, \dots, i_m \geq 0}} \frac{m!}{i_1! \cdots i_m!} g^{(j)}(b) \left(\frac{f^{(1)}(a)}{1!} \right)^{i_1} \cdots \left(\frac{f^{(m)}(a)}{m!} \right)^{i_m}.$$

When $m = 1$, the above simplifies to the familiar formula

$$(g \circ f)'(a) = g'(b) f'(a),$$

and for $m = 2$, we have

$$(g \circ f)^{(2)}(a) = g^{(2)}(b) (f^{(1)}(a))^2 + g^{(1)}(b) f^{(2)}(a).$$

28.7 Vector Fields, Covariant Derivatives, Lie Brackets

In this section, we briefly consider vector fields and covariant derivatives of vector fields. Such derivatives play an important role in continuous mechanics. Given a normed affine space (E, \vec{E}) , a *vector field over (E, \vec{E})* is a function $X: E \rightarrow \vec{E}$. Intuitively, a vector field assigns a vector to every point in E . Such vectors could be forces, velocities, accelerations, etc.

Given two vector fields X, Y defined on some open subset Ω of E , for every point $a \in \Omega$, we would like to define the derivative of X with respect to Y at a . This is a type of directional derivative that gives the variation of X as we move along Y , and we denote it by $D_Y X(a)$. The derivative $D_Y X(a)$ is defined as follows.

Definition 28.12. Let (E, \vec{E}) be a normed affine space. Given any open subset Ω of E , given any two vector fields X and Y defined over Ω , for any $a \in \Omega$, the *covariant derivative (or Lie derivative) of X w.r.t. the vector field Y at a* , denoted by $D_Y X(a)$, is the limit (if it exists)

$$\lim_{t \rightarrow 0, t \in U} \frac{X(a + tY(a)) - X(a)}{t},$$

where $U = \{t \in \mathbb{R} \mid a + tY(a) \in \Omega, t \neq 0\}$.

If Y is a constant vector field, it is immediately verified that the map

$$X \mapsto D_Y X(a)$$

is a linear map called the *derivative* of the vector field X , and denoted by $DX(a)$. If $f: E \rightarrow \mathbb{R}$ is a function, we define $D_Y f(a)$ as the limit (if it exists)

$$\lim_{t \rightarrow 0, t \in U} \frac{f(a + tY(a)) - f(a)}{t},$$

where $U = \{t \in \mathbb{R} \mid a + tY(a) \in \Omega, t \neq 0\}$. It is the *directional derivative of f w.r.t. the vector field Y at a* , and it is also often denoted by $Y(f)(a)$, or $Y(f)_a$.

From now on, we assume that all the vector fields and all the functions under consideration are smooth (C^∞). The set $C^\infty(\Omega)$ of smooth C^∞ -functions $f: \Omega \rightarrow \mathbb{R}$ is a ring. Given a smooth vector field X and a smooth function f (both over Ω), the vector field fX is defined such that $(fX)(a) = f(a)X(a)$, and it is immediately verified that it is smooth. Thus, the set $\mathcal{X}(\Omega)$ of smooth vector fields over Ω is a $C^\infty(\Omega)$ -module.

The following proposition is left as an exercise. It shows that $D_Y X(a)$ is a \mathbb{R} -bilinear map on $\mathcal{X}(\Omega)$, is $C^\infty(\Omega)$ -linear in Y , and satisfies the Leibniz derivation rules with respect to X .

Proposition 28.26. *The covariant derivative $D_Y X(a)$ satisfies the following properties:*

$$\begin{aligned} D_{(Y_1+Y_2)}X(a) &= D_{Y_1}X(a) + D_{Y_2}X(a), \\ D_{fY}X(a) &= f(a)D_YX(a), \\ D_Y(X_1 + X_2)(a) &= D_YX_1(a) + D_YX_2(a), \\ D_YfX(a) &= D_Yf(a)X(a) + f(a)D_YX(a), \end{aligned}$$

where X, Y, X_1, X_2, Y_1, Y_2 are smooth vector fields over Ω , and $f: E \rightarrow \mathbb{R}$ is a smooth function.

In differential geometry, the above properties are taken as the axioms of *affine connections*, in order to define covariant derivatives of vector fields over manifolds. In many cases, the vector field Y is the tangent field of some smooth curve $\gamma:]-\eta, \eta[\rightarrow E$. If so, the following proposition holds.

Proposition 28.27. *Given a smooth curve $\gamma:]-\eta, \eta[\rightarrow E$, letting Y be the vector field defined on $\gamma(]-\eta, \eta[)$ such that*

$$Y(\gamma(u)) = \frac{d\gamma}{dt}(u),$$

for any vector field X defined on $\gamma(]-\eta, \eta[)$, we have

$$D_Y X(a) = \frac{d}{dt} \left[X(\gamma(t)) \right] (0),$$

where $a = \gamma(0)$.

The derivative $D_Y X(a)$ is thus the derivative of the vector field X along the curve γ , and it is called the *covariant derivative of X along γ* .

Given an affine frame $(O, (u_1, \dots, u_n))$ for (E, \vec{E}) , it is easily seen that the covariant derivative $D_Y X(a)$ is expressed as follows:

$$D_Y X(a) = \sum_{i=1}^n \sum_{j=1}^n \left(Y_j \frac{\partial X_i}{\partial x_j} \right) (a) e_i.$$

Generally, $D_Y X(a) \neq D_X Y(a)$. The quantity

$$[X, Y] = D_X Y - D_Y X$$

is called the *Lie bracket* of the vector fields X and Y . The Lie bracket plays an important role in differential geometry. In terms of coordinates,

$$[X, Y] = \sum_{i=1}^n \sum_{j=1}^n \left(X_j \frac{\partial Y_i}{\partial x_j} - Y_j \frac{\partial X_i}{\partial x_j} \right) e_i.$$

28.8 Further Readings

A thorough treatment of differential calculus can be found in Munkres [85], Lang [70], Schwartz [94], Cartan [20], and Avez [6]. The techniques of differential calculus have many applications, especially to the geometry of curves and surfaces and to differential geometry in general. For this, we recommend do Carmo [30, 31] (two beautiful classics on the subject), Kreyszig [65], Stoker [102], Gray [50], Berger and Gostiaux [10], Milnor [81], Lang [68], Warner [114] and Choquet-Bruhat [23].

Chapter 29

Extrema of Real-Valued Functions

29.1 Local Extrema, Constrained Local Extrema, and Lagrange Multipliers

Let $J: E \rightarrow \mathbb{R}$ be a real-valued function defined on a normed vector space E (or more generally, any topological space). Ideally we would like to find where the function J reaches a minimum or a maximum value, at least locally. In this chapter, we will usually use the notations $dJ(u)$ or $J'(u)$ (or dJ_u or J'_u) for the derivative of J at u , instead of $DJ(u)$. Our presentation follows very closely that of Ciarlet [24] (Chapter 7), which we find to be one of the clearest.

Definition 29.1. If $J: E \rightarrow \mathbb{R}$ is a real-valued function defined on a normed vector space E , we say that J has a *local minimum* (or *relative minimum*) at the point $u \in E$ if there is some open subset $W \subseteq E$ containing u such that

$$J(u) \leq J(w) \quad \text{for all } w \in W.$$

Similarly, we say that J has a *local maximum* (or *relative maximum*) at the point $u \in E$ if there is some open subset $W \subseteq E$ containing u such that

$$J(u) \geq J(w) \quad \text{for all } w \in W.$$

In either case, we say that J has a *local extremum* (or *relative extremum*) at u . We say that J has a *strict local minimum* (resp. *strict local maximum*) at the point $u \in E$ if there is some open subset $W \subseteq E$ containing u such that

$$J(u) < J(w) \quad \text{for all } w \in W - \{u\}$$

(resp.

$$J(u) > J(w) \quad \text{for all } w \in W - \{u\}).$$

By abuse of language, we often say that the point u itself “is a local minimum” or a “local maximum,” even though, strictly speaking, this does not make sense.

We begin with a well-known necessary condition for a local extremum.

Proposition 29.1. *Let E be a normed vector space and let $J: \Omega \rightarrow \mathbb{R}$ be a function, with Ω some open subset of E . If the function J has a local extremum at some point $u \in \Omega$ and if J is differentiable at u , then*

$$dJ(u) = J'(u) = 0.$$

Proof. Pick any $v \in E$. Since Ω is open, for t small enough we have $u + tv \in \Omega$, so there is an open interval $I \subseteq \mathbb{R}$ such that the function φ given by

$$\varphi(t) = J(u + tv)$$

for all $t \in I$ is well-defined. By applying the chain rule, we see that φ is differentiable at $t = 0$, and we get

$$\varphi'(0) = dJ_u(v).$$

Without loss of generality, assume that u is a local minimum. Then we have

$$\varphi'(0) = \lim_{t \rightarrow 0_-} \frac{\varphi(t) - \varphi(0)}{t} \leq 0$$

and

$$\varphi'(0) = \lim_{t \rightarrow 0_+} \frac{\varphi(t) - \varphi(0)}{t} \geq 0,$$

which shows that $\varphi'(0) = dJ_u(v) = 0$. As $v \in E$ is arbitrary, we conclude that $dJ_u = 0$. \square

A point $u \in \Omega$ such that $J(u) = 0$ is called a *critical point* of J .

It is important to note that the fact that Ω is open is crucial. For example, if J is the identity function on $[0, 1]$, then $dJ(x) = 1$ for all $x \in [0, 1]$, even though J has a minimum at $x = 0$ and a maximum at $x = 1$. Also, if $E = \mathbb{R}^n$, then the condition $dJ(u) = 0$ is equivalent to the system

$$\begin{aligned} \frac{\partial J}{\partial x_1}(u_1, \dots, u_n) &= 0 \\ &\vdots \\ \frac{\partial J}{\partial x_n}(u_1, \dots, u_n) &= 0. \end{aligned}$$

In many practical situations, we need to look for local extrema of a function J under additional constraints. This situation can be formalized conveniently as follows: We have a function $J: \Omega \rightarrow \mathbb{R}$ defined on some open subset Ω of a normed vector space, but we also have some subset U of Ω and we are looking for the local extrema of J with respect to the set U . Note that in most cases, U is not open. In fact, U is usually closed.

Definition 29.2. If $J: \Omega \rightarrow \mathbb{R}$ is a real-valued function defined on some open subset Ω of a normed vector space E and if U is some subset of Ω , we say that J has a *local minimum* (or *relative minimum*) at the point $u \in U$ with respect to U if there is some open subset $W \subseteq \Omega$ containing u such that

$$J(u) \leq J(w) \quad \text{for all } w \in U \cap W.$$

Similarly, we say that J has a *local maximum* (or *relative maximum*) at the point $u \in U$ with respect to U if there is some open subset $W \subseteq \Omega$ containing u such that

$$J(u) \geq J(w) \quad \text{for all } w \in U \cap W.$$

In either case, we say that J has a *local extremum* at u with respect to U .

We will be particularly interested in the case where $\Omega \subseteq E_1 \times E_2$ is an open subset of a product of normed vector spaces and where U is the zero locus of some continuous function $\varphi: \Omega \rightarrow E_2$, which means that

$$U = \{(u_1, u_2) \in \Omega \mid \varphi(u_1, u_2) = 0\}.$$

For the sake of brevity, we say that J has a *constrained local extremum* at u instead of saying that J has a *local extremum* at the point $u \in U$ with respect to U . Fortunately, there is a necessary condition for constrained local extrema in terms of *Lagrange multipliers*.

