

Appendix to Lecture Notes on *Calculus II for Statistics Students*

Peter Philip

May 18, 2015

Contents

A	Linear Algebra	3
A.1	Vector Spaces	3
A.2	Linear Maps	9
A.3	Matrices	13
A.4	Determinants	18
B	Metric Spaces	26
B.1	Metric Subspaces	26
B.2	Norm-Preserving and Isometric Maps	28
B.3	Uniform Continuity and Lipschitz Continuity	29
B.4	Viewing \mathbb{C}^n as \mathbb{R}^{2n}	31
B.5	Banach Fixed Point Theorem	33
B.6	Unit Balls in Normed Spaces	34
C	Differential Calculus in \mathbb{R}^n	36
C.1	Proof of the Chain Rule	36
C.2	Bounded Derivatives Imply Lipschitz Continuity	38
C.3	Surjectivity of Directional Derivatives	40
C.4	Implicit Function Theorem	41
D	Riemann Integral for \mathbb{C}-Valued Functions	46
D.1	Riemann Integrability	46

<i>CONTENTS</i>	2
D.2 Fubini Theorem	49
D.3 Change of Variables	50
References	51

A Linear Algebra

A.1 Vector Spaces

In [Phi15a], we encountered the abstract definition of a field in [Phi15a, Def. 4.4], and we studied the fields \mathbb{Q} , \mathbb{R} , and \mathbb{C} . Even though we will formulate the following definitions and results using a general field F as defined in [Phi15a, Def. 4.4], for the purposes of the present lecture, you may always think of F as being the field of real numbers \mathbb{R} or the field of complex numbers \mathbb{C} .

Definition A.1. Let F be a field and let V be a nonempty set with two maps

$$\begin{aligned} + : V \times V &\longrightarrow V, & (x, y) &\mapsto x + y, \\ \cdot : F \times V &\longrightarrow V, & (\lambda, x) &\mapsto \lambda \cdot x \end{aligned} \tag{A.1}$$

($+$ is called *(vector) addition* and \cdot is called *scalar multiplication*; often one writes xy instead of $x \cdot y$ – please take care not to confuse the vector addition on V with the addition on F and, likewise, not to confuse the scalar multiplication with the multiplication on F , the symbol $+$ is used for both additions and \cdot is used for both multiplications, but you can always determine from the context which addition or multiplication is meant). Then V is called a *vector space* or a *linear space* over F (sometimes also an F -vector space) if, and only if, the following conditions are satisfied:

(i) V is a commutative group with respect to $+$. The neutral element with respect to $+$ is denoted 0 (do not confuse $0 \in V$ with $0 \in F$ – once again, the same symbol is used for different objects (both objects only coincide for $F = V$)).

(ii) Distributivity:

$$\forall_{\lambda \in F} \quad \forall_{x, y \in V} \quad \lambda(x + y) = \lambda x + \lambda y, \tag{A.2a}$$

$$\forall_{\lambda, \mu \in F} \quad \forall_{x \in V} \quad (\lambda + \mu)x = \lambda x + \mu x. \tag{A.2b}$$

(iii) Compatibility between Multiplication on F and Scalar Multiplication:

$$\forall_{\lambda, \mu \in F} \quad \forall_{x \in V} \quad (\lambda\mu)x = \lambda(\mu x). \tag{A.3}$$

(iv) The neutral element with respect to the multiplication on F is also neutral with respect to the scalar multiplication:

$$\forall_{x \in V} \quad 1x = x. \tag{A.4}$$

If V is a vector space over F , then one calls the elements of V *vectors* and the elements of F *scalars*.

Example A.2. (a) Every field F is a vector space over itself if one uses the field addition in F as the vector addition and the field multiplication in F as the scalar multiplication (as important special cases, we obtain that \mathbb{R} is a vector space over \mathbb{R} and \mathbb{C} is a vector space over \mathbb{C}): All the vector space laws are immediate consequences of the corresponding field laws: Def. A.1(i) holds as every field is a commutative group with respect to addition; Def. A.1(ii) follows from the field distributivity [Phi15a, Def. (4.5)] and multiplicative commutativity on F ; Def. A.1(iii) is merely the multiplicative associativity on F ; and Def. A.1(iv) holds, since scalar multiplication coincides with field multiplication on F .

(b) The reasoning in (a) actually shows that every field F is a vector space over every subfield E of F (over every $E \subseteq F$ that is a field with respect to the addition and multiplication defined on F). In particular, \mathbb{R} is a vector space over \mathbb{Q} .

(c) If A is any nonempty set, F is a field, and Y is a vector space over the field F , then we can make $V := \mathcal{F}(A, Y) = Y^A$ (the set of functions from A into Y) into a vector space over F by defining for each $f, g : A \rightarrow Y$:

$$(f + g) : A \rightarrow Y, \quad (f + g)(x) := f(x) + g(x), \quad (\text{A.5a})$$

$$(\lambda \cdot f) : A \rightarrow Y, \quad (\lambda \cdot f)(x) := \lambda \cdot f(x) \quad \text{for each } \lambda \in F \quad (\text{A.5b})$$

(note that, for $Y = F = \mathbb{K}$, the above definitions are the same as the ones in [Phi15a, (6.1a)] and [Phi15a, (6.1b)], respectively).

It is an exercise to verify that $(V, +, \cdot)$ is, indeed, a vector space over F .

(d) For each $n \in \mathbb{N}$, $(\mathbb{K}^n, +, \cdot)$, with vector addition and scalar multiplication as defined in (1.1a) and (1.1c), respectively, constitutes a vector space over \mathbb{K} . The validity of Def. A.1(i) – Def. A.1(iv) can easily be verified directly, but $(\mathbb{K}^n, +, \cdot)$ can also be seen as a special case of (c) with $A = \{1, \dots, n\}$ and $Y = F = \mathbb{K}$. To this end, recall that, according to [Phi15a, Ex. 2.15(c)], $\mathbb{K}^n = \mathbb{K}^{\{1, \dots, n\}} = \mathcal{F}(\{1, \dots, n\}, \mathbb{K})$ is the set of functions from $\{1, \dots, n\}$ into \mathbb{K} . Then $z = (z_1, \dots, z_n) \in \mathbb{K}^n$ is the same as the function $f : \{1, \dots, n\} \rightarrow \mathbb{K}$, $f(j) = z_j$. Thus, (1.1a) is, indeed, the same as (A.5a), and (1.1c) is, indeed, the same as (A.5b).

Definition and Remark A.3. Let $(V, +, \cdot)$ be a vector space over the field F . A subset $U \subseteq V$ is called a *subspace* of V if, and only if, U is a vector space over F with respect to operations $+$ and \cdot it inherits from V . Clearly, every law (commutativity, associativity, etc.) that holds on V must, in particular, hold on U , showing that $\emptyset \neq U \subseteq V$ is a subspace of V if, and only if,

$$\forall_{x, y \in U} \quad x + y \in U, \quad (\text{A.6a})$$

$$\wedge \quad \forall_{\lambda \in F} \quad \forall_{x \in U} \quad \lambda x \in U. \quad (\text{A.6b})$$

which holds if, and only if,

$$\forall_{\lambda, \mu \in F} \quad \forall_{x, y \in U} \quad \lambda x + \mu y \in U. \quad (\text{A.7})$$

Example A.4. (a) \mathbb{Q} is *not* a subspace of \mathbb{R} if \mathbb{R} is considered as a vector space over \mathbb{R} (for example, $\sqrt{2} \cdot 2 \notin \mathbb{Q}$). However, \mathbb{Q} *is* a subspace of \mathbb{R} if \mathbb{R} is considered as a vector space over \mathbb{Q} .

(b) From Ex. A.2(c), we know that, for each $\emptyset \neq A$, $\mathcal{F}(A, \mathbb{K})$ constitutes a vector space over \mathbb{K} . Thus, as a consequence of Def. and Rem. A.3, a subset of $\mathcal{F}(A, \mathbb{K})$ is a vector space over \mathbb{K} if, and only if, it is closed under addition and scalar multiplication. By using results from [Phi15a], we obtain the following examples:

- (i) The set $\mathcal{P}(\mathbb{K})$ of all polynomials mapping from \mathbb{K} into \mathbb{K} is a vector space over \mathbb{K} by [Phi15a, Rem. 6.4]; for each $n \in \mathbb{N}_0$, the set $\mathcal{P}_n(\mathbb{K})$ of all such polynomials of degree at most n is also a vector space over \mathbb{K} by [Phi15a, Rem. (6.4a),(6.4b)].
- (ii) If $\emptyset \neq M \subseteq \mathbb{C}$, then the set of continuous functions from M into \mathbb{K} , i.e. $C(M, \mathbb{K})$, is a vector space over \mathbb{K} by [Phi15a, Th. 7.38].
- (iii) If $a, b \in \mathbb{R} \cup \{-\infty, \infty\}$ and $a < b$, then the set of differentiable functions from $I :=]a, b[$ into \mathbb{K} is a vector space over \mathbb{K} by [Phi15a, Th. 9.6(a),(b)]. Moreover, [Phi15a, Th. 9.6(a),(b)] also implies that, for each $k \in \mathbb{N}$, the set of k times differentiable functions from I into \mathbb{K} is a vector space over \mathbb{K} , and so is each set $C^k(I, \mathbb{K})$ of k times continuously differentiable functions ([Phi15a, Th. 7.38] is also needed for the last conclusion). The set $C^\infty(I, \mathbb{K})$ is also a vector space over \mathbb{K} by Th. A.5(a) below.

Theorem A.5. *Let V be a vector space over the field F .*

- (a)** *Let $I \neq \emptyset$ be an index set and $(U_i)_{i \in I}$ a family of subspaces of V . Then the intersection $U := \bigcap_{i \in I} U_i$ is also a subspace of V .*
- (b)** *In contrast to intersections, unions of subspaces are almost never subspaces. More precisely, if U_1 and U_2 are subspaces of V , then*

$$U_1 \cup U_2 \text{ is subspace of } V \quad \Leftrightarrow \quad (U_1 \subseteq U_2 \vee U_2 \subseteq U_1). \quad (\text{A.8})$$

Proof. See, e.g., [Str08, Th. 8.7]. ■

Definition A.6. Let V be a vector space over the field F .

- (a)** Let $n \in \mathbb{N}$ and $v_1, \dots, v_n \in V$. A vector $v \in V$ is called a *linear combination* of v_1, \dots, v_n if, and only if, there exist $\lambda_1, \dots, \lambda_n \in F$ such that

$$v = \sum_{i=1}^n \lambda_i v_i. \quad (\text{A.9})$$

Moreover, $v \in V$ is called *linearly dependent* of a subset U of V if, and only if, there exists $n \in \mathbb{N}$ and $u_1, \dots, u_n \in U$ such that v is a linear combination of u_1, \dots, u_n .

(b) A subset U of V is called *linearly independent* if, and only if,

$$\left(n \in \mathbb{N} \quad \wedge \quad u_1, \dots, u_n \in U \quad \wedge \quad \lambda_1, \dots, \lambda_n \in F \quad \wedge \quad \sum_{i=1}^n \lambda_i u_i = 0 \right) \\ \Rightarrow \quad \lambda_1 = \dots = \lambda_n = 0. \quad (\text{A.10a})$$

Occasionally, one also wants to have the notion available for families of vectors instead of sets, and one calls a family $(u_i)_{i \in I}$ of vectors in V *linearly independent* if, and only if,

$$\left(n \in \mathbb{N} \quad \wedge \quad i_1, \dots, i_n \in I \quad \wedge \quad \lambda_1, \dots, \lambda_n \in F \quad \wedge \quad \sum_{k=1}^n \lambda_k u_{i_k} = 0 \right) \\ \Rightarrow \quad \lambda_1 = \dots = \lambda_n = 0. \quad (\text{A.10b})$$

Sets and families that are not linearly independent are called *linearly dependent*.

Definition A.7. Let V be a vector space over the field F , $A \subseteq V$, and

$$\mathcal{U} := \{U \in \mathcal{P}(V) : A \subseteq U \wedge U \text{ is subspace of } V\}. \quad (\text{A.11})$$

Then the set

$$\langle A \rangle := \text{span } A := \bigcap_{U \in \mathcal{U}} U \quad (\text{A.12})$$

is called the *span* of A . Moreover A is called a *spanning set* of V if, and only if, $\langle A \rangle = V$.

Proposition A.8. Let V be a vector space over the field F and $A \subseteq V$.

- (a) $\langle A \rangle$ is a subspace of V , namely the smallest subspace of V containing A .
- (b) If $A = \emptyset$, then $\langle A \rangle = \{0\}$; if $A \neq \emptyset$, then $\langle A \rangle$ is the set of all linear combinations of elements from A , i.e.

$$\langle A \rangle = \left\{ \sum_{i=1}^n \lambda_i a_i : n \in \mathbb{N} \wedge \lambda_1, \dots, \lambda_n \in F \wedge a_1, \dots, a_n \in A \right\}. \quad (\text{A.13})$$

Proof. (a) is immediate from Th. A.5(a).

(b): For the case $A = \emptyset$, note that $\{0\}$ is a subspace of V , and that $\{0\}$ is contained in every subspace of V . For $A \neq \emptyset$, let W denote the right-hand side of (A.13), and recall from (A.12) that $\langle A \rangle$ is the intersection of all subspaces U of V that contain A . If U is a subspace of V and $A \subseteq U$, then $W \subseteq U$, since U is closed under vector addition and scalar multiplication, showing $W \subseteq \langle A \rangle$. On the other hand, W is clearly a subspace of V that contains A , showing $\langle A \rangle \subseteq W$, completing the proof of $\langle A \rangle = W$. ■

Definition A.9. Let V be a vector space over the field F and $B \subseteq V$.

- (a) B is called a *generating set* for V if, and only if, $\langle B \rangle = V$. One then also says that V is *generated* or *spanned* by B .
- (b) B is called a *basis* for V if, and only if, B is a generating set for V that is also linearly independent (see Def. A.6(b)).

Theorem A.10. *Let V be a vector space over the field F and $B \subseteq V$. Then the following statements (i) – (iii) are equivalent:*

- (i) B is a basis for V .
- (ii) B is a maximal linearly independent subset of V , i.e. B is linearly independent and each set $A \subseteq V$ with $B \subsetneq A$ is linearly dependent.
- (iii) B is a minimal generating set for V , i.e. $\langle B \rangle = V$ and $\langle A \rangle \subsetneq V$ for each $A \subsetneq B$.

Proof. See, e.g., [Str08, Th. 9.6]. ■

Theorem A.11 (Coordinates). *Let V be a vector space over the field F and assume $B \subseteq V$ is a basis of V . Then each vector $v \in V$ has unique coordinates with respect to the basis B , i.e., for each $v \in V$, there exists a unique finite subset B_v of B and a unique map $c : B_v \rightarrow F \setminus \{0\}$ such that*

$$v = \sum_{b \in B_v} c(b) b. \quad (\text{A.14})$$

Note that, for $v = 0$, one has $B_v = \emptyset$, c is the empty map, and (A.14) becomes $0 = \sum_{b \in \emptyset} c(b) b$, employing the useful convention that sums over the empty set are defined as 0.

Proof. The existence of B_v and the map c follows from the fact that the basis B is a generating set, $\langle B \rangle = V$. For the uniqueness proof, consider finite sets $B_v, \tilde{B}_v \subseteq B$ and maps $c : B_v \rightarrow F \setminus \{0\}$, $\tilde{c} : \tilde{B}_v \rightarrow F \setminus \{0\}$ such that

$$v = \sum_{b \in B_v} c(b) b = \sum_{b \in \tilde{B}_v} \tilde{c}(b) b. \quad (\text{A.15})$$

Extend both c and \tilde{c} to $A := B_v \cup \tilde{B}_v$ by letting $c(b) := 0$ for $b \in \tilde{B}_v \setminus B_v$ and $\tilde{c}(b) := 0$ for $b \in B_v \setminus \tilde{B}_v$. Then

$$0 = \sum_{b \in A} (c(b) - \tilde{c}(b)) b, \quad (\text{A.16})$$

such that the linear independence of A implies $c(b) = \tilde{c}(b)$ for each $b \in A$, which, in turn, implies $B_v = \tilde{B}_v$ and $c = \tilde{c}$. ■

Remark A.12. If the basis B of V has finitely many elements, then one often enumerates the elements $B = \{b_1, \dots, b_n\}$, $n = \#B \in \mathbb{N}$, and writes $\lambda_i := c(b_i)$ for $b_i \in B_v$, $\lambda_i := 0$ for $b_i \notin B_v$, such that (A.14) takes the, perhaps more familiar looking, form

$$v = \sum_{i=1}^n \lambda_i b_i. \quad (\text{A.17})$$

Theorem A.13. *Every vector space V over a field F has a basis $B \subseteq V$. Moreover, bases of V have a unique cardinality, i.e. if $B \subseteq V$ and $\tilde{B} \subseteq V$ are both bases of V , then $\#B = \#\tilde{B}$.*

Proof. See, e.g., [Str08, Lem. 11.3, Th. 11.5]. ■

Definition A.14. According to Th. A.13, for each vector space V over a field F , the cardinality of its bases is unique. This unique cardinality is called the *dimension* of V and is denoted $\dim V$. If $\dim V < \infty$ (i.e. $\dim V \in \mathbb{N}_0$), then V is called *finite dimensional*, otherwise *infinite dimensional*.

Example A.15. Given a field F and a nonempty set I , let F_{fin}^I denote the set of functions $f : I \rightarrow F$ such that there exists a finite set $I_f \subseteq I$ satisfying

$$f(i) = 0 \quad \text{for each } i \in I \setminus I_f, \quad (\text{A.18a})$$

$$f(i) \neq 0 \quad \text{for each } i \in I_f. \quad (\text{A.18b})$$

Then $F_{\text{fin}}^I = F^I$ if, and only if, I is finite (for example $F_{\text{fin}}^n = F^n$ for $n \in \mathbb{N}$); in general F_{fin}^I is a strict subset of F^I . However, if $f, g \in F_{\text{fin}}^I$ and $\lambda \in F$, then $I_{\lambda f} = I_f$ for $\lambda \neq 0$, $I_{\lambda f} = \emptyset$ for $\lambda = 0$, and $I_{f+g} \subseteq I_f \cup I_g$, showing F_{fin}^I is always a subspace of F^I . Define

$$e_i : I \rightarrow F, \quad e_i(j) := \begin{cases} 1 & \text{if } j = i, \\ 0 & \text{if } j \neq i. \end{cases} \quad (\text{A.19})$$

Then $I_{e_i} = \{i\}$ for each $i \in I$, in particular, $e_i \in F_{\text{fin}}^I$ for each $i \in I$. We claim that $B := \{e_i : i \in I\}$ is a basis for F_{fin}^I . Indeed, if $f \in F_{\text{fin}}^I$, then

$$f = \sum_{i \in I_f} f(i)e_i, \quad (\text{A.20})$$

showing $\langle B \rangle = F_{\text{fin}}^I$. If J is a finite subset of I and $(\lambda_j)_{j \in J}$ is a family in F such that

$$\sum_{j \in J} \lambda_j e_j \equiv 0, \quad (\text{A.21})$$

then

$$\forall_{j \in J} \quad \lambda_j \stackrel{(\text{A.19})}{=} \sum_{j \in J} \lambda_j e_j(j) \stackrel{(\text{A.21})}{=} 0, \quad (\text{A.22})$$

proving B is linearly independent. Clearly, $\#B = \#I$, so we have shown

$$\dim F_{\text{fin}}^I = \#I. \quad (\text{A.23})$$

In particular, we have shown that, for each $n \in \mathbb{N}$, the set $\{e_j : j = 1, \dots, n\}$, where

$$e_1 := (1, 0, \dots, 0), \quad e_2 := (0, 1, \dots, 0), \quad \dots, \quad e_n := (0, \dots, 0, 1), \quad (\text{A.24})$$

forms a basis of F^n (of \mathbb{R}^n if $F = \mathbb{R}$ and of \mathbb{C}^n if $F = \mathbb{C}$),

$$\dim F^n = \dim \mathbb{R}^n = \dim \mathbb{C}^n = n. \quad (\text{A.25})$$

Remark A.16. We will see in Th. A.24 below that, in a certain sense, F_{fin}^I is the only vector space of dimension $\#I$ over F . In particular, for $n \in \mathbb{N}$, you can think of \mathbb{K}^n as the *standard model* of an n -dimensional vector space over \mathbb{K} .

A.2 Linear Maps

Definition A.17. Let V and W be vector spaces over the field F .

- (a) A map $A : V \longrightarrow W$ is called F -linear (or merely *linear* if the field F is understood) if, and only if,

$$\forall_{v_1, v_2 \in V} \quad A(v_1 + v_2) = A(v_1) + A(v_2), \quad (\text{A.26a})$$

$$\wedge \quad \forall_{\lambda \in F} \quad \forall_{v \in V} \quad A(\lambda v) = \lambda A(v) \quad (\text{A.26b})$$

or, equivalently, if, and only if,

$$\forall_{\lambda, \mu \in F} \quad \forall_{v_1, v_2 \in V} \quad A(\lambda v_1 + \mu v_2) = \lambda A(v_1) + \mu A(v_2) \quad (\text{A.27})$$

(note that, in general, vector addition and scalar multiplication will be different on the left-hand sides and right-hand sides of the above equations).

One also calls linear maps (vector space) *homomorphisms*. We denote the set of all F -linear maps from V into W by $\mathcal{L}(V, W)$.

- (b) A linear map $I : V \longrightarrow W$ is called a (vector space or linear) *isomorphism* if, and only if, it is bijective (i.e. invertible). The vector spaces V and W are called *isomorphic* (denoted $V \cong W$) if, and only if, there exists a vector space isomorphism $I : V \longrightarrow W$.

Theorem A.18. Let V and W be vector spaces over the field F . If $I : V \longrightarrow W$ is a linear isomorphism, then so is $I^{-1} : W \longrightarrow V$ (i.e. I^{-1} is not only bijective, but also linear).

Proof. See, e.g., [Str08, Th. 13.5]. ■

Definition A.19. Let V and W be vector spaces over the field F , and $A \in \mathcal{L}(V, W)$. Define the *kernel* and the *image* of A by

$$\ker A := A^{-1}\{0\} = \{v \in V : A(v) = 0\}, \quad (\text{A.28a})$$

$$\text{Im } A := A(V) = \{A(v) : v \in V\}, \quad (\text{A.28b})$$

respectively.

Theorem A.20. Let V and W be vector spaces over the field F , and $A \in \mathcal{L}(V, W)$.

- (a) $\ker A$ is a subspace of V and $\text{Im } A$ is a subspace of W .
 (b) A is injective if, and only if, $\ker A = \{0\}$.

Proof. (a): See, e.g., [Str08, Lem. 12.5].

(b): Since $A(0) = 0$, A being injective implies $\ker A = \{0\}$. Conversely, assume $\ker A = \{0\}$ and $A(v_1) = A(v_2)$ for $v_1, v_2 \in V$. Then $A(-v_1 + v_2) = -A(v_1) + A(v_2) = -A(v_1) + A(v_1) = 0$, i.e. $-v_1 + v_2 \in \ker A$, i.e. $-v_1 + v_2 = 0$, showing $v_1 = v_2$ and the injectivity of A . ■

Theorem A.21 (Dimension Formulas). *Let V and W be vector spaces over the field F , and let $A : V \longrightarrow W$ be linear.*

- (a) *If V is finite dimensional, then $\dim V = \dim \ker A + \dim \operatorname{Im} A$.*
- (b) *If V is finite dimensional, then $\dim \operatorname{Im} A \leq \dim V$.*
- (c) *If W is finite dimensional, then $\dim \operatorname{Im} A \leq \dim W$.*

Proof. See, e.g., [Str08, Th. 12.12]. ■

Proposition A.22. *Let V and W be vector spaces over the field F , and let $A : V \longrightarrow W$ be linear.*

- (a) *A is injective if, and only if, for each linearly independent subset S of V , $A(S)$ is a linearly independent subset of W .*
- (b) *A is surjective if, and only if, for each generating subset S of V , $A(S)$ is a generating subset of W .*
- (c) *A is bijective if, and only if, for each basis B of V , $A(B)$ is a basis of W .*

Proof. (a): If A is not injective, then, according to Th. A.20(b), there exists $0 \neq v \in V$ such that $A(v) = 0$. Then $S := \{v\}$ is linearly independent, whereas $A(S) = \{0\}$ is not. Conversely, if A is injective, $S \subseteq V$ is linearly independent, and $\lambda_1, \dots, \lambda_n \in F$; $s_1, \dots, s_n \in S$; $n \in \mathbb{N}$; such that

$$0 = \sum_{i=1}^n \lambda_i A(s_i) = A \left(\sum_{i=1}^n \lambda_i s_i \right), \quad (\text{A.29})$$

then $\sum_{i=1}^n \lambda_i s_i = 0$ by Th. A.20(b), implying $\lambda_1 = \dots = \lambda_n = 0$, showing that $A(S)$ is also linearly independent.

(b): If A is not surjective, then $\langle A(V) \rangle = \operatorname{Im} A \neq W$, since $\operatorname{Im} A$ is a subspace of W . Conversely, if A is surjective, $S \subseteq V$, $\langle S \rangle = V$, and $w \in W$, then there are $v \in V$; $\lambda_1, \dots, \lambda_n \in F$; $s_1, \dots, s_n \in S$; $n \in \mathbb{N}$; such that $A(v) = w$ and $v = \sum_{i=1}^n \lambda_i s_i$, i.e. $w = A(v) = \sum_{i=1}^n \lambda_i A(s_i)$, proving $w \in \langle A(S) \rangle$. Since $w \in W$ was arbitrary, we have shown $\langle A(S) \rangle = W$.