Theorem 29.2. (*Necessary condition for a constrained extremum*) Let $\Omega \subseteq E_1 \times E_2$ be an open subset of a product of normed vector spaces, with E_1 a Banach space (E_1 is complete), let $\varphi: \Omega \rightarrow E_2$ be a C^1 -function (which means that $d\varphi(\omega)$ exists and is continuous for all $\omega \in \Omega$), and let

$$U = \{(u_1, u_2) \in \Omega \mid \varphi(u_1, u_2) = 0\}.$$

Moreover, let $u = (u_1, u_2) \in U$ be a point such that

$$\frac{\partial \varphi}{\partial x_2}(u_1, u_2) \in \mathcal{L}(E_2; E_2) \quad \text{and} \quad \left(\frac{\partial \varphi}{\partial x_2}(u_1, u_2) \right)^{-1} \in \mathcal{L}(E_2; E_2),$$

and let $J: \Omega \rightarrow \mathbb{R}$ be a function which is differentiable at u . If J has a constrained local extremum at u , then there is a continuous linear form $\Lambda(u) \in \mathcal{L}(E_2; \mathbb{R})$ such that

$$dJ(u) + \Lambda(u) \circ d\varphi(u) = 0.$$

Proof. The plan of attack is to use the implicit function theorem; Theorem 28.13. Observe that the assumptions of Theorem 28.13 are indeed met. Therefore, there exist some open subsets $U_1 \subseteq E_1$, $U_2 \subseteq E_2$, and a continuous function $g: U_1 \rightarrow U_2$ with $(u_1, u_2) \in U_1 \times U_2 \subseteq \Omega$ and such that

$$\varphi(v_1, g(v_1)) = 0$$

for all $v_1 \in U_1$. Moreover, g is differentiable at $u_1 \in U_1$ and

$$dg(u_1) = -\left(\frac{\partial\varphi}{\partial x_2}(u)\right)^{-1} \circ \frac{\partial\varphi}{\partial x_1}(u).$$

It follows that the restriction of J to $(U_1 \times U_2) \cap U$ yields a function G of a single variable, with

$$G(v_1) = J(v_1, g(v_1))$$

for all $v_1 \in U_1$. Now, the function G is differentiable at u_1 and it has a local extremum at u_1 on U_1 , so Proposition 29.1 implies that

$$dG(u_1) = 0.$$

By the chain rule,

$$\begin{aligned} dG(u_1) &= \frac{\partial J}{\partial x_1}(u) + \frac{\partial J}{\partial x_2}(u) \circ dg(u_1) \\ &= \frac{\partial J}{\partial x_1}(u) - \frac{\partial J}{\partial x_2}(u) \circ \left(\frac{\partial\varphi}{\partial x_2}(u)\right)^{-1} \circ \frac{\partial\varphi}{\partial x_1}(u). \end{aligned}$$

From $dG(u_1) = 0$, we deduce

$$\frac{\partial J}{\partial x_1}(u) = \frac{\partial J}{\partial x_2}(u) \circ \left(\frac{\partial\varphi}{\partial x_2}(u)\right)^{-1} \circ \frac{\partial\varphi}{\partial x_1}(u),$$

and since we also have

$$\frac{\partial J}{\partial x_2}(u) = \frac{\partial J}{\partial x_2}(u) \circ \left(\frac{\partial\varphi}{\partial x_2}(u)\right)^{-1} \circ \frac{\partial\varphi}{\partial x_2}(u),$$

if we let

$$\Lambda(u) = -\frac{\partial J}{\partial x_2}(u) \circ \left(\frac{\partial\varphi}{\partial x_2}(u)\right)^{-1},$$

then we get

$$\begin{aligned} dJ(u) &= \frac{\partial J}{\partial x_1}(u) + \frac{\partial J}{\partial x_2}(u) \\ &= \frac{\partial J}{\partial x_2}(u) \circ \left(\frac{\partial\varphi}{\partial x_2}(u)\right)^{-1} \circ \left(\frac{\partial\varphi}{\partial x_1}(u) + \frac{\partial\varphi}{\partial x_2}(u)\right) \\ &= -\Lambda(u) \circ d\varphi(u), \end{aligned}$$

which yields $dJ(u) + \Lambda(u) \circ d\varphi(u) = 0$, as claimed. □

In most applications, we have $E_1 = \mathbb{R}^{n-m}$ and $E_2 = \mathbb{R}^m$ for some integers m, n such that $1 \leq m < n$, Ω is an open subset of \mathbb{R}^n , $J: \Omega \rightarrow \mathbb{R}$, and we have m functions $\varphi_i: \Omega \rightarrow \mathbb{R}$ defining the subset

$$U = \{v \in \Omega \mid \varphi_i(v) = 0, 1 \leq i \leq m\}.$$

Theorem 29.2 yields the following necessary condition:

Theorem 29.3. *(Necessary condition for a constrained extremum in terms of Lagrange multipliers) Let Ω be an open subset of \mathbb{R}^n , consider m C^1 -functions $\varphi_i: \Omega \rightarrow \mathbb{R}$ (with $1 \leq m < n$), let*

$$U = \{v \in \Omega \mid \varphi_i(v) = 0, 1 \leq i \leq m\},$$

and let $u \in U$ be a point such that the derivatives $d\varphi_i(u) \in \mathcal{L}(\mathbb{R}^n; \mathbb{R})$ are linearly independent; equivalently, assume that the $m \times n$ matrix $((\partial\varphi_i/\partial x_j)(u))$ has rank m . If $J: \Omega \rightarrow \mathbb{R}$ is a function which is differentiable at $u \in U$ and if J has a local constrained extremum at u , then there exist m numbers $\lambda_i(u) \in \mathbb{R}$, uniquely defined, such that

$$dJ(u) + \lambda_1(u)d\varphi_1(u) + \cdots + \lambda_m(u)d\varphi_m(u) = 0;$$

equivalently,

$$\nabla J(u) + \lambda_1(u)\nabla\varphi_1(u) + \cdots + \lambda_m(u)\nabla\varphi_m(u) = 0.$$

Proof. The linear independence of the m linear forms $d\varphi_i(u)$ is equivalent to the fact that the $m \times n$ matrix $A = ((\partial\varphi_i/\partial x_j)(u))$ has rank m . By reordering the columns, we may assume that the first m columns are linearly independent. If we let $\varphi: \Omega \rightarrow \mathbb{R}^m$ be the function defined by

$$\varphi(v) = (\varphi_1(v), \dots, \varphi_m(v))$$

for all $v \in \Omega$, then we see that $\partial\varphi/\partial x_2(u)$ is invertible and both $\partial\varphi/\partial x_2(u)$ and its inverse are continuous, so that Theorem 29.3 applies, and there is some (continuous) linear form $\Lambda(u) \in \mathcal{L}(\mathbb{R}^m; \mathbb{R})$ such that

$$dJ(u) + \Lambda(u) \circ d\varphi(u) = 0.$$

However, $\Lambda(u)$ is defined by some m -tuple $(\lambda_1(u), \dots, \lambda_m(u)) \in \mathbb{R}^m$, and in view of the definition of φ , the above equation is equivalent to

$$dJ(u) + \lambda_1(u)d\varphi_1(u) + \cdots + \lambda_m(u)d\varphi_m(u) = 0.$$

The uniqueness of the $\lambda_i(u)$ is a consequence of the linear independence of the $d\varphi_i(u)$. \square

The numbers $\lambda_i(u)$ involved in Theorem 29.3 are called the *Lagrange multipliers* associated with the constrained extremum u (again, with some minor abuse of language). The linear independence of the linear forms $d\varphi_i(u)$ is equivalent to the fact that the Jacobian matrix $((\partial\varphi_i/\partial x_j)(u))$ of $\varphi = (\varphi_1, \dots, \varphi_m)$ at u has rank m . If $m = 1$, the linear independence of the $d\varphi_i(u)$ reduces to the condition $\nabla\varphi_1(u) \neq 0$.

A fruitful way to reformulate the use of Lagrange multipliers is to introduce the notion of the *Lagrangian* associated with our constrained extremum problem. This is the function $L: \Omega \times \mathbb{R}^m \rightarrow \mathbb{R}$ given by

$$L(v, \lambda) = J(v) + \lambda_1 \varphi_1(v) + \cdots + \lambda_m \varphi_m(v),$$

with $\lambda = (\lambda_1, \dots, \lambda_m)$. Then, observe that there exists some $\mu = (\mu_1, \dots, \mu_m)$ and some $u \in U$ such that

$$dJ(u) + \mu_1 d\varphi_1(u) + \cdots + \mu_m d\varphi_m(u) = 0$$

if and only if

$$dL(u, \mu) = 0,$$

or equivalently

$$\nabla L(u, \mu) = 0;$$

that is, iff (u, λ) is a *critical point* of the Lagrangian L .

Indeed $dL(u, \mu) = 0$ if equivalent to

$$\begin{aligned} \frac{\partial L}{\partial v}(u, \mu) &= 0 \\ \frac{\partial L}{\partial \lambda_1}(u, \mu) &= 0 \\ &\vdots \\ \frac{\partial L}{\partial \lambda_m}(u, \mu) &= 0, \end{aligned}$$

and since

$$\frac{\partial L}{\partial v}(u, \mu) = dJ(u) + \mu_1 d\varphi_1(u) + \cdots + \mu_m d\varphi_m(u)$$

and

$$\frac{\partial L}{\partial \lambda_i}(u, \mu) = \varphi_i(u),$$

we get

$$dJ(u) + \mu_1 d\varphi_1(u) + \cdots + \mu_m d\varphi_m(u) = 0$$

and

$$\varphi_1(u) = \cdots = \varphi_m(u) = 0,$$

that is, $u \in U$.

If we write out explicitly the condition

$$dJ(u) + \mu_1 d\varphi_1(u) + \cdots + \mu_m d\varphi_m(u) = 0,$$

we get the $n \times m$ system

$$\begin{aligned} \frac{\partial J}{\partial x_1}(u) + \lambda_1 \frac{\partial \varphi_1}{\partial x_1}(u) + \cdots + \lambda_m \frac{\partial \varphi_m}{\partial x_1}(u) &= 0 \\ &\vdots \\ \frac{\partial J}{\partial x_n}(u) + \lambda_1 \frac{\partial \varphi_1}{\partial x_n}(u) + \cdots + \lambda_m \frac{\partial \varphi_m}{\partial x_n}(u) &= 0, \end{aligned}$$

and it is important to note that the matrix of this system is the *transpose* of the Jacobian matrix of φ at u . If we write $\text{Jac}(J)(u) = ((\partial \varphi_i / \partial x_j)(u))$ for the Jacobian matrix of J (at u), then the above system is written in matrix form as

$$\nabla J(u) + (\text{Jac}(J)(u))^\top \lambda = 0,$$

where λ is viewed as a column vector, and the Lagrangian is equal to

$$L(u, \lambda) = J(u) + (\varphi_1(u), \dots, \varphi_m(u))\lambda.$$

Remark: If the Jacobian matrix $\text{Jac}(J)(v) = ((\partial \varphi_i / \partial x_j)(v))$ has rank m for all $v \in U$ (which is equivalent to the linear independence of the linear forms $d\varphi_i(v)$), then we say that $0 \in \mathbb{R}^m$ is a *regular value* of φ . In this case, it is known that

$$U = \{v \in \Omega \mid \varphi(v) = 0\}$$

is a *smooth submanifold of dimension $n - m$ of \mathbb{R}^n* . Furthermore, the set

$$T_v U = \{w \in \mathbb{R}^n \mid d\varphi_i(v)(w) = 0, 1 \leq i \leq m\} = \bigcap_{i=1}^m \text{Ker } d\varphi_i(v)$$

is the *tangent space* to U at v (a vector space of dimension $n - m$). Then, the condition

$$dJ(v) + \mu_1 d\varphi_1(v) + \cdots + \mu_m d\varphi_m(v) = 0$$

implies that $dJ(v)$ vanishes on the tangent space $T_v U$. Conversely, if $dJ(v)(w) = 0$ for all $w \in T_v U$, this means that $dJ(v)$ is orthogonal (in the sense of Definition 4.7) to $T_v U$. Since (by Theorem 4.17 (b)) the orthogonal of $T_v U$ is the space of linear forms spanned by $d\varphi_1(v), \dots, d\varphi_m(v)$, it follows that $dJ(v)$ must be a linear combination of the $d\varphi_i(v)$. Therefore, when 0 is a regular value of φ , Theorem 29.3 asserts that if $u \in U$ is a local extremum of J , then $dJ(u)$ must vanish on the tangent space $T_u U$. We can say even more. The subset $Z(J)$ of Ω given by

$$Z(J) = \{v \in \Omega \mid J(v) = J(u)\}$$

(the *level set of level* $J(u)$) is a hypersurface in Ω , and if $dJ(u) \neq 0$, the zero locus of $dJ(u)$ is the tangent space $T_u Z(J)$ to $Z(J)$ at u (a vector space of dimension $n - 1$), where

$$T_u Z(J) = \{w \in \mathbb{R}^n \mid dJ(u)(w) = 0\}.$$

Consequently, Theorem 29.3 asserts that

$$T_u U \subseteq T_u Z(J);$$

this is a geometric condition.

The beauty of the Lagrangian is that the constraints $\{\varphi_i(v) = 0\}$ have been incorporated into the function $L(v, \lambda)$, and that the necessary condition for the existence of a constrained local extremum of J is reduced to the necessary condition for the existence of a local extremum of the *unconstrained* L .

However, one should be careful to check that the assumptions of Theorem 29.3 are satisfied (in particular, the linear independence of the linear forms $d\varphi_i$). For example, let $J: \mathbb{R}^3 \rightarrow \mathbb{R}$ be given by

$$J(x, y, z) = x + y + z^2$$

and $g: \mathbb{R}^3 \rightarrow \mathbb{R}$ by

$$g(x, y, z) = x^2 + y^2.$$

Since $g(x, y, z) = 0$ iff $x = y = 0$, we have $U = \{(0, 0, z) \mid z \in \mathbb{R}\}$ and the restriction of J to U is given by

$$J(0, 0, z) = z^2,$$

which has a minimum for $z = 0$. However, a “blind” use of Lagrange multipliers would require that there is some λ so that

$$\frac{\partial J}{\partial x}(0, 0, z) = \lambda \frac{\partial g}{\partial x}(0, 0, z), \quad \frac{\partial J}{\partial y}(0, 0, z) = \lambda \frac{\partial g}{\partial y}(0, 0, z), \quad \frac{\partial J}{\partial z}(0, 0, z) = \lambda \frac{\partial g}{\partial z}(0, 0, z),$$

and since

$$\frac{\partial g}{\partial x}(x, y, z) = 2x, \quad \frac{\partial g}{\partial y}(x, y, z) = 2y, \quad \frac{\partial g}{\partial z}(0, 0, z) = 0,$$

the partial derivatives above all vanish for $x = y = 0$, so at a local extremum we should also have

$$\frac{\partial J}{\partial x}(0, 0, z) = 0, \quad \frac{\partial J}{\partial y}(0, 0, z) = 0, \quad \frac{\partial J}{\partial z}(0, 0, z) = 0,$$

but this is absurd since

$$\frac{\partial J}{\partial x}(x, y, z) = 1, \quad \frac{\partial J}{\partial y}(x, y, z) = 1, \quad \frac{\partial J}{\partial z}(x, y, z) = 2z.$$

The reader should enjoy finding the reason for the flaw in the argument.

One should also keep in mind that Theorem 29.3 gives only a necessary condition. The (u, λ) may *not* correspond to local extrema! Thus, it is always necessary to analyze the local behavior of J near a critical point u . This is generally difficult, but in the case where J is affine or quadratic and the constraints are affine or quadratic, this is possible (although not always easy).

Let us apply the above method to the following example in which $E_1 = \mathbb{R}$, $E_2 = \mathbb{R}$, $\Omega = \mathbb{R}^2$, and

$$\begin{aligned} J(x_1, x_2) &= -x_2 \\ \varphi(x_1, x_2) &= x_1^2 + x_2^2 - 1. \end{aligned}$$

Observe that

$$U = \{(x_1, x_2) \in \mathbb{R}^2 \mid x_1^2 + x_2^2 = 1\}$$

is the unit circle, and since

$$\nabla \varphi(x_1, x_2) = \begin{pmatrix} 2x_1 \\ 2x_2 \end{pmatrix},$$

it is clear that $\nabla \varphi(x_1, x_2) \neq 0$ for every point $= (x_1, x_2)$ on the unit circle. If we form the Lagrangian

$$L(x_1, x_2, \lambda) = -x_2 + \lambda(x_1^2 + x_2^2 - 1),$$

Theorem 29.3 says that a necessary condition for J to have a constrained local extremum is that $\nabla L(x_1, x_2, \lambda) = 0$, so the following equations must hold:

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

The second equation implies that $\lambda \neq 0$, and then the first yields $x_1 = 0$, so the third yields $x_2 = \pm 1$, and we get two solutions:

$$\begin{aligned} \lambda &= \frac{1}{2}, & (x_1, x_2) &= (0, 1) \\ \lambda &= -\frac{1}{2}, & (x'_1, x'_2) &= (0, -1). \end{aligned}$$

We can check immediately that the first solution is a minimum and the second is a maximum. The reader should look for a geometric interpretation of this problem.

Let us now consider the case in which J is a quadratic function of the form

$$J(v) = \frac{1}{2}v^\top A v - v^\top b,$$

where A is an $n \times n$ symmetric matrix, $b \in \mathbb{R}^n$, and the constraints are given by a linear system of the form

$$Cv = d,$$

where C is an $m \times n$ matrix with $m < n$ and $d \in \mathbb{R}^m$. We also assume that C has rank m . In this case, the function φ is given by

$$\varphi(v) = (Cv - d)^\top,$$

because we view $\varphi(v)$ as a row vector (and v as a column vector), and since

$$d\varphi(v)(w) = C^\top w,$$

the condition that the Jacobian matrix of φ at u have rank m is satisfied. The Lagrangian of this problem is

$$L(v, \lambda) = \frac{1}{2}v^\top Av - v^\top b + (Cv - d)^\top \lambda = \frac{1}{2}v^\top Av - v^\top b + \lambda^\top (Cv - d),$$

where λ is viewed as a column vector. Now, because A is a symmetric matrix, it is easy to show that

$$\nabla L(v, \lambda) = \begin{pmatrix} Av - b + C^\top \lambda \\ Cv - d \end{pmatrix}.$$

Therefore, the necessary condition for constrained local extrema is

$$\begin{aligned} Av + C^\top \lambda &= b \\ Cv &= d, \end{aligned}$$

which can be expressed in matrix form as

$$\begin{pmatrix} A & C^\top \\ C & 0 \end{pmatrix} \begin{pmatrix} v \\ \lambda \end{pmatrix} = \begin{pmatrix} b \\ d \end{pmatrix},$$

where the matrix of the system is a symmetric matrix. We should not be surprised to find the system of Section 18, except for some renaming of the matrices and vectors involved. As we know from Section 18.2, the function J has a minimum iff A is positive definite, so in general, if A is only a symmetric matrix, the critical points of the Lagrangian do *not* correspond to extrema of J .

We now investigate conditions for the existence of extrema involving the second derivative of J .

29.2 Using Second Derivatives to Find Extrema

For the sake of brevity, we consider only the case of local minima; analogous results are obtained for local maxima (replace J by $-J$, since $\max_u J(u) = -\min_u -J(u)$). We begin with a necessary condition for an unconstrained local minimum.

Proposition 29.4. *Let E be a normed vector space and let $J: \Omega \rightarrow \mathbb{R}$ be a function, with Ω some open subset of E . If the function J is differentiable in Ω , if J has a second derivative $D^2J(u)$ at some point $u \in \Omega$, and if J has a local minimum at u , then*

$$D^2J(u)(w, w) \geq 0 \quad \text{for all } w \in E.$$

Proof. Pick any nonzero vector $w \in E$. Since Ω is open, for t small enough, $u + tw \in \Omega$ and $J(u + tw) \geq J(u)$, so there is some open interval $I \subseteq \mathbb{R}$ such that

$$u + tw \in \Omega \quad \text{and} \quad J(u + tw) \geq J(u)$$

for all $t \in I$. Using the Taylor–Young formula and the fact that we must have $dJ(u) = 0$ since J has a local minimum at u , we get

$$0 \leq J(u + tw) - J(u) = \frac{t^2}{2} D^2J(u)(w, w) + t^2 \|w\|^2 \epsilon(tw),$$

with $\lim_{t \rightarrow 0} \epsilon(tw) = 0$, which implies that

$$D^2J(u)(w, w) \geq 0.$$

Since the argument holds for all $w \in E$ (trivially if $w = 0$), the proposition is proved. \square

One should be cautioned that there is no converse to the previous proposition. For example, the function $f: x \mapsto x^3$ has no local minimum at 0, yet $df(0) = 0$ and $D^2f(0)(u, v) = 0$. Similarly, the reader should check that the function $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ given by

$$f(x, y) = x^2 - 3y^3$$

has no local minimum at $(0, 0)$; yet $df(0, 0) = 0$ and $D^2f(0, 0)(u, v) = 2u^2 \geq 0$.

When $E = \mathbb{R}^n$, Proposition 29.4 says that a necessary condition for having a local minimum is that the Hessian $\nabla^2 J(u)$ be positive semidefinite (it is always symmetric).

We now give sufficient conditions for the existence of a local minimum.

Theorem 29.5. *Let E be a normed vector space, let $J: \Omega \rightarrow \mathbb{R}$ be a function with Ω some open subset of E , and assume that J is differentiable in Ω and that $dJ(u) = 0$ at some point $u \in \Omega$. The following properties hold:*

(1) *If $D^2J(u)$ exists and if there is some number $\alpha \in \mathbb{R}$ such that $\alpha > 0$ and*

$$D^2J(u)(w, w) \geq \alpha \|w\|^2 \quad \text{for all } w \in E,$$

then J has a strict local minimum at u .

(2) If $D^2J(v)$ exists for all $v \in \Omega$ and if there is a ball $B \subseteq \Omega$ centered at u such that

$$D^2J(v)(w, w) \geq 0 \quad \text{for all } v \in B \text{ and all } w \in E,$$

then J has a local minimum at u .

Proof. (1) Using the formula of Taylor–Young, for every vector w small enough, we can write

$$\begin{aligned} J(u + w) - J(u) &= \frac{1}{2}D^2J(u)(w, w) + \|w\|^2 \epsilon(w) \\ &\geq \left(\frac{1}{2}\alpha + \epsilon(w) \right) \|w\|^2 \end{aligned}$$

with $\lim_{w \rightarrow 0} \epsilon(w) = 0$. Consequently if we pick $r > 0$ small enough that $|\epsilon(w)| < \alpha$ for all w with $\|w\| < r$, then $J(u + w) > J(u)$ for all $u + w \in B$, where B is the open ball of center u and radius r . This proves that J has a local strict minimum at u .

(2) The formula of Taylor–Maclaurin shows that for all $u + w \in B$, we have

$$J(u + w) = J(u) + \frac{1}{2}D^2J(v)(w, w) \geq J(u),$$

for some $v \in]u, u + w[$. □

There are no converses of the two assertions of Theorem 29.5. However, there is a condition on $D^2J(u)$ that implies the condition of part (1). Since this condition is easier to state when $E = \mathbb{R}^n$, we begin with this case.

Recall that a $n \times n$ symmetric matrix A is *positive definite* if $x^\top Ax > 0$ for all $x \in \mathbb{R}^n - \{0\}$. In particular, A must be invertible.

Proposition 29.6. *For any symmetric matrix A , if A is positive definite, then there is some $\alpha > 0$ such that*

$$x^\top Ax \geq \alpha \|x\|^2 \quad \text{for all } x \in \mathbb{R}^n.$$

Proof. Pick any norm in \mathbb{R}^n (recall that all norms on \mathbb{R}^n are equivalent). Since the unit sphere $S^{n-1} = \{x \in \mathbb{R}^n \mid \|x\| = 1\}$ is compact and since the function $f(x) = x^\top Ax$ is never zero on S^{n-1} , the function f has a minimum $\alpha > 0$ on S^{n-1} . Using the usual trick that $x = \|x\| (x/\|x\|)$ for every nonzero vector $x \in \mathbb{R}^n$ and the fact that the inequality of the proposition is trivial for $x = 0$, from

$$x^\top Ax \geq \alpha \quad \text{for all } x \text{ with } \|x\| = 1,$$

we get

$$x^\top Ax \geq \alpha \|x\|^2 \quad \text{for all } x \in \mathbb{R}^n,$$

as claimed. □

We can combine Theorem 29.5 and Proposition 29.6 to obtain a useful sufficient condition for the existence of a strict local minimum. First let us introduce some terminology.