(c) follows immediately by combining (a) and (b) (recalling that a basis is a linearly independent generating set). ■

Theorem A.23. *Let V and W be vector spaces over the field F . Then each linear map $A : V \longrightarrow W$ is uniquely determined by its values on a basis of V . More precisely, if B is a basis of V , $(w_b)_{b \in B}$ is a family in W , and, for each $v \in V$, B_v and $c_v : B_v \longrightarrow F \setminus \{0\}$ are as in Th. A.11 (we now write c_v instead of c to underline the dependence of c on v), then the map*

$$A : V \longrightarrow W, \quad A(v) = A \left(\sum_{b \in B_v} c_v(b) b \right) := \sum_{b \in B_v} c_v(b) w_b, \quad (\text{A.30})$$

is linear, and $\tilde{A} \in \mathcal{L}(V, W)$ with

$$\forall_{b \in B} \quad \tilde{A}(b) = w_b, \quad (\text{A.31})$$

implies $A = \tilde{A}$.

Proof. We first verify A is linear. Let $v \in V$ and $\lambda \in F$. If $\lambda = 0$, then $A(\lambda v) = A(0) = 0 = \lambda A(v)$. If $\lambda \neq 0$, then $B_{\lambda v} = B_v$, $c_{\lambda v} = \lambda c_v$, and

$$A(\lambda v) = A\left(\sum_{b \in B_{\lambda v}} c_{\lambda v}(b) b\right) = \sum_{b \in B_v} \lambda c_v(b) w_b = \lambda A\left(\sum_{b \in B_v} c_v(b) b\right) = \lambda A(v). \quad (\text{A.32a})$$

Now let $u, v \in V$. If $u = 0$, then $A(u + v) = A(v) = 0 + A(v) = A(u) + A(v)$, and analogously if $v = 0$. So assume $u, v \neq 0$. If $u + v = 0$, then $v = -u$ and $A(u + v) = A(0) = 0 = A(u) + A(-u) = A(u) + A(v)$. If $u + v \neq 0$, then $B_{u+v} \subseteq B_u \cup B_v$ and c_{u+v} can be extended to $B_u \cup B_v$ by letting

$$c_{u+v}(b) := \begin{cases} c_u(b) + c_v(b) & \text{if } b \in B_u \cap B_v, \\ c_u(b) & \text{if } b \in B_u \setminus B_v, \\ c_v(b) & \text{if } b \in B_v \setminus B_u. \end{cases} \quad (\text{A.32b})$$

One then obtains

$$\begin{aligned} A(u + v) &= A\left(\sum_{b \in B_{u+v}} c_{u+v}(b) b\right) = \sum_{b \in B_{u+v}} c_{u+v}(b) w_b \\ &= \sum_{b \in B_u} c_u(b) w_b + \sum_{b \in B_v} c_v(b) w_b = A(u) + A(v). \end{aligned} \quad (\text{A.32c})$$

If $v \in V$ and B_v and c_v are as before, then the linearity of \tilde{A} and (A.31) imply

$$\begin{aligned} \tilde{A}(v) &= \tilde{A}\left(\sum_{b \in B_v} c_v(b) b\right) \stackrel{\tilde{A} \in \mathcal{L}(V, W)}{=} \sum_{b \in B_v} c_v(b) \tilde{A}(b) \\ &= \sum_{b \in B_v} c_v(b) A(b) \stackrel{(*)}{=} \sum_{b \in B_v} c_v(b) w_b = A(v), \end{aligned} \quad (\text{A.33})$$

where, at $(*)$, it was used that, for each $b \in B$, one has $A(b) = w_b$ (note $B_b = \{b\}$, $c_b(b) = 1$). Since (A.33) establishes $\tilde{A} = A$, the proof is complete. \blacksquare

Theorem A.24. *Let V and W be vector spaces over the field F . Then $V \cong W$ (i.e. V and W are isomorphic) if, and only if, $\dim V = \dim W$.*

Proof. Suppose $\dim V = \dim W$. If B_V is a basis of V and B_W is a basis of W , then there exists a bijective map $i : B_V \rightarrow B_W$. According to Th. A.23, i defines a unique

linear map $A : V \longrightarrow W$ with $A(b) = i(b)$ for each $b \in B_V$. More precisely, letting, once again, for each $v \in V$, B_v and $c_v : B_v \longrightarrow F \setminus \{0\}$ be as in Th. A.11 (writing c_v instead of c to underline the dependence of c on v),

$$\forall_{v \in V} \quad A(v) = A \left(\sum_{b \in B_v} c_v(b) b \right) = \sum_{b \in B_v} c_v(b) i(b). \quad (\text{A.34})$$

It remains to show A is bijective. If $v \neq 0$, then $B_v \neq \emptyset$ and $A(v) = \sum_{b \in B_v} c_v(b) i(b) \neq 0$, since $c_v(b) \neq 0$ and $\{i(b) : b \in B_v\} \subseteq B_W$ is linearly independent, showing A is injective by Th. A.20(b). If $w \in W$, then there exists a finite set $\tilde{B}_w \subseteq B_W$ and $c_w : \tilde{B}_w \longrightarrow F$ such that $\sum_{\tilde{b} \in \tilde{B}_w} \tilde{c}_w(\tilde{b}) \tilde{b}$. Then

$$\begin{aligned} A \left(\sum_{\tilde{b} \in \tilde{B}_w} \tilde{c}_w(\tilde{b}) i^{-1}(\tilde{b}) \right) &\stackrel{\tilde{A} \in \mathcal{L}(V, W)}{=} \sum_{\tilde{b} \in \tilde{B}_w} \tilde{c}_w(\tilde{b}) A \left(i^{-1}(\tilde{b}) \right) \stackrel{i^{-1}(\tilde{b}) \in B_V}{=} \sum_{\tilde{b} \in \tilde{B}_w} \tilde{c}_w(\tilde{b}) i \left(i^{-1}(\tilde{b}) \right) \\ &= \sum_{\tilde{b} \in \tilde{B}_w} \tilde{c}_w(\tilde{b}) \tilde{b} = w, \end{aligned} \quad (\text{A.35})$$

showing $\text{Im } A = W$, completing the proof that A is bijective.

If $A : V \longrightarrow W$ is a linear isomorphism and B is a basis for V , then, by Prop. A.22(c), $A(B)$ is a basis for W . As A is bijective, so is $A|_B$, showing $\dim V = \#B = \#A(B) = \dim W$ as claimed. \blacksquare

Definition A.25. Let V and W be vector spaces over the field F . We define an addition and a scalar multiplication on $\mathcal{L}(V, W)$ by

$$(A + B) : V \longrightarrow W, \quad (A + B)(x) := A(x) + B(x), \quad (\text{A.36a})$$

$$(\lambda \cdot A) : V \longrightarrow W, \quad (\lambda \cdot A)(x) := \lambda \cdot A(x) \quad \text{for each } \lambda \in F. \quad (\text{A.36b})$$

Theorem A.26. Let V and W be vector spaces over the field F . The addition and scalar multiplication on $\mathcal{L}(V, W)$ given by (A.36) are well-defined in the sense that, if $A, B \in \mathcal{L}(V, W)$ and $\lambda \in F$, then $A + B \in \mathcal{L}(V, W)$ and $\lambda A \in \mathcal{L}(V, W)$. Moreover, with the operations defined in (A.36), $\mathcal{L}(V, W)$ forms a vector space over F .

Proof. See, e.g., [Str08, Th. 13.2]. \blacksquare

Theorem A.27. Let V and W be finite dimensional vector spaces over the field F , let $\{v_1, \dots, v_n\}$ and $\{w_1, \dots, w_m\}$ be bases of V and W , respectively; $m, n \in \mathbb{N}$. Using Th. A.23, define maps $A_{ji} \in \mathcal{L}(V, W)$ by letting

$$\forall_{(j,i,k) \in \{1, \dots, m\} \times \{1, \dots, n\}^2} \quad A_{ji}(v_k) := \begin{cases} w_j & \text{for } k = i, \\ 0 & \text{for } k \neq i. \end{cases} \quad (\text{A.37})$$

Then $\{A_{ji} : (j, i) \in \{1, \dots, m\} \times \{1, \dots, n\}\}$ constitutes a basis for $\mathcal{L}(V, W)$ and, in particular,

$$\dim \mathcal{L}(V, W) = \dim V \cdot \dim W = n \cdot m. \quad (\text{A.38})$$

Proof. See, e.g., [Str08, Th. 13.11]. ■

Theorem A.28. Let V, W, X be vector spaces over the field F .

(a) The composition of linear maps is linear, i.e. if $A \in \mathcal{L}(V, W)$ and $B \in \mathcal{L}(W, X)$, then $B \circ A \in \mathcal{L}(V, X)$.

(b) If $A \in \mathcal{L}(V, W)$ and $B, C \in \mathcal{L}(W, X)$, then

$$A \circ (B + C) = A \circ B + A \circ C. \quad (\text{A.39})$$

(c) If $A, B \in \mathcal{L}(V, W)$ and $C \in \mathcal{L}(W, X)$, then

$$(A + B) \circ C = A \circ C + B \circ C. \quad (\text{A.40})$$

Proof. See, e.g., [Str08, Th. 13.3]. ■

A.3 Matrices

Matrices provide a convenient representation for linear maps A between finite dimensional vector spaces V and W . Recall the basis $\{A_{ji} : (j, i) \in \{1, \dots, m\} \times \{1, \dots, n\}\}$ of $\mathcal{L}(V, W)$ that, in Th. A.27, was shown to arise from bases $\{v_1, \dots, v_n\}$ and $\{w_1, \dots, w_m\}$ of V and W , respectively; $m, n \in \mathbb{N}$. Thus, each $A \in \mathcal{L}(V, W)$ can be written in the form

$$A = \sum_{i=1}^n \sum_{j=1}^m a_{ji} A_{ji}, \quad (\text{A.41})$$

with coordinates $(a_{ji})_{(j,i) \in \{1, \dots, m\} \times \{1, \dots, n\}}$ in F . This motivates the following definition of matrices.

Definition A.29. Let F be a field and $m, n \in \mathbb{N}$. A family $(a_{ji})_{(j,i) \in \{1, \dots, m\} \times \{1, \dots, n\}}$ is called an m -by- n or an $m \times n$ *matrix* over F , where $m \times n$ is called the *size*, *dimension* or *type* of the matrix. The a_{ji} are called the *entries* or *elements* of the matrix. One also writes just (a_{ji}) instead of $(a_{ji})_{(j,i) \in \{1, \dots, m\} \times \{1, \dots, n\}}$ if the size of the matrix is understood. One usually thinks of the $m \times n$ matrix (a_{ji}) as the *rectangular array*

$$(a_{ji}) = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \vdots & \vdots \\ a_{m1} & \dots & a_{mn} \end{pmatrix} \quad (\text{A.42})$$

with m rows and n columns. One therefore also calls $1 \times n$ matrices *row vectors* and $m \times 1$ matrices *column vectors*, and one calls $n \times n$ matrices *quadratic*. The set of all $m \times n$ matrices over F is denoted by $\mathcal{M}(m, n, F)$, and for the set of all quadratic $n \times n$ matrices, one uses the abbreviation $\mathcal{M}(n, F) := \mathcal{M}(n, n, F)$.

Definition A.30 (Matrix Arithmetic). Let F be a field and $m, n, l \in \mathbb{N}$.

(a) *Matrix Addition:* For $m \times n$ matrices (a_{ji}) and (b_{ji}) over F , define the sum

$$(a_{ji}) + (b_{ji}) := (a_{ji} + b_{ji}). \quad (\text{A.43})$$

(b) *Scalar Multiplication:* For each $m \times n$ matrix (a_{ji}) and each $\lambda \in F$, define

$$\lambda(a_{ji}) := (\lambda a_{ji}). \quad (\text{A.44})$$

(c) *Matrix Multiplication:* For each $m \times n$ matrix (a_{ji}) and each $n \times l$ matrix (b_{ji}) over F , define the product

$$(a_{ji})(b_{ji}) := \left(\sum_{k=1}^n a_{jk} b_{ki} \right)_{(j,i) \in \{1, \dots, m\} \times \{1, \dots, l\}}, \quad (\text{A.45})$$

i.e. the product of an $m \times n$ matrix and an $n \times l$ matrix is an $m \times l$ matrix (cf. Th. A.34 below).

Remark A.31. We consider matrices over a field F .

(a) For each $m, n \in \mathbb{N}$, the set $\mathcal{M}(m, n, F)$ of $m \times n$ matrices over F with the operations of Def. A.30(a),(b) constitutes a vector space over F : An $m \times n$ matrix $A = (a_{ji})_{(j,i) \in \{1, \dots, m\} \times \{1, \dots, n\}}$ is defined as a family in F , i.e., recalling [Phi15a, Def. 2.14(a)], A is defined as the function $A : \{1, \dots, m\} \times \{1, \dots, n\} \rightarrow F$, $A(j, i) = a_{ji}$; and $\mathcal{M}(m, n, F) = \mathcal{F}(\{1, \dots, m\} \times \{1, \dots, n\}, F)$ (so we notice that matrices are nothing new in terms of objects, but just a new way of thinking about functions from $\{1, \dots, m\} \times \{1, \dots, n\}$ into F , that turns out to be convenient in certain contexts). Thus, the operations defined in Def. A.30(a),(b) are precisely the same operations that were defined in (A.5) and $\mathcal{M}(m, n, F)$ is a vector space according to Ex. A.2(c). Clearly, the map

$$I : \mathcal{M}(m, n, F) \rightarrow F^{m \cdot n}, \quad (a_{ji}) \mapsto (\lambda_1, \dots, \lambda_{m \cdot n}), \quad (\text{A.46})$$

where $\lambda_k = a_{ji}$ if, and only if, $k = (j - 1) \cdot n + i$,

constitutes a linear isomorphism. Other important linear isomorphisms between $\mathcal{M}(m, n, F)$ and vector spaces of linear maps will be provided in Th. A.32 below.