Given a function $J: \Omega \rightarrow \mathbb{R}$ as before, say that a point $u \in \Omega$ is a *nondegenerate critical point* if $dJ(u) = 0$ and if the Hessian matrix $\nabla^2 J(u)$ is invertible.

Proposition 29.7. *Let $J: \Omega \rightarrow \mathbb{R}$ be a function defined on some open subset $\Omega \subseteq \mathbb{R}^n$. If J is differentiable in Ω and if some point $u \in \Omega$ is a nondegenerate critical point such that $\nabla^2 J(u)$ is positive definite, then J has a strict local minimum at u .*

Remark: It is possible to generalize Proposition 29.7 to infinite-dimensional spaces by finding a suitable generalization of the notion of a nondegenerate critical point. Firstly, we assume that E is a Banach space (a complete normed vector space). Then, we define the dual E' of E as the set of continuous linear forms on E , so that $E' = \mathcal{L}(E; \mathbb{R})$. Following Lang, we use the notation E' for the space of continuous linear forms to avoid confusion with the space $E^* = \text{Hom}(E, \mathbb{R})$ of all linear maps from E to \mathbb{R} . A continuous bilinear map $\varphi: E \times E \rightarrow \mathbb{R}$ in $\mathcal{L}_2(E, E; \mathbb{R})$ yields a map Φ from E to E' given by

$$\Phi(u) = \varphi_u,$$

where $\varphi_u \in E'$ is the linear form defined by

$$\varphi_u(v) = \varphi(u, v).$$

It is easy to check that φ_u is continuous and that the map Φ is continuous. Then, we say that φ is *nondegenerate* iff $\Phi: E \rightarrow E'$ is an isomorphism of Banach spaces, which means that Φ is invertible and that both Φ and Φ^{-1} are continuous linear maps. Given a function $J: \Omega \rightarrow \mathbb{R}$ differentiable on Ω as before (where Ω is an open subset of E), if $D^2 J(u)$ exists for some $u \in \Omega$, we say that u is a *nondegenerate critical point* if $dJ(u) = 0$ and if $D^2 J(u)$ is nondegenerate. Of course, $D^2 J(u)$ is positive definite if $D^2 J(u)(w, w) > 0$ for all $w \in E - \{0\}$.

Using the above definition, Proposition 29.6 can be generalized to a nondegenerate positive definite bilinear form (on a Banach space) and Theorem 29.7 can also be generalized to the situation where $J: \Omega \rightarrow \mathbb{R}$ is defined on an open subset of a Banach space. For details and proofs, see Cartan [20] (Part I Chapter 8) and Avez [6] (Chapter 8 and Chapter 10).

In the next section, we make use of convexity; both on the domain Ω and on the function J itself.

29.3 Using Convexity to Find Extrema

We begin by reviewing the definition of a convex set and of a convex function.

Definition 29.3. Given any real vector space E , we say that a subset C of E is *convex* if either $C = \emptyset$ or if for every pair of points $u, v \in C$,

$$(1 - \lambda)u + \lambda v \in C \quad \text{for all } \lambda \in \mathbb{R} \text{ such that } 0 \leq \lambda \leq 1.$$

If C is a nonempty convex subset of E , a function $f: C \rightarrow \mathbb{R}$ is *convex* (on C) if for every pair of points $u, v \in C$,

$$f((1 - \lambda)u + \lambda v) \leq (1 - \lambda)f(u) + \lambda f(v) \quad \text{for all } \lambda \in \mathbb{R} \text{ such that } 0 \leq \lambda \leq 1;$$

the function f is *strictly convex* (on C) if for every pair of distinct points $u, v \in C$ ($u \neq v$),

$$f((1 - \lambda)u + \lambda v) < (1 - \lambda)f(u) + \lambda f(v) \quad \text{for all } \lambda \in \mathbb{R} \text{ such that } 0 < \lambda < 1.$$

A function $f: C \rightarrow \mathbb{R}$ defined on a convex subset C is *concave* (resp. *strictly concave*) if $(-f)$ is convex (resp. strictly convex).

Given any two points $u, v \in E$, the *line segment* $[u, v]$ is the set

$$[u, v] = \{(1 - \lambda)u + \lambda v \in E \mid \lambda \in \mathbb{R}, 0 \leq \lambda \leq 1\}.$$

Clearly, a nonempty set C is convex iff $[u, v] \subseteq C$ whenever $u, v \in C$. Subspaces $V \subseteq E$ of a vector space E are convex; *affine subspaces*, that is, sets of the form $u + V$, where V is a subspace of E and $u \in E$, are convex. Balls (open or closed) are convex. Given any linear form $\varphi: E \rightarrow \mathbb{R}$, for any scalar $c \in \mathbb{R}$, the *closed half-spaces*

$$H_{\varphi, c}^+ = \{u \in E \mid \varphi(u) \geq c\}, \quad H_{\varphi, c}^- = \{u \in E \mid \varphi(u) \leq c\},$$

are convex. Any intersection of half-spaces is convex. More generally, any intersection of convex sets is convex.

Linear forms are convex functions (but not strictly convex). Any norm $\|\cdot\|: E \rightarrow \mathbb{R}_+$ is a convex function. The max function,

$$\max(x_1, \dots, x_n) = \max\{x_1, \dots, x_n\}$$

is convex on \mathbb{R}^n . The exponential $x \mapsto e^{cx}$ is strictly convex for any $c \neq 0$ ($c \in \mathbb{R}$). The logarithm function is concave on $\mathbb{R}_+ - \{0\}$, and the *log-determinant function* $\log \det$ is concave on the set of symmetric positive definite matrices. This function plays an important role in convex optimization. An excellent exposition of convexity and its applications to optimization can be found in Boyd [17].

Here is a necessary condition for a function to have a local minimum with respect to a convex subset U .

Theorem 29.8. (Necessary condition for a local minimum on a convex subset) Let $J: \Omega \rightarrow \mathbb{R}$ be a function defined on some open subset Ω of a normed vector space E and let $U \subseteq \Omega$ be a nonempty convex subset. Given any $u \in U$, if $dJ(u)$ exists and if J has a local minimum in u with respect to U , then

$$dJ(u)(v - u) \geq 0 \quad \text{for all } v \in U.$$

Proof. Let $v = u + w$ be an arbitrary point in U . Since U is convex, we have $u + tw \in U$ for all t such that $0 \leq t \leq 1$. Since $dJ(u)$ exists, we can write

$$J(u + tw) - J(u) = dJ(u)(tw) + \|tw\| \epsilon(tw)$$

with $\lim_{t \rightarrow 0} \epsilon(tw) = 0$. However, because $0 \leq t \leq 1$,

$$J(u + tw) - J(u) = t(dJ(u)(w) + \|w\| \epsilon(tw))$$

and since u is a local minimum with respect to U , we have $J(u + tw) - J(u) \geq 0$, so we get

$$t(dJ(u)(w) + \|w\| \epsilon(tw)) \geq 0.$$

The above implies that $dJ(u)(w) \geq 0$, because otherwise we could pick $t > 0$ small enough so that

$$dJ(u)(w) + \|w\| \epsilon(tw) < 0,$$

a contradiction. Since the argument holds for all $v = u + w \in U$, the theorem is proved. \square

Observe that the convexity of U is a substitute for the use of Lagrange multipliers, but we now have to deal with an *inequality* instead of an equality.

Consider the special case where U is a subspace of E . In this case, since $u \in U$ we have $2u \in U$, and for any $u + w \in U$, we must have $2u - (u + w) = u - w \in U$. The previous theorem implies that $dJ(u)(w) \geq 0$ and $dJ(u)(-w) \geq 0$, that is, $dJ(u)(w) \leq 0$, so $dJ(u)(w) = 0$. Since the argument holds for $w \in U$ (because U is a subspace, if $u, w \in U$, then $u + w \in U$), we conclude that

$$dJ(u)(w) = 0 \quad \text{for all } w \in U.$$

We will now characterize convex functions when they have a first derivative or a second derivative.

Proposition 29.9. (Convexity and first derivative) Let $f: \Omega \rightarrow \mathbb{R}$ be a function differentiable on some open subset Ω of a normed vector space E and let $U \subseteq \Omega$ be a nonempty convex subset.

(1) The function f is convex on U iff

$$f(v) \geq f(u) + df(u)(v - u) \quad \text{for all } u, v \in U.$$

(2) The function f is strictly convex on U iff

$$f(v) > f(u) + df(u)(v - u) \quad \text{for all } u, v \in U \text{ with } u \neq v.$$

Proof. Let $u, v \in U$ be any two distinct points and pick $\lambda \in \mathbb{R}$ with $0 < \lambda < 1$. If the function f is convex, then

$$f((1 - \lambda)u + \lambda v) \leq (1 - \lambda)f(u) + \lambda f(v),$$

which yields

$$\frac{f((1 - \lambda)u + \lambda v) - f(u)}{\lambda} \leq f(v) - f(u).$$

It follows that

$$df(u)(v - u) = \lim_{\lambda \rightarrow 0} \frac{f((1 - \lambda)u + \lambda v) - f(u)}{\lambda} \leq f(v) - f(u).$$

If f is strictly convex, the above reasoning does not work, because a strict inequality is not necessarily preserved by “passing to the limit.” We have recourse to the following trick: For any ω such that $0 < \omega < 1$, observe that

$$(1 - \lambda)u + \lambda v = u + \lambda(v - u) = \frac{\omega - \lambda}{\omega}u + \frac{\lambda}{\omega}(u + \omega(v - u)).$$

If we assume that $0 < \lambda \leq \omega$, the convexity of f yields

$$f(u + \lambda(v - u)) \leq \frac{\omega - \lambda}{\omega}f(u) + \frac{\lambda}{\omega}f(u + \omega(v - u)).$$

If we subtract $f(u)$ to both sides, we get

$$\frac{f(u + \lambda(v - u)) - f(u)}{\lambda} \leq \frac{f(u + \omega(v - u)) - f(u)}{\omega}.$$

Now, since $0 < \omega < 1$ and f is strictly convex,

$$f(u + \omega(v - u)) = f((1 - \omega)u + \omega v) < (1 - \omega)f(u) + \omega f(v),$$

which implies that

$$\frac{f(u + \omega(v - u)) - f(u)}{\omega} < f(v) - f(u),$$

and thus we get

$$\frac{f(u + \lambda(v - u)) - f(u)}{\lambda} \leq \frac{f(u + \omega(v - u)) - f(u)}{\omega} < f(v) - f(u).$$

If we let λ go to 0, by passing to the limit we get

$$df(u)(v - u) \leq \frac{f(u + \omega(v - u)) - f(u)}{\omega} < f(v) - f(u),$$

which yields the desired strict inequality.

Let us now consider the converse of (1); that is, assume that

$$f(v) \geq f(u) + df(u)(v - u) \quad \text{for all } u, v \in U.$$

For any two distinct points $u, v \in U$ and for any λ with $0 < \lambda < 1$, we get

$$\begin{aligned} f(v) &\geq f(v + \lambda(v - u)) - \lambda df(v + \lambda(u - v))(u - v) \\ f(u) &\geq f(v + \lambda(u - v)) + (1 - \lambda)df(v + \lambda(u - v))(u - v), \end{aligned}$$

and if we multiply the first inequality by $1 - \lambda$ and the second inequality by λ and then add up the resulting inequalities, we get

$$(1 - \lambda)f(v) + \lambda f(u) \geq f(v + \lambda(u - v)) = f((1 - \lambda)v + \lambda u),$$

which proves that f is convex.

The proof of the converse of (2) is similar, except that the inequalities are replaced by strict inequalities. \square

We now establish a convexity criterion using the second derivative of f . This criterion is often easier to check than the previous one.

Proposition 29.10. (*Convexity and second derivative*) Let $f: \Omega \rightarrow \mathbb{R}$ be a function twice differentiable on some open subset Ω of a normed vector space E and let $U \subseteq \Omega$ be a nonempty convex subset.

(1) The function f is convex on U iff

$$D^2f(u)(v - u, v - u) \geq 0 \quad \text{for all } u, v \in U.$$

(2) If

$$D^2f(u)(v - u, v - u) > 0 \quad \text{for all } u, v \in U \text{ with } u \neq v,$$

then f is strictly convex.

Proof. First, assume that the inequality in condition (1) is satisfied. For any two distinct points $u, v \in U$, the formula of Taylor–Maclaurin yields

$$\begin{aligned} f(v) - f(u) - df(u)(v - u) &= \frac{1}{2}D^2(w)(v - u, v - u) \\ &= \frac{\rho^2}{2}D^2(w)(v - w, v - w), \end{aligned}$$

for some $w = (1 - \lambda)u + \lambda v = u + \lambda(v - u)$ with $0 < \lambda < 1$, and with $\rho = 1/(1 - \lambda) > 0$, so that $v - u = \rho(v - w)$. Since $D^2f(u)(v - w, v - w) \geq 0$ for all $u, w \in U$, we conclude by applying Theorem 29.9(1).

Similarly, if (2) holds, the above reasoning and Theorem 29.9(2) imply that f is strictly convex.

To prove the necessary condition in (1), define $g: \Omega \rightarrow \mathbb{R}$ by

$$g(v) = f(v) - df(u)(v),$$

where $u \in U$ is any point considered fixed. If f is convex and if f has a local minimum at u with respect to U , since

$$g(v) - g(u) = f(v) - f(u) - df(u)(v - u),$$

Theorem 29.9 implies that $f(v) - f(u) - df(u)(v - u) \geq 0$, which implies that g has a local minimum at u with respect to all $v \in U$. Therefore, we have $dg(u) = 0$. Observe that g is twice differentiable in Ω and $D^2g(u) = D^2f(u)$, so the formula of Taylor–Young yields for every $v = u + w \in U$ and all t with $0 \leq t \leq 1$,

$$\begin{aligned} 0 \leq g(u + tw) - g(u) &= \frac{t^2}{2} D^2(u)(tw, tw) + \|tw\|^2 \epsilon(tw) \\ &= \frac{t^2}{2} (D^2(u)(w, w) + 2\|w\|^2 \epsilon(wt)), \end{aligned}$$

with $\lim_{t \rightarrow 0} \epsilon(wt) = 0$, and for t small enough, we must have $D^2(u)(w, w) \geq 0$, as claimed. \square

The converse of Theorem 29.10 (2) is false as we see by considering the function f given by $f(x) = x^4$. On the other hand, if f is a quadratic function of the form

$$f(u) = \frac{1}{2} u^\top A u - u^\top b$$

where A is a symmetric matrix, we know that

$$df(u)(v) = v^\top (A u - b),$$

so

$$\begin{aligned} f(v) - f(u) - df(u)(v - u) &= \frac{1}{2} v^\top A v - v^\top b - \frac{1}{2} u^\top A u + u^\top b - (v - u)^\top (A u - b) \\ &= \frac{1}{2} v^\top A v - \frac{1}{2} u^\top A u - (v - u)^\top A u \\ &= \frac{1}{2} v^\top A v + \frac{1}{2} u^\top A u - v^\top A u \\ &= \frac{1}{2} (v - u)^\top A (v - u). \end{aligned}$$

Therefore, Theorem 29.9 implies that if A is positive semidefinite, then f is convex and if A is positive definite, then f is strictly convex. The converse follows by Theorem 29.10.

We conclude this section by applying our previous theorems to convex functions defined on convex subsets. In this case, local minima (resp. local maxima) are global minima (resp. global maxima).

Definition 29.4. Let $f: E \rightarrow \mathbb{R}$ be any function defined on some normed vector space (or more generally, any set). For any $u \in E$, we say that f has a *minimum* in u (resp. *maximum* in u) if

$$f(u) \leq f(v) \text{ (resp. } f(u) \geq f(v)) \text{ for all } v \in E.$$

We say that f has a *strict minimum* in u (resp. *strict maximum* in u) if

$$f(u) < f(v) \text{ (resp. } f(u) > f(v)) \text{ for all } v \in E - \{u\}.$$

If $U \subseteq E$ is a subset of E and $u \in U$, we say that f has a *minimum* in u (resp. *strict minimum* in u) *with respect to* U if

$$f(u) \leq f(v) \text{ for all } v \in U \text{ (resp. } f(u) < f(v) \text{ for all } v \in U - \{u\}),$$

and similarly for a *maximum* in u (resp. *strict maximum* in u) *with respect to* U with \leq changed to \geq and $<$ to $>$.

Sometimes, we say *global* maximum (or minimum) to stress that a maximum (or a minimum) is not simply a local maximum (or minimum).

Theorem 29.11. *Given any normed vector space E , let U be any nonempty convex subset of E .*

- (1) *For any convex function $J: U \rightarrow \mathbb{R}$, for any $u \in U$, if J has a local minimum at u in U , then J has a (global) minimum at u in U .*
- (2) *Any strict convex function $J: U \rightarrow \mathbb{R}$ has at most one minimum (in U), and if it does, then it is a strict minimum (in U).*
- (3) *Let $J: \Omega \rightarrow \mathbb{R}$ be any function defined on some open subset Ω of E with $U \subseteq \Omega$ and assume that J is convex on U . For any point $u \in U$, if $dJ(u)$ exists, then J has a minimum in u with respect to U iff*

$$dJ(u)(v - u) \geq 0 \text{ for all } v \in U.$$

- (4) *If the convex subset U in (3) is open, then the above condition is equivalent to*

$$dJ(u) = 0.$$

Proof. (1) Let $v = u + w$ be any arbitrary point in U . Since J is convex, for all t with $0 \leq t \leq 1$, we have

$$J(u + tw) = J(u + t(v - u)) \leq (1 - t)J(u) + tJ(v),$$

which yields

$$J(u + tw) - J(u) \leq t(J(v) - J(u)).$$

Because J has a local minimum in u , there is some t_0 with $0 < t_0 < 1$ such that

$$0 \leq J(u + t_0 w) - J(u),$$

which implies that $J(v) - J(u) \geq 0$.

(2) If J is strictly convex, the above reasoning with $w \neq 0$ shows that there is some t_0 with $0 < t_0 < 1$ such that

$$0 \leq J(u + t_0 w) - J(u) < t_0(J(v) - J(u)),$$

which shows that u is a strict global minimum (in U), and thus that it is unique.

(3) We already know from Theorem 29.9 that the condition $dJ(u)(v - u) \geq 0$ for all $v \in U$ is necessary (even if J is not convex). Conversely, because J is convex, careful inspection of the proof of part (1) of Proposition 29.9 shows that only the fact that $dJ(u)$ exists is needed to prove that

$$J(v) - J(u) \geq dJ(u)(v - u) \quad \text{for all } v \in U,$$

and if

$$dJ(u)(v - u) \geq 0 \quad \text{for all } v \in U,$$

then

$$J(v) - J(u) \geq 0 \quad \text{for all } v \in U,$$

as claimed.

(4) If U is open, then for every $u \in U$ we can find an open ball B centered at u of radius ϵ small enough so that $B \subseteq U$. Then, for any $w \neq 0$ such that $\|w\| < \epsilon$, we have both $v = u + w \in B$ and $v' = u - w \in B$, so condition (3) implies that

$$dJ(u)(w) \geq 0 \quad \text{and} \quad dJ(u)(-w) \geq 0,$$

which yields

$$dJ(u)(w) = 0.$$

Since the above holds for all $w \neq 0$ such that $\|w\| < \epsilon$ and since $dJ(u)$ is linear, we leave it to the reader to fill in the details of the proof that $dJ(u) = 0$. \square

Theorem 29.11 can be used to rederive the fact that the least squares solutions of a linear system $Ax = b$ (where A is an $m \times n$ matrix) are given by the normal equation

$$A^\top Ax = A^\top b.$$

For this, we consider the quadratic function

$$J(v) = \frac{1}{2} \|Av - b\|_2^2 - \frac{1}{2} \|b\|_2^2,$$

and our least squares problem is equivalent to finding the minima of J on \mathbb{R}^n . A computation reveals that

$$J(v) = \frac{1}{2}v^\top A^\top Av - v^\top B^\top b,$$

and so

$$dJ(u) = A^\top Au - B^\top b.$$

Since $B^\top B$ is positive semidefinite, the function J is convex, and Theorem 29.11(4) implies that the minima of J are the solutions of the equation

$$A^\top Au - A^\top b = 0.$$

The considerations in this chapter reveal the need to find methods for finding the zeros of the derivative map

$$dJ: \Omega \rightarrow E',$$

where Ω is some open subset of a normed vector space E and E' is the space of all continuous linear forms on E (a subspace of E^*). Generalizations of *Newton's method* yield such methods and they are the object of the next chapter.

29.4 Summary

The main concepts and results of this chapter are listed below:

- *Local minimum, local maximum, local extremum, strict local minimum, strict local maximum.*
- Necessary condition for a local extremum involving the derivative; *critical point*.
- *Local minimum with respect to a subset U , local maximum with respect to a subset U , local extremum with respect to a subset U .*
- *Constrained local extremum.*
- Necessary condition for a constrained extremum.
- Necessary condition for a constrained extremum in terms of *Lagrange multipliers*.
- *Lagrangian.*
- *Critical points of a Lagrangian.*
- Necessary condition of an unconstrained local minimum involving the second-order derivative.
- Sufficient condition for a local minimum involving the second-order derivative.

- A sufficient condition involving *nondegenerate critical points*.
- *Convex sets, convex functions, concave functions, strictly convex functions, strictly concave functions,*
- Necessary condition for a local minimum on a convex set involving the derivative.
- Convexity of a function involving a condition on its first derivative.
- Convexity of a function involving a condition on its second derivative.
- Minima of convex functions on convex sets.

Chapter 30

Newton's Method and its Generalizations

30.1 Newton's Method for Real Functions of a Real Argument

In Chapter 29 we investigated the problem of determining when a function $J: \Omega \rightarrow \mathbb{R}$ defined on some open subset Ω of a normed vector space E has a local extremum. Proposition 29.1 gives a necessary condition when J is differentiable: if J has a local extremum at $u \in \Omega$, then we must have

$$J'(u) = 0.$$

Thus, we are led to the problem of finding the zeros of the derivative

$$J': \Omega \rightarrow E',$$

where $E' = \mathcal{L}(E; \mathbb{R})$ is the set of linear continuous functions from E to \mathbb{R} ; that is, the *dual* of E , as defined in the Remark after Proposition 29.7.

This leads us to consider the problem in a more general form, namely: Given a function $f: \Omega \rightarrow Y$ from an open subset Ω of a normed vector space X to a normed vector space Y , find

- (i) Sufficient conditions which guarantee the *existence of a zero* of the function f ; that is, an element $a \in \Omega$ such that $f(a) = 0$.
- (ii) An *algorithm* for approximating such an a , that is, a sequence (x_k) of points of Ω whose limit is a .

When $X = Y = \mathbb{R}$, we can use *Newton's method*. We pick some initial element $x_0 \in \mathbb{R}$ “close enough” to a zero a of f , and we define the sequence (x_k) by

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)},$$

for all $k \geq 0$, provided that $f'(x_k) \neq 0$. The idea is to define x_{k+1} as the intersection of the x -axis with the tangent line to the graph of the function $x \mapsto f(x)$ at the point $(x_k, f(x_k))$. Indeed, the equation of this tangent line is

$$y - f(x_k) = f'(x_k)(x - x_k),$$

and its intersection with the x -axis is obtained for $y = 0$, which yields

$$x = x_k - \frac{f(x_k)}{f'(x_k)},$$

as claimed.

For example, if $\alpha > 0$ and $f(x) = x^2 - \alpha$, Newton's method yields the sequence

$$x_{k+1} = \frac{1}{2} \left(x_k + \frac{\alpha}{x_k} \right)$$

to compute the square root $\sqrt{\alpha}$ of α . It can be shown that the method converges to $\sqrt{\alpha}$ for any $x_0 > 0$. Actually, the method also converges when $x_0 < 0$! Find out what is the limit.

The case of a real function suggests the following method for finding the zeros of a function $f: \Omega \rightarrow Y$, with $\Omega \subseteq X$: given a starting point $x_0 \in \Omega$, the sequence (x_k) is defined by

$$x_{k+1} = x_k - (f'(x_k))^{-1}(f(x_k))$$

for all $k \geq 0$.