(b) Matrix multiplication is associative whenever all relevant multiplications are defined. More precisely, if A is an $m \times n$ matrix, B is an $n \times l$, and C is an $l \times p$ matrix, then

$$(AB)C = A(BC) : \quad (\text{A.47})$$

Indeed, one has $m \times p$ matrices $(AB)C = (d_{ji})$ and $A(BC) = (e_{ji})$, where

$$d_{ji} = \sum_{\alpha=1}^l \left(\sum_{k=1}^n a_{jk} b_{k\alpha} \right) c_{\alpha i} = \sum_{\alpha=1}^l \sum_{k=1}^n a_{jk} b_{k\alpha} c_{\alpha i} = \sum_{k=1}^n a_{jk} \left(\sum_{\alpha=1}^l b_{k\alpha} c_{\alpha i} \right) = e_{ji}. \quad (\text{A.48})$$

- (c) Matrix multiplication is, in general, *not* commutative: If A is an $m \times n$ matrix and B is an $n \times l$ with $m \neq l$, then BA is not even defined. If $m = l$, but $m \neq n$, then AB has dimension $m \times m$, but BA has different dimension, namely $n \times n$. And even if $m = n = l > 1$, then commutativity is, in general not true – for example

$$\begin{pmatrix} 1 & 1 & \dots & 1 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 \end{pmatrix} = \begin{pmatrix} \lambda & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix}, \quad (\text{A.49a})$$

$$\begin{pmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 \end{pmatrix} \begin{pmatrix} 1 & 1 & \dots & 1 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix} = \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{pmatrix}. \quad (\text{A.49b})$$

Note that $\lambda = m$ for $F = \mathbb{R}$, but, in general, λ will depend on F , e.g. for $F = \{0, 1\}$, one obtains $\lambda = m \pmod{2}$.

—

Let us come back to the situation discussed at the beginning of the section above, resulting in (A.41). Let $v = \sum_{i=1}^n \lambda_i v_i \in V$ with $\lambda_1, \dots, \lambda_n \in F$. Then

$$\begin{aligned} A(v) &= \sum_{i=1}^n \lambda_i A(v_i) = \sum_{i=1}^n \lambda_i \sum_{k=1}^n \sum_{j=1}^m a_{jk} A_{jk}(v_i) \stackrel{(\text{A.37})}{=} \sum_{i=1}^n \lambda_i \sum_{j=1}^m a_{ji} w_j \\ &= \sum_{j=1}^m \left(\sum_{i=1}^n a_{ji} \lambda_i \right) w_j. \end{aligned} \quad (\text{A.50})$$

Thus, if we represent v by a column vector \tilde{v} (an $n \times 1$ matrix) containing its coordinates $\lambda_1, \dots, \lambda_n$ with respect to the basis $\{v_1, \dots, v_n\}$ and $A(v)$ by a column vector \tilde{w} (an $m \times 1$ matrix) containing its coordinates with respect to the basis $\{w_1, \dots, w_m\}$, then (A.50) shows

$$\tilde{w} = M\tilde{v} = M \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_n \end{pmatrix}, \quad \text{where } M := (a_{ji}). \quad (\text{A.51})$$

For finite dimensional vector spaces, the precise relationship between linear maps, bases, and matrices is provided by the following theorem:

Theorem A.32. *Let V and W be finite dimensional vector spaces over the field F , let $\{v_1, \dots, v_n\}$ and $\{w_1, \dots, w_m\}$ be bases of V and W , respectively; $m, n \in \mathbb{N}$. Then the map*

$$I : \mathcal{L}(V, W) \longrightarrow \mathcal{M}(m, n, F), \quad A \mapsto (a_{ji}), \quad (\text{A.52})$$

where the a_{ji} are given by (A.41) constitutes a linear isomorphism.

Proof. According to Th. A.27, $\{A_{ji} : (j, i) \in \{1, \dots, m\} \times \{1, \dots, n\}\}$ forms a basis of $\mathcal{L}(V, W)$. Thus, to every family of coordinates $\{a_{ji} : (j, i) \in \{1, \dots, m\} \times \{1, \dots, n\}\}$ in F , (A.41) defines a unique element of $\mathcal{L}(V, W)$, i.e. I is bijective. It remains to verify that I is linear. To this end, let $\lambda, \mu \in F$ and $A, B \in \mathcal{M}(m, n, F)$ with

$$A = \sum_{i=1}^n \sum_{j=1}^m a_{ji} A_{ji}, \quad (a_{ji}) = I(A) \in \mathcal{M}(m, n, F), \quad (\text{A.53a})$$

$$B = \sum_{i=1}^n \sum_{j=1}^m b_{ji} A_{ji}, \quad (b_{ji}) = I(B) \in \mathcal{M}(m, n, F). \quad (\text{A.53b})$$

Then

$$\lambda A + \mu B = \lambda \sum_{i=1}^n \sum_{j=1}^m a_{ji} A_{ji} + \mu \sum_{i=1}^n \sum_{j=1}^m b_{ji} A_{ji} = \sum_{i=1}^n \sum_{j=1}^m (\lambda a_{ji} + \mu b_{ji}) A_{ji}, \quad (\text{A.54})$$

showing

$$I(\lambda A + \mu B) = (\lambda a_{ji} + \mu b_{ji}) A_{ji} = \lambda (a_{ji}) + \mu (b_{ji}) = \lambda I(A) + \mu I(B), \quad (\text{A.55})$$

proving the linearity of I . ■

Definition and Remark A.33. In the situation of Th. A.32, for each $A \in \mathcal{L}(V, W)$, one calls the matrix $I(A) = (a_{ji}) \in \mathcal{M}(m, n, F)$ the *(transformation) matrix corresponding to A with respect to* the basis $\{v_1, \dots, v_n\}$ of V and the basis $\{w_1, \dots, w_m\}$ of W . If the bases are understood, then one often tends to identify the map with its corresponding matrix.

However, as $I(A)$ depends on the bases, identifying A and $I(A)$ is only admissible as long as one keeps the bases of V and W fixed! Moreover, if one represents matrices as rectangular arrays as in (A.42) (which one usually does), then one actually considers the basis vectors of $\{v_1, \dots, v_n\}$ and $\{w_1, \dots, w_m\}$ as *ordered* from 1 to n (resp. m), i.e. $I(A)$ actually depends on the so-called *ordered bases* (v_1, \dots, v_n) and (w_1, \dots, w_m) (ordered bases are tuples rather than sets and the matrix corresponding to A changes if the order of the basis vectors changes).

Similarly, we had seen in (A.51) that it can be useful to identify a vector $v = \sum_{i=1}^n \lambda_i v_i$ with its coordinates $(\lambda_1, \dots, \lambda_n)$, typically represented as an $n \times 1$ matrix (a column vector, as in (A.51)) or a $1 \times n$ matrix (a row vector). Obviously, this identification is also only admissible as long as the basis $\{v_1, \dots, v_n\}$ and its order is kept fixed.

—

The following Th. A.34 is the justification for defining matrix multiplication according to Def. A.30(c).

Theorem A.34. *Let F be a field, let $n, m, l \in \mathbb{N}$, and let V, W, X be finite dimensional vector spaces over F such that V has basis $\{v_1, \dots, v_n\}$, W has basis $\{w_1, \dots, w_m\}$, and*

X has basis $\{x_1, \dots, x_l\}$. If $A \in \mathcal{L}(V, W)$, $B \in \mathcal{L}(W, X)$, $M = (a_{ji}) \in \mathcal{M}(m, n, F)$ is the matrix corresponding to A with respect to $\{v_1, \dots, v_n\}$ and $\{w_1, \dots, w_m\}$, and $N = (b_{ji}) \in \mathcal{M}(l, m, F)$ is the matrix corresponding to B with respect to $\{w_1, \dots, w_m\}$ and $\{x_1, \dots, x_l\}$, then $NM = (\sum_{k=1}^m b_{jk}a_{ki}) \in \mathcal{M}(l, n, F)$ is the matrix corresponding to BA with respect to $\{v_1, \dots, v_n\}$ and $\{x_1, \dots, x_l\}$.

Proof. For each $i \in \{1, \dots, n\}$, one computes

$$\begin{aligned} (BA)(v_i) &= B(A(v_i)) = B\left(\sum_{k=1}^m a_{ki}w_k\right) = \sum_{k=1}^m a_{ki}B(w_k) = \sum_{k=1}^m a_{ki} \sum_{j=1}^l b_{jk}x_j \\ &= \sum_{j=1}^l \sum_{k=1}^m b_{jk}a_{ki}x_j = \sum_{j=1}^l \left(\sum_{k=1}^m b_{jk}a_{ki}\right) x_j, \end{aligned} \quad (\text{A.56})$$

proving $NM = (\sum_{k=1}^m b_{jk}a_{ki})$ is the matrix corresponding to BA with respect to the bases $\{v_1, \dots, v_n\}$ and $\{x_1, \dots, x_l\}$. ■

Definition and Remark A.35. Let F be a field, $A := (a_{ji})_{(j,i) \in \{1, \dots, m\} \times \{1, \dots, n\}} \in \mathcal{M}(m, n, F)$, and $m, n \in \mathbb{N}$. Then we define the *transpose* of A , denoted A^t , by

$$A^t := (a_{ji})_{(i,j) \in \{1, \dots, n\} \times \{1, \dots, m\}}. \quad (\text{A.57})$$

Thus, if A is an $m \times n$ matrix, then its transpose is an $n \times m$ matrix, where one obtains A^t from A by switching rows and columns. One has to use care when using the notation of (A.57), as one often implicitly assumes that, when writing (a_{ji}) , the first index is for rows and the second index for columns. However, this is actually determined by the order of the factors of the cartesian product that determines the domain of the family. Whereas A is the map $f : \{1, \dots, m\} \times \{1, \dots, n\} \rightarrow F$, $f(j, i) = a_{ji}$, its transpose A^t is the map $f^t : \{1, \dots, n\} \times \{1, \dots, m\} \rightarrow F$, $f^t(i, j) = f(j, i) = a_{ji}$. To emphasize this in the notation, one can rewrite (A.57) in the form

$$A^t = (b_{ij})_{(i,j) \in \{1, \dots, n\} \times \{1, \dots, m\}}, \quad \text{where} \quad \forall_{(i,j) \in \{1, \dots, n\} \times \{1, \dots, m\}} \quad b_{ij} := a_{ji}. \quad (\text{A.58})$$

For the transpose of A , one also finds the notation A' instead of A^t .

Theorem A.36. Let F be a field and $m, n, l \in \mathbb{N}$.

(a) The map

$$I : \mathcal{M}(m, n, F) \rightarrow \mathcal{M}(n, m, F), \quad A \mapsto A^t, \quad (\text{A.59})$$

is a linear isomorphism and

$$\forall_{A \in \mathcal{M}(m, n, F)} \quad (A^t)^t = A. \quad (\text{A.60})$$

(b) If $A \in \mathcal{M}(m, n, F)$ and $B \in \mathcal{M}(n, l, F)$, then

$$(AB)^t = B^t A^t. \quad (\text{A.61})$$

Proof. (a): It is immediate from (A.57) that (A.60) is valid, showing I has an inverse map and is, hence, bijective. So it just remains to verify I is linear. However, if $A, B \in \mathcal{M}(m, n, F)$, $A = (a_{ji})$, $B = (b_{ji})$, and $\mu, \lambda \in F$, then

$$\begin{aligned} (\lambda A + \mu B)^t &= (\lambda a_{ji} + \mu b_{ji})_{(j,i) \in \{1, \dots, m\} \times \{1, \dots, n\}}^t = (\lambda a_{ji} + \mu b_{ji})_{(i,j) \in \{1, \dots, n\} \times \{1, \dots, m\}} \\ &= \lambda (a_{ji})_{(i,j) \in \{1, \dots, n\} \times \{1, \dots, m\}} + \mu (b_{ji})_{(i,j) \in \{1, \dots, n\} \times \{1, \dots, m\}} \\ &= \lambda (a_{ji})_{(j,i) \in \{1, \dots, m\} \times \{1, \dots, n\}}^t + \mu (b_{ji})_{(j,i) \in \{1, \dots, m\} \times \{1, \dots, n\}}^t \\ &= \lambda A^t + \mu B^t, \end{aligned} \tag{A.62}$$

thereby establishing the case.

(b): Let $A = (a_{ji})$, $B = (b_{ji})$, $A^t = (a_{ji}^t)$, $B^t = (b_{ji}^t)$. Then

$$\begin{aligned} (AB)^t &\stackrel{(A.45)}{=} \left(\sum_{k=1}^n a_{jk} b_{ki} \right)_{(j,i) \in \{1, \dots, m\} \times \{1, \dots, l\}}^t = \left(\sum_{k=1}^n b_{ki} a_{jk} \right)_{(i,j) \in \{1, \dots, l\} \times \{1, \dots, m\}} \\ &= \left(\sum_{k=1}^n b_{kj} a_{ik} \right)_{(j,i) \in \{1, \dots, l\} \times \{1, \dots, m\}} = \left(\sum_{k=1}^n b_{jk}^t a_{ki}^t \right)_{(j,i) \in \{1, \dots, l\} \times \{1, \dots, m\}} \\ &\stackrel{(A.45)}{=} B^t A^t, \end{aligned} \tag{A.63}$$

proving (A.61). ■

A.4 Determinants

For each quadratic matrix $A \in \mathcal{M}(n, F)$, one can define its determinant $\det(A) \in F$, resulting in a function $\det : \mathcal{M}(n, F) \rightarrow F$ that is often useful when studying matrices and linear maps. One can characterize the determinant function axiomatically (see Def. A.38 below), and, with some preparation, one can also provide an explicit formula (see (A.82) below).

One important feature of the determinant is its being nonzero if, and only if, the matrix is invertible (cf. Def. A.46 and Th. A.48(a) below). Another is the fact that the determinant's value only depends on the linear map defined by A and an arbitrary basis of F^n . This allows to define $\det : \mathcal{L}(V, V) \rightarrow F$ as in Def. and Rem. A.53 below. One can show that, if $A \in \mathcal{L}(V, V) \rightarrow F$, then $\det(A)$ is a measure of the n -dimensional volume distortion caused by applying A , but, here, we will not pursue this aspect.

Definition A.37. Let F be a field, $n \in \mathbb{N}$. Then the $n \times n$ matrix

$$\text{Id} := \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix} = (e_{ji}), \quad \text{where} \quad e_{ji} := \begin{cases} 1 & \text{for } j = i, \\ 0 & \text{for } j \neq i, \end{cases} \tag{A.64}$$

is called *identity matrix* or just *identity* or *unit matrix*. The dependence of Id on n is suppressed in the notation, but n should always be clear from the context. In the literature, one also finds the notation E or I instead of Id .

Definition A.38. Let F be a field, $n \in \mathbb{N}$. A map $\det : \mathcal{M}(n, F) \rightarrow F$ is called *determinant* if, and only if, it satisfies the following conditions (i) – (iii):

- (i) \det is *multilinear* with regard to matrix columns, i.e., for each $A \in \mathcal{M}(n, F)$, $b \in \mathcal{M}(n, 1, F)$, $i \in \{1, \dots, n\}$, and $\lambda, \mu \in F$:

$$\begin{aligned} & \det(a_1, \dots, \lambda a_i + \mu b, \dots, a_n) \\ &= \lambda \det(A) + \mu \det(a_1, \dots, a_{i-1}, b, a_{i+1}, \dots, a_n), \end{aligned} \quad (\text{A.65})$$

where a_1, \dots, a_n denote the columns of A .

- (ii) If the columns a_1, \dots, a_n of $A = (a_1, \dots, a_n) \in \mathcal{M}(n, F)$ are linearly dependent, then $\det(A) = 0$.

- (iii) $\det(\text{Id}) = 1$.

Notation A.39. If F is a field, $n \in \mathbb{N}$, and $\det : \mathcal{M}(n, F) \rightarrow F$ is a determinant, then, for $A = (a_{ji}) \in \mathcal{M}(n, F)$, one commonly uses the notation

$$|A| := \begin{vmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \vdots & \vdots \\ a_{m1} & \dots & a_{mn} \end{vmatrix} := \det(A). \quad (\text{A.66})$$

In Th. A.45 below, it will be stated that, for each $n \in \mathbb{N}$, there exists a unique determinant. To also state an explicit formula for this determinant, we need to know a few things about permutations.

Definition and Remark A.40. Let $n \in \mathbb{N}$. Each bijective map $\pi : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$ is called a *permutation* of $\{1, \dots, n\}$. The set of permutations of $\{1, \dots, n\}$ forms a group with respect to the composition of maps, the so-called *symmetric group* S_n : Indeed, the composition of maps is associative by [Phi15a, Prop. 2.9(a)]; the neutral element is the identity map $e : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$, $e(i) = i$; and, for each $\sigma \in S_n$, its inverse map σ^{-1} is also its inverse element in the group S_n . Caveat: Simple examples show that S_n is *not* commutative.

Definition A.41. Let $k, n \in \mathbb{N}$, $k \leq n$. A permutation $\pi \in S_n$ is called a *k-cycle* if, and only if, there exist k distinct numbers $i_1, \dots, i_k \in \{1, \dots, n\}$ such that

$$\pi(i) = \begin{cases} i_{j+1} & \text{if } j \in \{1, \dots, k-1\}, \\ i_1 & \text{if } i = i_k, \\ i & \text{if } i \notin \{i_1, \dots, i_k\}. \end{cases} \quad (\text{A.67})$$

If π is a cycle as in (A.67), then one writes

$$\pi = (i_1 \ i_2 \ \dots \ i_k). \quad (\text{A.68})$$

Each 2-cycle is also known as a *transposition*.

Theorem A.42. *Let $n \in \mathbb{N}$.*

- (a) *Each permutation can be decomposed into finitely many disjoint cycles: For each $\pi \in S_n$ there exists a decomposition of $\{1, \dots, n\}$ into disjoint sets A_1, \dots, A_N , $N \in \mathbb{N}$, i.e.*

$$\{1, \dots, n\} = \bigcup_{i=1}^N A_i \quad \text{and} \quad A_i \cap A_j = \emptyset \quad \text{for } i \neq j, \quad (\text{A.69})$$

such that A_i consists of the distinct elements a_{i1}, \dots, a_{i,N_i} and

$$\pi = (a_{N1} \dots a_{N,N_N}) \cdots (a_{11} \dots a_{1,N_1}). \quad (\text{A.70})$$

The decomposition (A.70) is unique up to the order of the cycles.

- (b) *If $n \geq 2$, then every permutation $\pi \in S_n$ is the composition of finitely many transpositions, where each transposition permutes two juxtaposed elements, i.e.*

$$\forall_{\pi \in S_n} \quad \exists_{N \in \mathbb{N}} \quad \exists_{\tau_1, \dots, \tau_N \in T} \quad \pi = \tau_N \circ \cdots \circ \tau_1, \quad (\text{A.71})$$

where $T := \{(i \ i+1) : i \in \{1, \dots, n-1\}\}$.

Proof. (a): We prove the statement by induction on n . For $n = 1$, there is nothing to prove. Let $n > 1$ and choose $i \in \{1, \dots, n\}$. We claim that

$$\exists_{k \in \mathbb{N}} \quad \left(\pi^k(i) = i \wedge \forall_{l \in \{1, \dots, k-1\}} \pi^l(i) \neq i \right). \quad (\text{A.72})$$

Indeed, since $\{1, \dots, n\}$ is finite, there must be a smallest $k \in \mathbb{N}$ such that $\pi^k(i) \in A_1 := \{i, \pi(i), \dots, \pi^{k-1}(i)\}$. Since π is bijective, it must be $\pi^k(i) = i$ and $(i \ \pi(i) \ \dots, \pi^{k-1}(i))$ is a k -cycle. We are already done in case $k = n$. If $k < n$, then consider $B := \{1, \dots, n\} \setminus A_1$. Then, again using the bijectivity of π , $\pi|_B$ is a permutation on B with $1 \leq \#B < n$. By induction, there are disjoint sets A_2, \dots, A_N such that $B = \bigcup_{j=2}^N A_j$, A_j consists of the distinct elements a_{j1}, \dots, a_{j,N_j} and

$$\pi|_B = (a_{N1} \dots a_{N,N_N}) \cdots (a_{21} \dots a_{2,N_2}).$$

Since $\pi = (i \ \pi(i) \ \dots, \pi^{k-1}(i)) \circ \pi|_B$, this finishes the proof of (A.70). If there were another, different, decomposition of π into cycles, say, given by disjoint sets B_1, \dots, B_M , $\{1, \dots, n\} = \bigcup_{i=1}^M B_i$, $M \in \mathbb{N}$, then there were $A_i \neq B_j$ and $k \in A_i \cap B_j$. But then k were in the cycle given by A_i and in the cycle given by B_j , implying $A_i = \{\pi^l(k) : l \in \mathbb{N}\} = B_j$, in contradiction to $A_i \neq B_j$.

(b): We first show that every $\pi \in S_n$ is a composition of finitely many transpositions (not necessarily transpositions from the set T): According to (a), it suffices to show that every cycle is a composition of finitely many transpositions. Since each 1-cycle is

the identity, it is $(i) = \text{Id} = (1\ 2)(1\ 2)$ for each $i \in \{1, \dots, n\}$. If $(i_1 \dots i_k)$ is a k -cycle, $k \in \{2, \dots, n\}$, then

$$(i_1 \dots i_k) = (i_1\ i_2)(i_2\ i_3) \cdots (i_{k-1}\ i_k) : \quad (\text{A.73})$$

Indeed,

$$\forall_{i \in \{1, \dots, n\}} (i_1\ i_2)(i_2\ i_3) \cdots (i_{k-1}\ i_k)(i) = \begin{cases} i_1 & \text{for } i = i_k, \\ i_{l+1} & \text{for } i = i_l, l \in \{1, \dots, k-1\}, \\ i & \text{for } i \notin \{i_1, \dots, i_k\}, \end{cases} \quad (\text{A.74})$$

proving (A.73). To finish the proof of (b), we observe that every transposition is a composition of finitely many elements of T : If $i, j \in \{1, \dots, n\}$, $i < j$, then

$$(i\ j) = (i\ i+1) \cdots (j-2\ j-1)(j-1\ j) \cdots (i+1\ i+2)(i\ i+1) : \quad (\text{A.75})$$

Indeed,

$$\begin{aligned} \forall_{k \in \{1, \dots, n\}} (i\ i+1) \cdots (j-2\ j-1)(j-1\ j) \cdots (i+1\ i+2)(i\ i+1)(k) \\ = \begin{cases} j & \text{for } k = i, \\ i & \text{for } k = j, \\ k & \text{for } i < k < j, \\ k & \text{for } k \notin \{i, i+1, \dots, j\}, \end{cases} \end{aligned} \quad (\text{A.76})$$

proving (A.75). ■

Definition A.43. Let $n \in \mathbb{N}$. For each permutation $\pi \in S_n$, one defines its *sign*, $\text{sgn}(\pi)$, via the map

$$\text{sgn} : S_n \longrightarrow \{-1, 1\}, \quad \text{sgn}(\pi) := \prod_{1 \leq i < j \leq n} \frac{\pi(i) - \pi(j)}{i - j}. \quad (\text{A.77})$$

Note that, for $n = 1$, $\text{sgn} : S_1 = \{e\} \longrightarrow \{-1, 1\}$, $\text{sgn}(e) = 1$, as the product in (A.77) is empty.

Proposition A.44. Let $n \in \mathbb{N}$.

- (a) The sign is well-defined by (A.77), i.e. the map is, indeed, $\{-1, 1\}$ -valued.
- (b) The function $\text{sgn} : S_n \longrightarrow \{-1, 1\}$ is a group homomorphism (note that $\{-1, 1\}$ forms a multiplicative subgroup of \mathbb{R}), i.e.

$$\forall_{\pi_1, \pi_2 \in S_n} \text{sgn}(\pi_1 \circ \pi_2) = \text{sgn}(\pi_1) \text{sgn}(\pi_2). \quad (\text{A.78})$$

- (c) For $n \geq 2$, if a permutation $\pi \in S_n$ is the composition of k transpositions, then the parity of k is uniquely determined by π (i.e., for a given π , k is either always even or always odd) and

$$\operatorname{sgn}(\pi) = (-1)^k = \begin{cases} 1 & \text{if } k \text{ is even,} \\ -1 & \text{if } k \text{ is odd.} \end{cases} \quad (\text{A.79})$$

Proof. (a): The map sgn is $\{-1, 1\}$ -valued, since the bijectivity of $\pi \in S_n$ implies that the factor $i - j$ appears in the denominator of $\operatorname{sgn}(\pi)$ as defined in (A.77) if, and only if, the factor $i - j$ or the factor $j - i$ appears in the numerator.