For the above to make sense, it must be ensured that

- (1) All the points x_k remain within Ω .
- (2) The function f is differentiable within Ω .
- (3) The derivative $f'(x)$ is a bijection from X to Y for all $x \in \Omega$.

These are rather demanding conditions but there are sufficient conditions that guarantee that they are met. Another practical issue is that it may be very costly to compute $(f'(x_k))^{-1}$ at every iteration step. In the next section, we investigate generalizations of Newton's method which address the issues that we just discussed.

30.2 Generalizations of Newton's Method

Suppose that $f: \Omega \rightarrow \mathbb{R}^n$ is given by n functions $f_i: \Omega \rightarrow \mathbb{R}$, where $\Omega \subseteq \mathbb{R}^n$. In this case, finding a zero a of f is equivalent to solving the system

$$\begin{aligned} f_1(a_1, \dots, a_n) &= 0 \\ f_2(a_1, \dots, a_n) &= 0 \\ &\vdots \\ f_n(a_1, \dots, a_n) &= 0. \end{aligned}$$

A single iteration of Newton's method consists in solving the linear system

$$(J(f)(x_k))\epsilon_k = -f(x_k),$$

and then setting

$$x_{k+1} = x_k + \epsilon_k,$$

where $J(f)(x_k) = (\frac{\partial f_i}{\partial x_j}(x_k))$ is the Jacobian matrix of f at x_k .

In general, it is very costly to compute $J(f)(x_k)$ at each iteration and then to solve the corresponding linear system. If the method converges, the consecutive vectors x_k should differ only a little, as also the corresponding matrices $J(f)(x_k)$. Thus, we are led to a variant of Newton's method which consists in keeping the same matrix for p consecutive steps (where p is some fixed integer ≥ 2):

$$\begin{aligned} x_{k+1} &= x_k - (f'(x_0))^{-1}(f(x_k)), & 0 \leq k \leq p-1 \\ x_{k+1} &= x_k - (f'(x_p))^{-1}(f(x_k)), & p \leq k \leq 2p-1 \\ &\vdots \\ x_{k+1} &= x_k - (f'(x_{rp}))^{-1}(f(x_k)), & rp \leq k \leq (r+1)p-1 \\ &\vdots \end{aligned}$$

It is also possible to set $p = \infty$, that is, to use the same matrix $f'(x_0)$ for all iterations, which leads to iterations of the form

$$x_{k+1} = x_k - (f'(x_0))^{-1}(f(x_k)), \quad k \geq 0,$$

or even to replace $f'(x_0)$ by a particular matrix A_0 which is easy to invert:

$$x_{k+1} = x_k - A_0^{-1}f(x_k), \quad k \geq 0.$$

In the last two cases, if possible, we use an LU factorization of $f'(x_0)$ or A_0 to speed up the method. In some cases, it may even be possible to set $A_0 = I$.

The above considerations lead us to the definition of a *generalized Newton method*, as in Ciarlet [24] (Chapter 7). Recall that a linear map $f \in \mathcal{L}(E; F)$ is called an *isomorphism* iff f is continuous, bijective, and f^{-1} is also continuous.

Definition 30.1. If X and Y are two normed vector spaces and if $f: \Omega \rightarrow Y$ is a function from some open subset Ω of X , a *generalized Newton method* for finding zeros of f consists of

- (1) A sequence of families $(A_k(x))$ of linear isomorphisms from X to Y , for all $x \in \Omega$ and all integers $k \geq 0$;
- (2) Some starting point $x_0 \in \Omega$;

(3) A sequence (x_k) of points of Ω defined by

$$x_{k+1} = x_k - (A_k(x_\ell))^{-1}(f(x_k)), \quad k \geq 0,$$

where for every integer $k \geq 0$, the integer ℓ satisfies the condition

$$0 \leq \ell \leq k.$$

The function $A_k(x)$ usually depends on f' .

Definition 30.1 gives us enough flexibility to capture all the situations that we have previously discussed:

$$\begin{aligned} A_k(x) &= f'(x), & \ell &= k \\ A_k(x) &= f'(x), & \ell &= \min\{rp, k\}, \text{ if } rp \leq k \leq (r+1)p-1, r \geq 0 \\ A_k(x) &= f'(x), & \ell &= 0 \\ A_k(x) &= A_0, \end{aligned}$$

where A_0 is a linear isomorphism from X to Y . The first case corresponds to Newton's original method and the others to the variants that we just discussed. We could also have $A_k(x) = A_k$, a fixed linear isomorphism independent of $x \in \Omega$.

The following theorem inspired by the *Newton–Kantorovich theorem* gives sufficient conditions that guarantee that the sequence (x_k) constructed by a generalized Newton method converges to a zero of f close to x_0 . Although quite technical, these conditions are not very surprising.

Theorem 30.1. *Let X be a Banach space, let $f: \Omega \rightarrow Y$ be differentiable on the open subset $\Omega \subseteq X$, and assume that there are constants $r, M, \beta > 0$ such that if we let*

$$B = \{x \in X \mid \|x - x_0\| \leq r\} \subseteq \Omega,$$

then

(1)

$$\sup_{k \geq 0} \sup_{x \in B} \|A_k^{-1}(x)\|_{\mathcal{L}(Y;X)} \leq M,$$

(2) $\beta < 1$ and

$$\sup_{k \geq 0} \sup_{x, x' \in B} \|f'(x) - A_k(x')\|_{\mathcal{L}(X;Y)} \leq \frac{\beta}{M}$$

(3)

$$\|f(x_0)\| \leq \frac{r}{M}(1 - \beta).$$

Then, the sequence (x_k) defined by

$$x_{k+1} = x_k - A_k^{-1}(x_k)(f(x_k)), \quad 0 \leq k \leq \ell$$

is entirely contained within B and converges to a zero a of f , which is the only zero of f in B . Furthermore, the convergence is geometric, which means that

$$\|x_k - a\| \leq \frac{\|x_1 - x_0\|}{1 - \beta} \beta^k.$$

A proof of Theorem 30.1 can be found in Ciarlet [24] (Section 7.5). It is not really difficult but quite technical.

If we assume that we already know that some element $a \in \Omega$ is a zero of f , the next theorem gives sufficient conditions for a special version of a generalized Newton method to converge. For this special method, the linear isomorphisms $A_k(x)$ are independent of $x \in \Omega$.

Theorem 30.2. *Let X be a Banach space, and let $f: \Omega \rightarrow Y$ be differentiable on the open subset $\Omega \subseteq X$. If $a \in \Omega$ is a point such that $f(a) = 0$, if $f'(a)$ is a linear isomorphism, and if there is some λ with $0 < \lambda < 1/2$ such that*

$$\sup_{k \geq 0} \|A_k - f'(a)\|_{\mathcal{L}(X;Y)} \leq \frac{\lambda}{\|(f'(a))^{-1}\|_{\mathcal{L}(Y;X)}},$$

then there is a closed ball B of center a such that for every $x_0 \in B$, the sequence (x_k) defined by

$$x_{k+1} = x_k - A_k^{-1}(f(x_k)), \quad k \geq 0,$$

is entirely contained within B and converges to a , which is the only zero of f in B . Furthermore, the convergence is geometric, which means that

$$\|x_k - a\| \leq \beta^k \|x_0 - a\|,$$

for some $\beta < 1$.

A proof of Theorem 30.2 can be also found in Ciarlet [24] (Section 7.5).

For the sake of completeness, we state a version of the Newton–Kantorovich theorem, which corresponds to the case where $A_k(x) = f'(x)$. In this instance, a stronger result can be obtained especially regarding upper bounds, and we state a version due to Gragg and Tapia which appears in Problem 7.5-4 of Ciarlet [24].

Theorem 30.3. *(Newton–Kantorovich) Let X be a Banach space, and let $f: \Omega \rightarrow Y$ be differentiable on the open subset $\Omega \subseteq X$. Assume that there exist three positive constants λ, μ, ν and a point $x_0 \in \Omega$ such that*

$$0 < \lambda\mu\nu \leq \frac{1}{2},$$

and if we let

$$\begin{aligned}\rho^- &= \frac{1 - \sqrt{1 - 2\lambda\mu\nu}}{\mu\nu} \\ \rho^+ &= \frac{1 + \sqrt{1 - 2\lambda\mu\nu}}{\mu\nu} \\ B &= \{x \in X \mid \|x - x_0\| < \rho^-\} \\ \Omega^+ &= \{x \in \Omega \mid \|x - x_0\| < \rho^+\},\end{aligned}$$

then $\overline{B} \subseteq \Omega$, $f'(x_0)$ is an isomorphism of $\mathcal{L}(X; Y)$, and

$$\begin{aligned}\|(f'(x_0))^{-1}\| &\leq \mu, \\ \|(f'(x_0))^{-1}f(x_0)\| &\leq \lambda, \\ \sup_{x, y \in \Omega^+} \|f'(x) - f'(y)\| &\leq \nu \|x - y\|.\end{aligned}$$

Then, $f'(x)$ is isomorphism of $\mathcal{L}(X; Y)$ for all $x \in B$, and the sequence defined by

$$x_{k+1} = x_k - (f'(x_k))^{-1}(f(x_k)), \quad k \geq 0$$

is entirely contained within the ball B and converges to a zero a of f which is the only zero of f in Ω^+ . Finally, if we write $\theta = \rho^-/\rho^+$, then we have the following bounds:

$$\begin{aligned}\|x_k - a\| &\leq \frac{2\sqrt{1 - 2\lambda\mu\nu}}{\lambda\mu\nu} \frac{\theta^{2k}}{1 - \theta^{2k}} \|x_1 - x_0\| && \text{if } \lambda\mu\nu < \frac{1}{2} \\ \|x_k - a\| &\leq \frac{\|x_1 - x_0\|}{2^{k-1}} && \text{if } \lambda\mu\nu = \frac{1}{2},\end{aligned}$$

and

$$\frac{2\|x_{k+1} - x_k\|}{1 + \sqrt{(1 + 4\theta^{2k}(1 + \theta^{2k})^{-2})}} \leq \|x_k - a\| \leq \theta^{2k-1} \|x_k - x_{k-1}\|.$$

We can now specialize Theorems 30.1 and 30.2 to the search of zeros of the derivative $f': \Omega \rightarrow E'$, of a function $f: \Omega \rightarrow \mathbb{R}$, with $\Omega \subseteq E$. The second derivative J'' of J is a continuous bilinear form $J'': E \times E \rightarrow \mathbb{R}$, but is convenient to view it as a linear map in $\mathcal{L}(E, E')$; the continuous linear form $J''(u)$ is given by $J''(u)(v) = J''(u, v)$. In our next theorem, we assume that the $A_k(x)$ are isomorphisms in $\mathcal{L}(E, E')$.

Theorem 30.4. *Let E be a Banach space, let $J: \Omega \rightarrow \mathbb{R}$ be twice differentiable on the open subset $\Omega \subseteq E$, and assume that there are constants $r, M, \beta > 0$ such that if we let*

$$B = \{x \in E \mid \|x - x_0\| \leq r\} \subseteq \Omega,$$

then

(1)

$$\sup_{k \geq 0} \sup_{x \in B} \|A_k^{-1}(x)\|_{\mathcal{L}(E'; E)} \leq M,$$

(2) $\beta < 1$ and

$$\sup_{k \geq 0} \sup_{x, x' \in B} \|J''(x) - A_k(x')\|_{\mathcal{L}(E; E')} \leq \frac{\beta}{M}$$

(3)

$$\|J'(x_0)\| \leq \frac{r}{M}(1 - \beta).$$

Then, the sequence (x_k) defined by

$$x_{k+1} = x_k - A_k^{-1}(x_\ell)(J'(x_k)), \quad 0 \leq \ell \leq k$$

is entirely contained within B and converges to a zero a of J' , which is the only zero of J' in B . Furthermore, the convergence is geometric, which means that

$$\|x_k - a\| \leq \frac{\|x_1 - x_0\|}{1 - \beta} \beta^k.$$

In the next theorem, we assume that the $A_k(x)$ are isomorphisms in $\mathcal{L}(E, E')$ that are independent of $x \in \Omega$.

Theorem 30.5. *Let E be a Banach space, and let $J: \Omega \rightarrow \mathbb{R}$ be twice differentiable on the open subset $\Omega \subseteq E$. If $a \in \Omega$ is a point such that $J'(a) = 0$, if $J''(a)$ is a linear isomorphism, and if there is some λ with $0 < \lambda < 1/2$ such that*

$$\sup_{k \geq 0} \|A_k - J''(a)\|_{\mathcal{L}(E; E')} \leq \frac{\lambda}{\|(J''(a))^{-1}\|_{\mathcal{L}(E'; E)}},$$

then there is a closed ball B of center a such that for every $x_0 \in B$, the sequence (x_k) defined by

$$x_{k+1} = x_k - A_k^{-1}(J'(x_k)), \quad k \geq 0,$$

is entirely contained within B and converges to a , which is the only zero of J' in B . Furthermore, the convergence is geometric, which means that

$$\|x_k - a\| \leq \beta^k \|x_0 - a\|,$$

for some $\beta < 1$.

When $E = \mathbb{R}^n$, the Newton method given by Theorem 30.4 yield an iteration step of the form

$$x_{k+1} = x_k - A_k^{-1}(x_\ell) \nabla J(x_k), \quad 0 \leq \ell \leq k,$$

where $\nabla J(x_k)$ is the gradient of J at x_k (here, we identify E' with \mathbb{R}^n). In particular, Newton's original method picks $A_k = J''$, and the iteration step is of the form

$$x_{k+1} = x_k - (\nabla^2 J(x_k))^{-1} \nabla J(x_k), \quad k \geq 0,$$

where $\nabla^2 J(x_k)$ is the Hessian of J at x_k .

As remarked in [24] (Section 7.5), generalized Newton methods have a very wide range of applicability. For example, various versions of gradient descent methods can be viewed as instances of Newton methods.

Newton's method also plays an important role in convex optimization, in particular, interior-point methods. A variant of Newton's method dealing with equality constraints has been developed. We refer the reader to Boyd and Vandenberghe [17], Chapters 10 and 11, for a comprehensive exposition of these topics.

30.3 Summary

The main concepts and results of this chapter are listed below:

- Newton's method for functions $f: \mathbb{R} \rightarrow \mathbb{R}$.
- Generalized Newton methods.
- The *Newton-Kantorovich* theorem.

Chapter 31

Appendix: Zorn's Lemma; Some Applications

31.1 Statement of Zorn's Lemma

Zorn's lemma is a particularly useful form of the axiom of choice, especially for algebraic applications. Readers who want to learn more about Zorn's lemma and its applications to algebra should consult either Lang [67], Appendix 2, §2 (pp. 878-884) and Chapter III, §5 (pp. 139-140), or Artin [4], Appendix §1 (pp. 588-589). For the logical ramifications of Zorn's lemma and its equivalence with the axiom of choice, one should consult Schwartz [93], (Vol. 1), Chapter I, §6, or a text on set theory such as Enderton [34], Suppes [107], or Kuratowski and Mostowski [66].

Given a set, S , a *partial order*, \leq , on S is a binary relation on S (i.e., $\leq \subseteq S \times S$) which is

- (1) *reflexive*, i.e., $x \leq x$, for all $x \in S$,
- (2) *transitive*, i.e, if $x \leq y$ and $y \leq z$, then $x \leq z$, for all $x, y, z \in S$, and
- (3) *antisymmetric*, i.e, if $x \leq y$ and $y \leq x$, then $x = y$, for all $x, y \in S$.

A pair (S, \leq) , where \leq is a partial order on S , is called a *partially ordered set* or *poset*. Given a poset, (S, \leq) , a subset, C , of S is *totally ordered* or a *chain* if for every pair of elements $x, y \in C$, either $x \leq y$ or $y \leq x$. The empty set is trivially a chain. A subset, P , (empty or not) of S is *bounded* if there is some $b \in S$ so that $x \leq b$ for all $x \in P$. Observe that the empty subset of S is bounded if and only if S is nonempty. A *maximal element* of P is an element, $m \in P$, so that $m \leq x$ implies that $m = x$, for all $x \in P$. Zorn's lemma can be stated as follows:

Lemma 31.1. *Given a partially ordered set, (S, \leq) , if every chain is bounded, then S has a maximal element.*

Proof. See any of Schwartz [93], Enderton [34], Suppes [107], or Kuratowski and Mostowski [66]. \square

Remark: As we noted, the hypothesis of Zorn's lemma implies that S is nonempty (since the empty set must be bounded). A partially ordered set such that every chain is bounded is sometimes called *inductive*.

We now give some applications of Zorn's lemma.

31.2 Proof of the Existence of a Basis in a Vector Space

Using Zorn's lemma, we can prove that Theorem 2.9 holds for arbitrary vector spaces, and not just for finitely generated vector spaces, as promised in Chapter 2.

Theorem 31.2. *Given any family, $S = (u_i)_{i \in I}$, generating a vector space E and any linearly independent subfamily, $L = (u_j)_{j \in J}$, of S (where $J \subseteq I$), there is a basis, B , of E such that $L \subseteq B \subseteq S$.*

Proof. Consider the set \mathcal{L} of linearly independent families, B , such that $L \subseteq B \subseteq S$. Since $L \in \mathcal{L}$, this set is nonempty. We claim that \mathcal{L} is inductive. Consider any chain, $(B_l)_{l \in \Lambda}$, of linearly independent families B_l in \mathcal{L} , and look at $B = \bigcup_{l \in \Lambda} B_l$. The family B is of the form $B = (v_h)_{h \in H}$, for some index set H , and it must be linearly independent. Indeed, if this was not true, there would be some family $(\lambda_h)_{h \in H}$ of scalars, of finite support, so that

$$\sum_{h \in H} \lambda_h v_h = 0,$$

where not all λ_h are zero. Since $B = \bigcup_{l \in \Lambda} B_l$ and only finitely many λ_h are nonzero, there is a finite subset, F , of Λ , so that $v_h \in B_{f_h}$ iff $\lambda_h \neq 0$. But $(B_l)_{l \in \Lambda}$ is a chain, and if we let $f = \max\{f_h \mid f_h \in F\}$, then $v_h \in B_f$, for all v_h for which $\lambda_h \neq 0$. Thus,

$$\sum_{h \in H} \lambda_h v_h = 0$$

would be a nontrivial linear dependency among vectors from B_f , a contradiction. Therefore, $B \in \mathcal{L}$, and since B is obviously an upper bound for the B_l 's, we have proved that \mathcal{L} is inductive. By Zorn's lemma (Lemma 31.1), the set \mathcal{L} has some maximal element, say $B = (u_h)_{h \in H}$. The rest of the proof is the same as in the proof of Theorem 2.9, but we repeat it for the reader's convenience. We claim that B generates E . Indeed, if B does not generate E , then there is some $u_p \in S$ that is not a linear combination of vectors in B (since S generates E), with $p \notin H$. Then, by Lemma 2.8, the family $B' = (u_h)_{h \in H \cup \{p\}}$ is linearly independent, and since $L \subseteq B \subset B' \subseteq S$, this contradicts the maximality of B . Thus, B is a basis of E such that $L \subseteq B \subseteq S$. \square

Another important application of Zorn's lemma is the existence of maximal ideals.

31.3 Existence of Maximal Ideals Containing a Given Proper Ideal

Let A be a commutative ring with identity element. Recall that an ideal \mathfrak{A} in A is a *proper ideal* if $\mathfrak{A} \neq A$. The following theorem holds:

Theorem 31.3. *Given any proper ideal, $\mathfrak{A} \subseteq A$, there is a maximal ideal, \mathfrak{B} , containing \mathfrak{A} .*

Proof. Let \mathcal{I} be the set of all proper ideals, \mathfrak{B} , in A that contain \mathfrak{A} . The set \mathcal{I} is nonempty, since $\mathfrak{A} \in \mathcal{I}$. We claim that \mathcal{I} is inductive. Consider any chain $(\mathfrak{A}_i)_{i \in I}$ of ideals \mathfrak{A}_i in A . One can easily check that $\mathfrak{B} = \bigcup_{i \in I} \mathfrak{A}_i$ is an ideal. Furthermore, \mathfrak{B} is a proper ideal, since otherwise, the identity element 1 would belong to $\mathfrak{B} = A$, and so, we would have $1 \in \mathfrak{A}_i$ for some i , which would imply $\mathfrak{A}_i = A$, a contradiction. Also, \mathfrak{B} is obviously an upper bound for all the \mathfrak{A}_i 's. By Zorn's lemma (Lemma 31.1), the set \mathcal{I} has a maximal element, say \mathfrak{B} , and \mathfrak{B} is a maximal ideal containing \mathfrak{A} . \square

Bibliography

- [1] Lars V. Ahlfors and Leo Sario. *Riemann Surfaces*. Princeton Math. Series, No. 2. Princeton University Press, 1960.
- [2] George E. Andrews, Richard Askey, and Ranjan Roy. *Special Functions*. Cambridge University Press, first edition, 2000.
- [3] Emil Artin. *Geometric Algebra*. Wiley Interscience, first edition, 1957.
- [4] Michael Artin. *Algebra*. Prentice Hall, first edition, 1991.
- [5] M. F. Atiyah and I. G. Macdonald. *Introduction to Commutative Algebra*. Addison Wesley, third edition, 1969.
- [6] A. Avez. *Calcul Différentiel*. Masson, first edition, 1991.
- [7] Sheldon Axler. *Linear Algebra Done Right*. Undergraduate Texts in Mathematics. Springer Verlag, second edition, 2004.
- [8] Marcel Berger. *Géométrie 1*. Nathan, 1990. English edition: Geometry 1, Universitext, Springer Verlag.
- [9] Marcel Berger. *Géométrie 2*. Nathan, 1990. English edition: Geometry 2, Universitext, Springer Verlag.
- [10] Marcel Berger and Bernard Gostiaux. *Géométrie différentielle: variétés, courbes et surfaces*. Collection Mathématiques. Puf, second edition, 1992. English edition: Differential geometry, manifolds, curves, and surfaces, GTM No. 115, Springer Verlag.
- [11] Rolf Berndt. *An Introduction to Symplectic Geometry*. Graduate Studies in Mathematics, Vol. 26. AMS, first edition, 2001.
- [12] J.E. Bertin. *Algèbre linéaire et géométrie classique*. Masson, first edition, 1981.
- [13] Nicolas Bourbaki. *Algèbre, Chapitre 9*. Eléments de Mathématiques. Hermann, 1968.
- [14] Nicolas Bourbaki. *Algèbre, Chapitres 1-3*. Eléments de Mathématiques. Hermann, 1970.

- [15] Nicolas Bourbaki. *Algèbre, Chapitres 4-7*. Éléments de Mathématiques. Masson, 1981.
- [16] Nicolas Bourbaki. *Espaces Vectoriels Topologiques*. Éléments de Mathématiques. Masson, 1981.
- [17] Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, first edition, 2004.
- [18] Glen E Bredon. *Topology and Geometry*. GTM No. 139. Springer Verlag, first edition, 1993.
- [19] G. Cagnac, E. Ramis, and J. Commeau. *Mathématiques Spéciales, Vol. 3, Géométrie*. Masson, 1965.
- [20] Henri Cartan. *Cours de Calcul Différentiel*. Collection Méthodes. Hermann, 1990.
- [21] Henri Cartan. *Differential Forms*. Dover, first edition, 2006.
- [22] Claude Chevalley. *The Algebraic Theory of Spinors and Clifford Algebras. Collected Works, Vol. 2*. Springer, first edition, 1997.
- [23] Yvonne Choquet-Bruhat, Cécile DeWitt-Morette, and Margaret Dillard-Bleick. *Analysis, Manifolds, and Physics, Part I: Basics*. North-Holland, first edition, 1982.
- [24] P.G. Ciarlet. *Introduction to Numerical Matrix Analysis and Optimization*. Cambridge University Press, first edition, 1989. French edition: Masson, 1994.
- [25] Timothée Cour and Jianbo Shi. Solving markov random fields with spectral relaxation. In Marita Meila and Xiaotong Shen, editors, *Artificial Intelligence and Statistics*. Society for Artificial Intelligence and Statistics, 2007.
- [26] H.S.M. Coxeter. *Introduction to Geometry*. Wiley, second edition, 1989.
- [27] James W. Demmel. *Applied Numerical Linear Algebra*. SIAM Publications, first edition, 1997.
- [28] Jean Dieudonné. *Algèbre Linéaire et Géométrie Élémentaire*. Hermann, second edition, 1965.
- [29] Jacques Dixmier. *General Topology*. UTM. Springer Verlag, first edition, 1984.
- [30] Manfredo P. do Carmo. *Differential Geometry of Curves and Surfaces*. Prentice Hall, 1976.
- [31] Manfredo P. do Carmo. *Riemannian Geometry*. Birkhäuser, second edition, 1992.
- [32] David S. Dummit and Richard M. Foote. *Abstract Algebra*. Wiley, second edition, 1999.