(b): Let $\pi_1, \pi_2 \in S_n$. One computes

$$\begin{aligned} \operatorname{sgn}(\pi_1 \circ \pi_2) &= \prod_{1 \leq i < j \leq n} \frac{\pi_1(\pi_2(i)) - \pi_1(\pi_2(j))}{i - j} \\ &= \prod_{1 \leq i < j \leq n} \left(\frac{\pi_1(\pi_2(i)) - \pi_1(\pi_2(j))}{\pi_2(i) - \pi_2(j)} \cdot \frac{\pi_2(i) - \pi_2(j)}{i - j} \right) \\ &\stackrel{\pi_2 \text{ bij.}}{=} \operatorname{sgn}(\pi_1) \operatorname{sgn}(\pi_2). \end{aligned} \quad (\text{A.80})$$

(c): If $\tau \in S_n$ is a transposition, then there are elements $i, j \in \{1, \dots, n\}$ such that $i < j$ and $\tau = (i \ j)$. Thus,

$$\operatorname{sgn}(\tau) = \frac{\tau(i) - \tau(j)}{i - j} = \frac{j - i}{i - j} = -1 \quad (\text{A.81})$$

holds for every transposition τ . In consequence, if $\pi \in S_n$ is the composition of k transpositions, $k \in \mathbb{N}$, then (A.79) must hold and, in particular, k is always even if $\operatorname{sgn}(\pi) = 1$ and k is always odd if $\operatorname{sgn}(\pi) = -1$. ■

Theorem A.45. *Let F be a field. For each $n \in \mathbb{N}$, there exists a unique determinant, i.e. there is a unique map $\det : \mathcal{M}(n, F) \rightarrow F$, satisfying (i) – (iii) of Def. A.38. Moreover, this map is given by*

$$\det : \mathcal{M}(n, F) \rightarrow F, \quad \det((a_{ji})) := \sum_{\pi \in S_n} \operatorname{sgn}(\pi) a_{1\pi(1)} \cdots a_{n\pi(n)}. \quad (\text{A.82})$$

Proof. See, e.g., [Str08, Th. 17.5, Th. 17.11(a)]. ■

Definition A.46. Let F be a field, $n \in \mathbb{N}$. A quadratic matrix $A \in \mathcal{M}(n, F)$ is called *invertible* or *regular* if, and only if,

$$\exists_{B \in \mathcal{M}(n, F)} AB = \operatorname{Id}. \quad (\text{A.83})$$

One then usually writes A^{-1} instead of B and calls A^{-1} the *inverse matrix* of A . If A is not regular, then it is called *singular*.

Remark A.47. If V is a finite dimensional vector space over a field F , and $\{v_1, \dots, v_n\}$ is a basis of V , $n \in \mathbb{N}$, then, due to Th. A.32, a linear map $A \in \mathcal{L}(V, V)$ is bijective if, and only if, its transformation matrix $I(V)$ with respect to the given basis is invertible.

—

Important properties of the determinant are compiled in the following Th. A.48.

Theorem A.48. Let F be a field, $n \in \mathbb{N}$, let $A \in \mathcal{M}(n, F)$, and let c_1, \dots, c_n denote the columns of A , whereas r_1, \dots, r_n denote the rows of A , i.e.

$$A = (c_1, \dots, c_n) = \begin{pmatrix} r_1 \\ \vdots \\ r_n \end{pmatrix}. \quad (\text{A.84})$$

- (a) $\det(A) = 0$ if, and only if, A is singular. If A is invertible, then $\det(A^{-1}) = (\det(A))^{-1}$.
- (b) If $B \in \mathcal{M}(n, F)$, then $\det(AB) = \det(A) \det(B)$.
- (c) $\det(A^t) = \det(A)$.
- (d) If $\lambda \in F$, then $\det(\lambda A) = \lambda^n \det(A)$.
- (e) The value of the determinant remains the same if one column of a matrix is replaced by the sum of that column and a scalar multiple of another column. More generally, the determinant remains the same if one column of a matrix is replaced by the sum of that column and a linear combination of the other columns, i.e., if $\lambda_1, \dots, \lambda_n \in F$ and $i \in \{1, \dots, n\}$, then

$$\det(A) = \det(c_1, \dots, c_n) = \det \left(c_1, \dots, c_{i-1}, c_i + \sum_{\substack{j=1 \\ j \neq i}}^n \lambda_j c_j, c_{i+1}, \dots, c_n \right). \quad (\text{A.85})$$

- (f) Switching columns i and j , where $i, j \in \{1, \dots, n\}$, $i \neq j$, changes the sign of the determinant, i.e.

$$\det(c_1, \dots, c_i, \dots, c_j, \dots, c_n) = -\det(c_1, \dots, c_j, \dots, c_i, \dots, c_n). \quad (\text{A.86})$$

- (g) \det is multilinear with regard to matrix rows, i.e., for each $b \in \mathcal{M}(1, n, F)$, $i \in \{1, \dots, n\}$, and $\lambda, \mu \in F$:

$$\det \begin{pmatrix} r_1 \\ \vdots \\ r_{i-1} \\ \lambda r_i + \mu b \\ r_{i+1} \\ \vdots \\ r_n \end{pmatrix} = \lambda \det(A) + \mu \det \begin{pmatrix} r_1 \\ \vdots \\ r_{i-1} \\ b \\ r_{i+1} \\ \vdots \\ r_n \end{pmatrix}. \quad (\text{A.87})$$

- (h) *The value of the determinant remains the same if one row of a matrix is replaced by the sum of that row and a scalar multiple of another row. More generally, the determinant remains the same if one row of a matrix is replaced by the sum of that row and a linear combination of the other rows, i.e., if $\lambda_1, \dots, \lambda_n \in F$ and $i \in \{1, \dots, n\}$, then*

$$\det(A) = \det \begin{pmatrix} r_1 \\ \vdots \\ r_i \\ \vdots \\ r_n \end{pmatrix} = \det \begin{pmatrix} r_1 \\ \vdots \\ r_i + \sum_{\substack{j=1 \\ j \neq i}}^n \lambda_j r_j \\ \vdots \\ r_n \end{pmatrix}. \quad (\text{A.88})$$

- (i) *Switching rows i and j , where $i, j \in \{1, \dots, n\}$, $i \neq j$, changes the sign of the determinant, i.e.*

$$\det \begin{pmatrix} r_1 \\ \vdots \\ r_i \\ \vdots \\ r_j \\ \vdots \\ r_n \end{pmatrix} = - \det \begin{pmatrix} r_1 \\ \vdots \\ r_j \\ \vdots \\ r_i \\ \vdots \\ r_n \end{pmatrix}. \quad (\text{A.89})$$

Proof. (a): See, e.g., [Str08, Th. 17.7(b), Th. 17.11(a)].

(b): See, e.g., [Str08, Th. 17.11(b)].

(c): See, e.g., [Str08, Lem. 18.1].

(d) is an immediate consequence of Def. A.38(i).

(e): One computes, for $i < j$,

$$\begin{aligned} & \det(c_1, \dots, c_{i-1}, c_i + \lambda_j c_j, c_{i+1}, \dots, c_n) \\ & \stackrel{\text{Def. A.38(i)}}{=} \lambda_j^{-1} \det(c_1, \dots, c_{i-1}, c_i + \lambda_j c_j, c_{i+1}, \dots, \lambda_j c_j, \dots, c_n) \\ & \stackrel{\text{Def. A.38(i)}}{=} \lambda_j^{-1} \det(c_1, \dots, c_{i-1}, c_i, c_{i+1}, \dots, \lambda_j c_j, \dots, c_n) \\ & \quad + \lambda_j^{-1} \det(c_1, \dots, c_{i-1}, \lambda_j c_j, c_{i+1}, \dots, \lambda_j c_j, \dots, c_n) \\ & \stackrel{\text{Def. A.38(ii)}}{=} \lambda_j^{-1} \det(c_1, \dots, c_{i-1}, c_i, c_{i+1}, \dots, \lambda_j c_j, \dots, c_n) + 0 \\ & \stackrel{\text{Def. A.38(i)}}{=} \det(c_1, \dots, c_n) = \det(A). \end{aligned} \quad (\text{A.90})$$

The general case of (A.85) then follows by induction.

(f): We compute

$$\begin{aligned}
 & \det(c_1, \dots, c_i, \dots, c_j, \dots, c_n) + \det(c_1, \dots, c_j, \dots, c_i, \dots, c_n) \\
 & \stackrel{(e)}{=} \det(c_1, \dots, c_i + c_j, \dots, c_j, \dots, c_n) + \det(c_1, \dots, c_j + c_i, \dots, c_i, \dots, c_n) \\
 & \stackrel{\text{Def. A.38(i)}}{=} \det(c_1, \dots, c_i + c_j, \dots, c_i + c_j, \dots, c_n) \stackrel{\text{Def. A.38(ii)}}{=} 0,
 \end{aligned} \tag{A.91}$$

proving (f).

(g) is inferred by combining Def. A.38(i) with (c).

(h) is inferred by combining (e) with (c).

(i) is inferred by combining (f) with (c). ■

Theorem A.49 (Block Matrices). *The determinant of so-called block matrices, where one block is a zero matrix (all entries 0), can be computed as the product of the determinants of the corresponding blocks. More precisely, if $n, m \in \mathbb{N}$, then*

$$\begin{vmatrix}
 a_{11} & \dots & a_{1n} & & & \\
 \vdots & \vdots & \vdots & & * & \\
 a_{n1} & \dots & a_{nn} & & & \\
 0 & \dots & 0 & b_{11} & \dots & b_{1m} \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
 0 & \dots & 0 & b_{m1} & \dots & b_{mm}
 \end{vmatrix} = \det(a_{ji}) \det(b_{ji}). \tag{A.92}$$

Proof. See, e.g., [Str08, Th. 18.3]. ■

Definition A.50. Let F be a field, $n \in \mathbb{N}$, $n \geq 2$, $A = (a_{ji}) \in \mathcal{M}(n, F)$. For each $j, i \in \{1, \dots, n\}$, let M_{ji} the $n \times n$ submatrix of A obtained by deleting the j th row and the i th column of A – the M_{ji} are sometimes called the *minor matrices* of A ; define

$$A_{ji} := (-1)^{i+j} \det(M_{ij}), \tag{A.93}$$

where the A_{ji} are called *cofactors* of A and the $\det(M_{ij})$ are called the *minors* of A . Let $\tilde{A} := (A_{ji})$ denote the matrix of cofactors.

Theorem A.51. *Let F be a field, $n \in \mathbb{N}$, $n \geq 2$, $A = (a_{ji}) \in \mathcal{M}(n, F)$. Moreover, let $\tilde{A} := (A_{ji})$ be the matrix of cofactors according to Def. A.50.*

(a) $A\tilde{A} = (\det A) \text{Id}$.

(b) If $\det A \neq 0$, then $\det \tilde{A} = (\det A)^{n-1}$.

(c) If $\det A \neq 0$, then $A^{-1} = (\det A)^{-1} \tilde{A}$.

(d) Laplace Expansion by Rows: $\det A = \sum_{j=1}^n a_{ij} A_{ji}$ (expansion with respect to the i th row).

(e) Laplace Expansion by Columns: $\det A = \sum_{j=1}^n A_{ij}a_{ji}$ (expansion with respect to the i th column).

Proof. See, e.g., [Str08, Th. 18.6]. ■

Theorem A.52. Let F be a field, $n \in \mathbb{N}$, let V be an n -dimensional vector space over F , and $A \in \mathcal{L}(V, V)$. Moreover, let $B_1 = \{v_1, \dots, v_n\}$ and $B_2 = \{w_1, \dots, w_n\}$ be bases of V . If $M = (m_{ji})$ is the transformation matrix corresponding to A with respect to B_1 and $N = (n_{ji})$ is the transformation matrix corresponding to A with respect to B_2 (i.e., for each $i \in \{1, \dots, n\}$, $A(v_i) = \sum_{j=1}^n m_{ji} v_j$ and $A(w_i) = \sum_{j=1}^n n_{ji} w_j$, cf. Def. and Rem. A.33), then $\det(M) = \det(N)$.

Proof. See, e.g., [Str08, Th. 17.11(a)]. ■

Definition and Remark A.53. Let F be a field, $n \in \mathbb{N}$, and let V be an n -dimensional vector space over F . Then Th. A.52 allows to define a determinant function for linear maps by

$$\det : \mathcal{L}(V, V) \longrightarrow F, \quad \det(A) := \det(M), \quad (\text{A.94})$$

where M is a transformation matrix for A with respect to an arbitrary basis of V . Then Th. A.32 shows that Th. A.48(a),(b),(d) yield the following properties of the new determinant function defined in (A.94):

- (a) If $A \in \mathcal{L}(V, V)$, then $\det(A) = 0$ if, and only if, A is not bijective. If A is bijective, then $\det(A^{-1}) = (\det(A))^{-1}$.
- (b) If $A, B \in \mathcal{M}(n, F)$, then $\det(AB) = \det(A) \det(B)$.
- (c) If $A \in \mathcal{L}(V, V)$ and $\lambda \in F$, then $\det(\lambda A) = \lambda^n \det(A)$. In particular, \det is *not* linear for $n > 1$.

In [Str08, §17], the author actually defines the determinant first for linear maps $A \in \mathcal{L}(V, V)$, establishes properties including the above properties, and only then defines the determinant function for square matrices. Several alternative, but equivalent, approaches are possible and can be found in the literature.

B Metric Spaces

B.1 Metric Subspaces

Definition B.1. If (X, d) is a metric space, $M \subseteq X$, then (M, d) is called a *metric subspace* of (X, d) (if d is understood, one also speaks of M as a metric subspace of X). Thus, the metric on the subspace M is just the metric on X restricted to M .

Remark B.2. One sees immediately that a metric subspace (M, d) of a metric space (X, d) is, indeed, a metric space: Since d satisfies the laws (i) – (iii) from Def. 1.17 for all $x, y, z \in X$, in particular, d satisfies the same laws for all $x, y, z \in M \subseteq X$.