- [33] Gerald A. Edgar. *Measure, Topology, and Fractal Geometry*. Undergraduate Texts in Mathematics. Springer Verlag, first edition, 1992.
- [34] Herbert B. Enderton. *Elements of Set Theory*. Academic Press, 1997.
- [35] Charles L. Epstein. *Introduction to the Mathematics of Medical Imaging*. SIAM, second edition, 2007.
- [36] Gerald Farin. *Curves and Surfaces for CAGD*. Academic Press, fourth edition, 1998.
- [37] Olivier Faugeras. *Three-Dimensional Computer Vision, A geometric Viewpoint*. the MIT Press, first edition, 1996.
- [38] James Foley, Andries van Dam, Steven Feiner, and John Hughes. *Computer Graphics. Principles and Practice*. Addison-Wesley, second edition, 1993.
- [39] David A. Forsyth and Jean Ponce. *Computer Vision: A Modern Approach*. Prentice Hall, first edition, 2002.
- [40] Jean Fresnel. *Méthodes Modernes En Géométrie*. Hermann, first edition, 1998.
- [41] William Fulton. *Algebraic Topology, A first course*. GTM No. 153. Springer Verlag, first edition, 1995.
- [42] William Fulton and Joe Harris. *Representation Theory, A first course*. GTM No. 129. Springer Verlag, first edition, 1991.
- [43] Jean H. Gallier. *Curves and Surfaces In Geometric Modeling: Theory And Algorithms*. Morgan Kaufmann, 1999.
- [44] Jean H. Gallier. *Geometric Methods and Applications, For Computer Science and Engineering*. TAM, Vol. 38. Springer, second edition, 2011.
- [45] Walter Gander, Gene H. Golub, and Urs von Matt. A constrained eigenvalue problem. *Linear Algebra and its Applications*, 114/115:815–839, 1989.
- [46] F.R. Gantmacher. *The Theory of Matrices, Vol. I*. AMS Chelsea, first edition, 1977.
- [47] Roger Godement. *Cours d'Algèbre*. Hermann, first edition, 1963.
- [48] Gene H. Golub. Some modified eigenvalue problems. *SIAM Review*, 15(2):318–334, 1973.
- [49] H. Golub, Gene and F. Van Loan, Charles. *Matrix Computations*. The Johns Hopkins University Press, third edition, 1996.
- [50] A. Gray. *Modern Differential Geometry of Curves and Surfaces*. CRC Press, second edition, 1997.

- [51] Donald T. Greenwood. *Principles of Dynamics*. Prentice Hall, second edition, 1988.
- [52] Larry C. Grove. *Classical Groups and Geometric Algebra*. Graduate Studies in Mathematics, Vol. 39. AMS, first edition, 2002.
- [53] Jacques Hadamard. *Leçons de Géométrie Élémentaire. I Géométrie Plane*. Armand Colin, thirteenth edition, 1947.
- [54] Jacques Hadamard. *Leçons de Géométrie Élémentaire. II Géométrie dans l'Espace*. Armand Colin, eighth edition, 1949.
- [55] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, second edition, 2009.
- [56] D. Hilbert and S. Cohn-Vossen. *Geometry and the Imagination*. Chelsea Publishing Co., 1952.
- [57] Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, first edition, 1990.
- [58] Roger A. Horn and Charles R. Johnson. *Topics in Matrix Analysis*. Cambridge University Press, first edition, 1994.
- [59] Nathan Jacobson. *Basic Algebra I*. Freeman, second edition, 1985.
- [60] Ramesh Jain, Rangachar Katsuri, and Brian G. Schunck. *Machine Vision*. McGraw-Hill, first edition, 1995.
- [61] Jürgen Jost. *Riemannian Geometry and Geometric Analysis*. Universitext. Springer Verlag, fourth edition, 2005.
- [62] Hoffman Kenneth and Kunze Ray. *Linear Algebra*. Prentice Hall, second edition, 1971.
- [63] D. Kincaid and W. Cheney. *Numerical Analysis*. Brooks/Cole Publishing, second edition, 1996.
- [64] Anthony W. Knap. *Lie Groups Beyond an Introduction*. Progress in Mathematics, Vol. 140. Birkhäuser, second edition, 2002.
- [65] Erwin Kreyszig. *Differential Geometry*. Dover, first edition, 1991.
- [66] K. Kuratowski and A. Mostowski. *Set Theory*. Studies in Logic, Vol. 86. Elsevier, 1976.
- [67] Serge Lang. *Algebra*. Addison Wesley, third edition, 1993.
- [68] Serge Lang. *Differential and Riemannian Manifolds*. GTM No. 160. Springer Verlag, third edition, 1995.

- [69] Serge Lang. *Real and Functional Analysis*. GTM 142. Springer Verlag, third edition, 1996.
- [70] Serge Lang. *Undergraduate Analysis*. UTM. Springer Verlag, second edition, 1997.
- [71] Peter Lax. *Linear Algebra and Its Applications*. Wiley, second edition, 2007.
- [72] N. N. Lebedev. *Special Functions and Their Applications*. Dover, first edition, 1972.
- [73] Saunders Mac Lane and Garrett Birkhoff. *Algebra*. Macmillan, first edition, 1967.
- [74] Ib Madsen and Jorgen Tornehave. *From Calculus to Cohomology. De Rham Cohomology and Characteristic Classes*. Cambridge University Press, first edition, 1998.
- [75] M.-P. Malliavin. *Algèbre Commutative. Applications en Géométrie et Théorie des Nombres*. Masson, first edition, 1985.
- [76] Jerrold E. Marsden and J.R. Hughes, Thomas. *Mathematical Foundations of Elasticity*. Dover, first edition, 1994.
- [77] William S. Massey. *Algebraic Topology: An Introduction*. GTM No. 56. Springer Verlag, second edition, 1987.
- [78] William S. Massey. *A Basic Course in Algebraic Topology*. GTM No. 127. Springer Verlag, first edition, 1991.
- [79] Dimitris N. Metaxas. *Physics-Based Deformable Models*. Kluwer Academic Publishers, first edition, 1997.
- [80] Carl D. Meyer. *Matrix Analysis and Applied Linear Algebra*. SIAM, first edition, 2000.
- [81] John W. Milnor. *Topology from the Differentiable Viewpoint*. The University Press of Virginia, second edition, 1969.
- [82] John W. Milnor and James D. Stasheff. *Characteristic Classes*. Annals of Math. Series, No. 76. Princeton University Press, first edition, 1974.
- [83] Shigeyuki Morita. *Geometry of Differential Forms*. Translations of Mathematical Monographs No 201. AMS, first edition, 2001.
- [84] James R. Munkres. *Topology, a First Course*. Prentice Hall, first edition, 1975.
- [85] James R. Munkres. *Analysis on Manifolds*. Addison Wesley, 1991.
- [86] Ivan Niven, Herbert S. Zuckerman, and Hugh L. Montgomery. *An Introduction to the Theory of Numbers*. Wiley, fifth edition, 1991.

- [87] Joseph O'Rourke. *Computational Geometry in C*. Cambridge University Press, second edition, 1998.
- [88] Dan Pedoe. *Geometry, A comprehensive Course*. Dover, first edition, 1988.
- [89] Eugène Rouché and Charles de Comberousse. *Traité de Géométrie*. Gauthier-Villars, seventh edition, 1900.
- [90] Pierre Samuel. *Projective Geometry*. Undergraduate Texts in Mathematics. Springer Verlag, first edition, 1988.
- [91] Giovanni Sansone. *Orthogonal Functions*. Dover, first edition, 1991.
- [92] Laurent Schwartz. *Topologie Générale et Analyse Fonctionnelle*. Collection Enseignement des Sciences. Hermann, 1980.
- [93] Laurent Schwartz. *Analyse I. Théorie des Ensembles et Topologie*. Collection Enseignement des Sciences. Hermann, 1991.
- [94] Laurent Schwartz. *Analyse II. Calcul Différentiel et Equations Différentielles*. Collection Enseignement des Sciences. Hermann, 1992.
- [95] H. Seifert and W. Threlfall. *A Textbook of Topology*. Academic Press, first edition, 1980.
- [96] Denis Serre. *Matrices, Theory and Applications*. GTM No. 216. Springer Verlag, second edition, 2010.
- [97] Jean-Pierre Serre. *A Course in Arithmetic*. Graduate Text in Mathematics, No. 7. Springer, first edition, 1973.
- [98] Igor R. Shafarevich. *Basic Algebraic Geometry 1*. Springer Verlag, second edition, 1994.
- [99] Ernst Snapper and Troyer Robert J. *Metric Affine Geometry*. Dover, first edition, 1989.
- [100] Harold M. Stark. *An Introduction to Number Theory*. MIT Press, first edition, 1994. Eighth Printing.
- [101] G.W. Stewart. On the early history of the singular value decomposition. *SIAM review*, 35(4):551–566, 1993.
- [102] J.J. Stoker. *Differential Geometry*. Wiley Classics. Wiley-Interscience, first edition, 1989.
- [103] Eric J. Stollnitz, Tony D. DeRose, and David H. Salesin. *Wavelets for Computer Graphics Theory and Applications*. Morgan Kaufmann, first edition, 1996.

- [104] Gilbert Strang. *Introduction to Applied Mathematics*. Wellesley-Cambridge Press, first edition, 1986.
- [105] Gilbert Strang. *Linear Algebra and its Applications*. Saunders HBJ, third edition, 1988.
- [106] Gilbert Strang and Nguyen Truong. *Wavelets and Filter Banks*. Wellesley-Cambridge Press, second edition, 1997.
- [107] Patrick Suppes. *Axiomatic Set Theory*. Dover, 1972.
- [108] Donald E. Taylor. *The Geometry of the Classical Groups*. Sigma Series in Pure Mathematics, Vol. 9. Heldermann Verlag Berlin, 1992.
- [109] Claude Tisseron. *Géométries affines, projectives, et euclidiennes*. Hermann, first edition, 1994.
- [110] L.N. Trefethen and D. Bau III. *Numerical Linear Algebra*. SIAM Publications, first edition, 1997.
- [111] Emanuele Trucco and Alessandro Verri. *Introductory Techniques for 3D Computer Vision*. Prentice-Hall, first edition, 1998.
- [112] B.L. Van Der Waerden. *Algebra, Vol. 1*. Ungar, seventh edition, 1973.
- [113] J.H. van Lint and R.M. Wilson. *A Course in Combinatorics*. Cambridge University Press, second edition, 2001.
- [114] Frank Warner. *Foundations of Differentiable Manifolds and Lie Groups*. GTM No. 94. Springer Verlag, first edition, 1983.
- [115] Ernst Witt. Theorie der quadratischen Formen in beliebigen Körpern. *J. Reine Angew. Math.*, 176:31–44, 1936.
- [116] Stella X. Yu and Jianbo Shi. Grouping with bias. In Thomas G. Dietterich, Sue Becker, and Zoubin Ghahramani, editors, *Neural Information Processing Systems, Vancouver, Canada, 3-8 Dec. 2001*. MIT Press, 2001.
- [117] Oscar Zariski and Pierre Samuel. *Commutative Algebra, Vol I*. GTM No. 28. Springer Verlag, first edition, 1975.