Definition B.3. Let (X, d) be a metric space, and let (M, d) be a metric subspace of (X, d) . Is $A \subseteq M$ open with respect to (M, d) , then one says that A is *open* in M or *M-open* or *relatively open*. For $A \subseteq M$ closed with respect to (M, d) , one introduces analogous terms. Moreover, for $x \in M$, $r > 0$, call

$$B_{r,M}(x) := M \cap B_r(x) = \{y \in M : d(x, y) < r\}, \quad (\text{B.1})$$

the *open M-ball* with radius r and center x .

Caveat B.4. One has to use care when working with a subspace (M, d) of a metric space (X, d) : As will be seen in Ex. B.5, the notions and properties with respect to M are in general very different from the corresponding notions and properties with respect to X . For example, a set that is M -open might not be X -open and a set that is M -closed might not be X -closed!

Example B.5. (a) If (M, d) is a metric subspace of a metric space (X, d) , then, according to Lem. 1.27(b), M is always both M -open and M -closed (irrespective of M being X -open or X -closed).

(b) Let $X = \mathbb{R}$ with the usual metric, i.e. $d(x, y) = |x - y|$ for each $x, y \in \mathbb{R}$. Let $M = [0, 1]$. According to (a), is both M -closed and M -open, even though $[0, 1]$ is not open in X . When noting before that \mathbb{Q} and $]0, 1]$ are metric spaces that are not complete, we already considered metric subspaces of \mathbb{R} without making use of the term subspace. If $M =]0, 1]$, then, again, M is both M -closed and M -open, even though $]0, 1]$ is neither closed nor open in X . Moreover, $]0, \frac{1}{2}]$ is M -closed (but not X -closed) and $[\frac{1}{2}, 1]$ is M -open (but not X -open).

Proposition B.6. Let (M, d) be a metric subspace of a metric space (X, d) .

- (a) A subset A of M is M -open if, and only if, there is a set $O \subseteq X$ which is X -open and $A = O \cap M$.
- (b) A subset A of M is M -closed if, and only if, there is a set $C \subseteq X$ which is X -closed and $A = C \cap M$.

Proof. (a): Suppose A is M -open. Then, for each $a \in A$, there is $\epsilon_a > 0$ such that the open M -ball $B_{\epsilon_a, M}(a)$ is contained in A , i.e.

$$B_{\epsilon_a, M}(a) = M \cap B_{\epsilon_a}(a) \subseteq A. \quad (\text{B.2})$$

Let $O := \bigcup_{a \in A} B_{\epsilon_a}(a)$. Then O is X -open by Th. 1.29(a). Moreover,

$$O \cap M = \bigcup_{a \in A} (M \cap B_{\epsilon_a}(a)) = \bigcup_{a \in A} B_{\epsilon_a, M}(a) = A, \quad (\text{B.3})$$

where the last equality is due to (B.2) and the fact that $a \in A$ implies $a \in B_{\epsilon_a, M}(a)$.

Conversely, if $O \subseteq X$ is X -open and $A = O \cap M$, then each $a \in A$ is an X -interior point of O , i.e. there is $\epsilon > 0$ such that the open X -ball $B_\epsilon(a)$ is contained in O , i.e.

$B_\epsilon(a) \subseteq O$. Intersecting with M yields $B_{\epsilon,M}(a) = M \cap B_\epsilon(a) \subseteq M \cap O = A$, i.e. the open M -ball $B_{\epsilon,M}(a)$ is contained in A , showing that a is an M -interior point of A . As a was an arbitrary point of A , A is M -open.

(b): If A is M -closed, then $M \setminus A$ is M -open. According to (a), there is an X -open set $O \subseteq X$ such that $M \setminus A = M \cap O$. Then $C := X \setminus O$ is an X -closed set and $M \cap C = M \cap (X \setminus O) = M \setminus (M \cap O) = M \setminus (M \setminus A) = A$.

Conversely, if there is an X -closed set $C \subseteq X$ with $A = C \cap M$, then $O := X \setminus C$ is an X -open set satisfying $O \cap M = M \cap (X \setminus C) = M \setminus (C \cap M) = M \setminus A$. Thus, according to (a), $M \setminus A$ is M -open, i.e. A is M -closed. ■

B.2 Norm-Preserving and Isometric Maps

Definition B.7. (a) Given normed vector spaces $(X, \|\cdot\|_X)$ and $(Y, \|\cdot\|_Y)$ over \mathbb{K} , a function $f : X \rightarrow Y$ is called *norm-preserving* if, and only if,

$$\|f(x)\|_Y = \|x\|_X \quad \text{for each } x \in X. \quad (\text{B.4})$$

(b) Given metric spaces (X, d_X) and (Y, d_Y) , a function $f : X \rightarrow Y$ is called *distance-preserving* or *isometric* if, and only if,

$$d_Y(f(x), f(y)) = d_X(x, y) \quad \text{for each } x, y \in X. \quad (\text{B.5})$$

Lemma B.8. *Given normed vector spaces $(X, \|\cdot\|_X)$ and $(Y, \|\cdot\|_Y)$ over \mathbb{K} , a \mathbb{K} -linear function $f : X \rightarrow Y$ is norm-preserving if, and only if, f is isometric with respect to the induced metrics.*

Proof. The function f is norm-preserving if, and only if, $\|f(x)\|_Y = \|x\|_X$ for each $x \in X$. This, in turn is the case if, and only if,

$$\|f(x) - f(y)\|_Y = \|f(x - y)\|_Y = \|x - y\|_X \quad \text{for each } x, y \in X, \quad (\text{B.6})$$

where it was used that f is linear. As (B.6) states that f is isometric with respect to the induced metrics, the proof is complete. ■

The following examples show that the assertion of Lem. B.8 becomes false if the word “linear” is omitted.

Example B.9. (a) Let $(X, \|\cdot\|_X)$ be a normed vector space over \mathbb{K} , and $f : X \rightarrow \mathbb{K}$, $f(x) := \|x\|_X$. If we take $\|\cdot\|_Y$ to be the usual norm on \mathbb{K} , i.e. $\|y\|_Y := |y|$, then, for each $x \in X$, $\|f(x)\|_Y = \|\|x\|_X\| = \|x\|_X$, i.e. f is norm-preserving. However, if $\dim X > 0$ (i.e. if $X \neq \{0\}$), then f is *not* isometric with respect to the induced metrics: Take any $0 \neq x \in X$. One computes

$$\|f(x) - f(-x)\|_Y = \|\|x\|_X - \|x\|_X\| = 0 \neq \|x - (-x)\|_X = 2\|x\|_X. \quad (\text{B.7})$$

- (b) Consider $(X, \|\cdot\|_X)$, $(Y, \|\cdot\|_Y)$, where $X = Y = \mathbb{K}$ and $\|x\|_X = \|x\|_Y = |x|$ for each $x \in \mathbb{K}$. Then $f : X \rightarrow Y$, $f(x) := 1 + x$, is isometric due to $|f(x) - f(y)| = |1 + x - (1 + y)| = |x - y|$, but f is *not* norm-preserving, since $0 = |0| \neq |f(0)| = 1$.

Lemma B.10. *Isometric functions between metric spaces are one-to-one (in particular, isometric functions between normed spaces are one-to-one).*

Proof. Let (X, d_X) and (Y, d_Y) be metric spaces, and let $f : X \rightarrow Y$ be an isometric function, i.e. $d_Y(f(x), f(y)) = d_X(x, y)$ for each $x, y \in X$. If $x \neq y$, then $0 \neq d_X(x, y) = d_Y(f(x), f(y))$. Thus, $f(x) \neq f(y)$, showing that f is one-to-one. ■

Example B.11. If a function between normed spaces is just norm-preserving, but not isometric, then this function is not necessarily one-to-one: To see this, we reemploy the function f from Ex. B.9(a), i.e. let $(X, \|\cdot\|_X)$ be a normed vector space over \mathbb{K} , $\dim X > 0$, and $f : X \rightarrow \mathbb{K}$, $f(x) := \|x\|_X$. In Ex. B.9(a), we saw that f is norm-preserving, but not isometric. Since, for $x \neq 0$, one has $x \neq -x$, but $f(x) = \|x\|_X = f(-x)$, f is not one-to-one.

Remark B.12. If $(X, \|\cdot\|)$ is a normed space, d is the induced metric, and $M \subseteq X$, then (M, d) can be considered as the metric subspace of (X, d) according to Def. B.3. Thus, every subset of a normed space is turned into a metric space in a natural way. It is quite remarkable that actually *every* metric space arises in this way. That means, given any metric space (M, d) , there exists a normed space $(X, \|\cdot\|)$ and an isometric (in particular, one-to-one) function $f : M \rightarrow X$: One can choose X as the \mathbb{R} -vector space of bounded functions from M into \mathbb{R} with the sup-norm (for $F \in X$, define $\|F\| := \sup\{|F(x)| : x \in M\}$) and $f : M \rightarrow X$, $f(x) = f_x$, where $f_x : M \rightarrow \mathbb{R}$, $f_x(y) = d(x, y) - d(x_0, y)$ with some fixed $x_0 \in M$. However, the normed space X can be very large (i.e. much larger than M), and, thus, in practice, it is not always useful to study X in order to learn more about the metric space M .

B.3 Uniform Continuity and Lipschitz Continuity

This section provides some additional important results regarding uniformly continuous functions (see Def. 1.49(b)) and Lipschitz continuous functions (see Def. 1.49(c)). We start with an auxiliary result:

Lemma B.13. *If f, g are real-valued functions on a set X , i.e. if $f, g : X \rightarrow \mathbb{R}$, then, for each $x, y \in X$,*

$$|\max(f, g)(x) - \max(f, g)(y)| \leq \max\{|f(x) - f(y)|, |g(x) - g(y)|\}, \quad (\text{B.8a})$$

$$|\min(f, g)(x) - \min(f, g)(y)| \leq \max\{|f(x) - f(y)|, |g(x) - g(y)|\}. \quad (\text{B.8b})$$

Proof. By possibly switching the names of f and g , one can assume, without loss of generality, that $\max(f, g)(x) = f(x)$, i.e. $g(x) \leq f(x)$. If $g(y) \leq f(y)$ as well, then

$|\max(f, g)(x) - \max(f, g)(y)| = |f(x) - f(y)|$ and $|\min(f, g)(x) - \min(f, g)(y)| = |g(x) - g(y)|$, i.e. (B.8) is true. If $g(y) > f(y)$, then

$$|\max(f, g)(x) - \max(f, g)(y)| = |f(x) - g(y)| \leq \begin{cases} \leq |g(x) - g(y)| & \text{for } f(x) \leq g(y), \\ < f(x) - f(y) & \text{for } f(x) > g(y), \end{cases} \quad (\text{B.9a})$$

$$|\min(f, g)(x) - \min(f, g)(y)| = |g(x) - f(y)| \leq \begin{cases} < |g(x) - g(y)| & \text{for } g(x) \leq f(y), \\ \leq f(x) - f(y) & \text{for } g(x) > f(y), \end{cases} \quad (\text{B.9b})$$

showing that (B.8) holds in all cases. ■

Theorem B.14. *Let (X, d) be a metric space (e.g. a normed space), $(Y, \|\cdot\|)$ a normed vector space over \mathbb{K} , and assume that $f, g : X \rightarrow Y$ are uniformly continuous. Then $f + g$ and λf are uniformly continuous for each $\lambda \in \mathbb{K}$, i.e. the set of all uniformly continuous functions from X into Y constitutes a subspace of the vector space $\mathcal{F}(X, Y)$ over \mathbb{K} . Moreover, if $Y = \mathbb{K} = \mathbb{R}$, then $\max(f, g)$, $\min(f, g)$, f^+ , f^- , $|f|$ are all uniformly continuous.*

Proof. As f and g are uniformly continuous, given $\epsilon > 0$, there exist $\delta_f > 0$ and $\delta_g > 0$ such that, for each $x, y \in X$,

$$d(x, y) < \delta_f \Rightarrow \|f(x) - f(y)\| < \epsilon/2, \quad (\text{B.10a})$$

$$d(x, y) < \delta_g \Rightarrow \|g(x) - g(y)\| < \epsilon/2. \quad (\text{B.10b})$$

Thus, if $d(x, y) < \min\{\delta_f, \delta_g\}$, then

$$\|(f + g)(x) - (f + g)(y)\| \leq \|f(x) - f(y)\| + \|g(x) - g(y)\| < \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon, \quad (\text{B.10c})$$

showing that $f + g$ is uniformly continuous. Next, if $\lambda = 0$, then $\lambda f \equiv 0$, and obviously uniformly continuous. For $\lambda \neq 0$, choose $\delta > 0$ such that $d(x, y) < \delta$ implies $\|f(x) - f(y)\| < \epsilon/|\lambda|$. Then

$$\|(\lambda f)(x) - (\lambda f)(y)\| = |\lambda| \|f(x) - f(y)\| < |\lambda| \frac{\epsilon}{|\lambda|} = \epsilon, \quad (\text{B.10d})$$

showing that λf is uniformly continuous. If $Y = \mathbb{K} = \mathbb{R}$, then $d(x, y) < \min\{\delta_f, \delta_g\}$ together with Lem. B.13 implies

$$|\max(f, g)(x) - \max(f, g)(y)| < \epsilon/2 < \epsilon, \quad (\text{B.10e})$$

$$|\min(f, g)(x) - \min(f, g)(y)| < \epsilon/2 < \epsilon, \quad (\text{B.10f})$$

showing the uniform continuity of $\max(f, g)$ and $\min(f, g)$ and, in turn, also of f^+ , f^- , and $|f|$. ■

Theorem B.15. *Let (X, d) be a metric space (e.g. a normed space), $(Y, \|\cdot\|)$ a normed vector space over \mathbb{K} , and assume that $f, g : X \rightarrow Y$ are Lipschitz continuous. Then $f+g$ and λf are Lipschitz continuous for each $\lambda \in \mathbb{K}$, i.e. the set $\text{Lip}(X, Y)$ constitutes a subspace of the vector space $\mathcal{F}(X, Y)$ over \mathbb{K} . Moreover, if $Y = \mathbb{K} = \mathbb{R}$, then $\max(f, g)$, $\min(f, g)$, f^+ , f^- , $|f|$ are all Lipschitz continuous.*

Proof. As f and g are Lipschitz continuous, there exist $L_f \geq 0$ and $L_g \geq 0$ such that, for each $x, y \in X$,

$$\|f(x) - f(y)\| \leq L_f d(x, y), \quad (\text{B.11a})$$

$$\|g(x) - g(y)\| \leq L_g d(x, y). \quad (\text{B.11b})$$

Thus,

$$\begin{aligned} \|(f+g)(x) - (f+g)(y)\| &\leq \|f(x) - f(y)\| + \|g(x) - g(y)\| \\ &\leq L_f d(x, y) + L_g d(x, y) = (L_f + L_g)d(x, y), \end{aligned} \quad (\text{B.11c})$$

showing that $f+g$ is Lipschitz continuous with Lipschitz constant $L_f + L_g$. Next, for $\lambda \in \mathbb{K}$,

$$\|(\lambda f)(x) - (\lambda f)(y)\| = |\lambda| \|f(x) - f(y)\| \leq |\lambda| L_f d(x, y), \quad (\text{B.11d})$$

showing that λf is Lipschitz continuous with Lipschitz constant $|\lambda|L_f$. For $Y = \mathbb{K} = \mathbb{R}$, Lem. B.13 shows $\max(f, g)$ and $\min(f, g)$ are Lipschitz continuous with Lipschitz constant $\max\{L_f, L_g\}$, f^+ and f^- are Lipschitz continuous with Lipschitz constant L_f , and $|f|$ is Lipschitz continuous with Lipschitz constant $2L_f$. ■

Caveat B.16. Products and quotients of uniformly continuous functions are not necessarily uniformly continuous; products and quotients of Lipschitz continuous functions are not necessarily Lipschitz continuous: Even though $f \equiv 1$ and $g(x) = x$ are Lipschitz continuous, it was shown in Examples 1.52(a),(b), respectively, that f/g and g^2 are not even uniformly continuous on \mathbb{R}^+ .

B.4 Viewing \mathbb{C}^n as \mathbb{R}^{2n}

Remark B.17. Recall that the set of complex numbers \mathbb{C} is *defined* to be \mathbb{R}^2 , where the imaginary unit is $i := (0, 1) \in \mathbb{R}^2$, which allows to write each $z = (x, y) \in \mathbb{C} = \mathbb{R}^2$ as $z = x + iy$, where $x = \text{Re } z$ and $y = \text{Im } z$. This, for each $n \in \mathbb{N}$, gives rise to the \mathbb{R} -linear bijective map

$$I : \mathbb{C}^n \rightarrow \mathbb{R}^{2n}, \quad I((x_1, y_1), \dots, (x_n, y_n)) := (x_1, y_1, \dots, x_n, y_n), \quad (\text{B.12})$$

allowing to canonically identify \mathbb{C}^n with \mathbb{R}^{2n} .

—

The identification (B.12) allows the identification of metric structures on \mathbb{C}^n and \mathbb{R}^{2n} due to the following general result:

Proposition B.18. *Let X, Y be sets, let $d : X \times X \longrightarrow \mathbb{R}_0^+$ be a metric on X , and let $I : X \longrightarrow Y$ be bijective. Then*

$$d_Y : Y \times Y \longrightarrow \mathbb{R}_0^+, \quad d_Y(x, y) := d(I^{-1}(x), I^{-1}(y)), \quad (\text{B.13})$$

defines a metric on Y such that (X, d) and (Y, d_Y) are isometric (with the map I providing the isometry).

Proof. Let $x, y, z \in Y$. Then

$$d_Y(x, y) = 0 \quad \Leftrightarrow \quad d(I^{-1}(x), I^{-1}(y)) = 0 \quad \Leftrightarrow \quad I^{-1}(x) = I^{-1}(y) \quad \Leftrightarrow \quad x = y, \quad (\text{B.14})$$

showing that d_Y is positive definite. Moreover,

$$d_Y(x, y) = d(I^{-1}(x), I^{-1}(y)) = d(I^{-1}(y), I^{-1}(x)) = d_Y(y, x), \quad (\text{B.15})$$

showing d_Y is symmetric. Finally,

$$\begin{aligned} d_Y(x, z) &= d(I^{-1}(x), I^{-1}(z)) \leq d(I^{-1}(x), I^{-1}(y)) + d(I^{-1}(y), I^{-1}(z)) \\ &= d_Y(x, y) + d_Y(y, z), \end{aligned} \quad (\text{B.16})$$

proving the triangle inequality for d_Y and completing the proof that d_Y constitutes a metric. That I provides an isometry between (X, d) and (Y, d_Y) is immediate from (B.13). ■

Corollary B.19. *Let $n \in \mathbb{N}$, let $d : \mathbb{C}^n \times \mathbb{C}^n \longrightarrow \mathbb{R}_0^+$ be a metric, and let I be the map from (B.12). Then*

$$d_r : \mathbb{R}^{2n} \times \mathbb{R}^{2n} \longrightarrow \mathbb{R}_0^+, \quad d_r(x, y) := d(I^{-1}(x), I^{-1}(y)), \quad (\text{B.17})$$

defines a metric on \mathbb{R}^{2n} such that (\mathbb{C}^n, d) and (\mathbb{R}^{2n}, d_r) are isometric (with the map I providing the isometry). Moreover, the map $d \mapsto d_r$ is bijective between the set of metrics on \mathbb{C}^n and the set of metrics on \mathbb{R}^{2n} . ■

Proposition B.20. *Let $n \in \mathbb{N}$. If $\|\cdot\|$ constitutes a norm on the vector space \mathbb{C}^n over \mathbb{C} in the sense of Def. 1.19, then*

$$\|\cdot\|_r : \mathbb{R}^{2n} \longrightarrow \mathbb{R}_0^+, \quad \|(x_1, y_1, \dots, x_n, y_n)\|_r := \|((x_1, y_1), \dots, (x_n, y_n))\| \quad (\text{B.18})$$

defines a norm on the vector space \mathbb{R}^{2n} over \mathbb{R} such that $(\mathbb{C}^n, \|\cdot\|)$ and $(\mathbb{R}^{2n}, \|\cdot\|_r)$ are isometric (with the map I from (B.12) providing the isometry – even more precisely, if d and d_r denote the respective induced metrics, then the relation between d and d_r is given by (B.17)).

Proof. Exercise. ■

Example B.21. Let $n \in \mathbb{N}$, $p \in [1, \infty]$, and let $\|\cdot\|$ denote the p -norm on the vector space \mathbb{R}^n over \mathbb{R} , i.e. $\|x\| := (\sum_{j=1}^n |x_j|^p)^{1/p}$ for $p < \infty$ and $\|x\| = \max\{|x_j| : j = 1, \dots, n\}$ for $p = \infty$. Then it is an exercise to show

$$\|\cdot\|_c : \mathbb{C}^n \longrightarrow \mathbb{R}_0^+, \quad \|(z_1, \dots, z_n)\|_c := \|(|z_1|, \dots, |z_n|)\| \quad (\text{B.19})$$

defines a norm on the vector space \mathbb{C}^n over \mathbb{C} .

Remark B.22. As a consequence of Th. 1.95, every norm on the normed vector space \mathbb{C}^n over \mathbb{C} generates precisely the same open subsets of \mathbb{C}^n – in other words, there is only one *norm topology* on \mathbb{C}^n . Analogously, there is only one norm topology on \mathbb{R}^n as every norm on the normed vector space \mathbb{R}^n over \mathbb{R} generates precisely the same open subsets of \mathbb{R}^n . Moreover, Prop. B.20 shows that the open sets of the norm topology on \mathbb{C}^n are actually precisely the same as the open sets of the norm topology on \mathbb{R}^{2n} .

Theorem B.23. Let $n \in \mathbb{N}$, $A \subseteq \mathbb{C}^n$. Then A is bounded in the normed vector space \mathbb{C}^n over \mathbb{C} if, and only if, A is bounded in the normed vector space \mathbb{R}^{2n} over \mathbb{R} .

Proof. Exercise. ■

B.5 Banach Fixed Point Theorem a.k.a. Contraction Mapping Principle

Definition B.24. Let $\emptyset \neq A$ be a subset of a metric space (X, d) . A map $\varphi : A \longrightarrow A$ is called a *contraction* if, and only if, there exists $0 \leq L < 1$ satisfying

$$d(\varphi(x), \varphi(y)) \leq L d(x, y) \quad \text{for each } x, y \in A. \quad (\text{B.20})$$

Remark B.25. According to Def. B.24, $\varphi : A \longrightarrow A$ is a contraction if, and only if, φ is Lipschitz continuous with Lipschitz constant $L < 1$.

The following Th. B.26 constitutes the Banach fixed point theorem. It is also known as the contraction mapping principle. Its proof is surprisingly simple, e.g. about an order of magnitude easier than the proof of the Brouwer fixed point theorem.

Theorem B.26 (Banach Fixed Point Theorem). Let $\emptyset \neq A$ be a closed subset of a complete metric space (X, d) (for example, a Banach space). If $\varphi : A \longrightarrow A$ is a contraction with Lipschitz constant $0 \leq L < 1$, then φ has a unique fixed point $x_* \in A$. Moreover, for each initial value $x_0 \in A$, the sequence $(x_n)_{n \in \mathbb{N}_0}$, defined by

$$x_{n+1} := \varphi(x_n) \quad \text{for each } n \in \mathbb{N}_0, \quad (\text{B.21})$$

converges to x_* :

$$\lim_{n \rightarrow \infty} \varphi^n(x_0) = x_*. \quad (\text{B.22})$$

Furthermore, for each such sequence, we have the error estimate

$$d(x_n, x_*) \leq \frac{L}{1-L} d(x_n, x_{n-1}) \leq \frac{L^n}{1-L} d(x_1, x_0) \quad (\text{B.23})$$

for each $n \in \mathbb{N}$.

Proof. We start with uniqueness: Let $x_*, x_{**} \in A$ be fixed points of φ . Then

$$d(x_*, x_{**}) = d(\varphi(x_*), \varphi(x_{**})) \leq L d(x_*, x_{**}), \quad (\text{B.24})$$

which implies $1 \leq L$ for $d(x_*, x_{**}) > 0$. Thus, $L < 1$ implies $d(x_*, x_{**}) = 0$ and $x_* = x_{**}$.

Next, we turn to existence. A simple induction on $m - n$ shows

$$\begin{aligned} d(x_{m+1}, x_m) &\leq L d(x_m, x_{m-1}) \leq L^{m-n} d(x_{n+1}, x_n) \\ &\text{for each } m, n \in \mathbb{N}_0, m > n. \end{aligned} \quad (\text{B.25})$$

This, in turn, allows us to estimate, for each $n, k \in \mathbb{N}_0$:

$$\begin{aligned} d(x_{n+k}, x_n) &\leq \sum_{m=n}^{n+k-1} d(x_{m+1}, x_m) \stackrel{(\text{B.25})}{\leq} \sum_{m=n}^{n+k-1} L^{m-n} d(x_{n+1}, x_n) \\ &\leq \frac{1}{1-L} d(x_{n+1}, x_n) \stackrel{(\text{B.25})}{\leq} \frac{L^n}{1-L} d(x_1, x_0) \rightarrow 0 \quad \text{for } n \rightarrow \infty, \end{aligned} \quad (\text{B.26})$$

establishing that $(x_n)_{n \in \mathbb{N}_0}$ constitutes a Cauchy sequence. Since X is complete, this Cauchy sequence must have a limit $x_* \in X$, and since the sequence is in A and A is closed, $x_* \in A$. The continuity of φ allows to take limits in (B.21), resulting in $x_* = \varphi(x_*)$, showing that x_* is a fixed point and proving existence.

Finally, the error estimate (B.23) follows from (B.26) by fixing n and taking the limit for $k \rightarrow \infty$. ■

Example B.27. Suppose, we are looking for a fixed point of the map $\varphi(x) = \cos x$ (or, equivalently, for a zero of $f(x) = \cos x - x$). To apply the Banach fixed point theorem, we need to restrict φ to a set A such that $\varphi(A) \subseteq A$. This is the case for $A := [0, 1]$. Moreover, $\varphi : A \rightarrow A$ is a contraction, due to $\sin 1 < 1$ and the mean value theorem providing $\tau \in]0, 1[$, satisfying

$$|\varphi(x) - \varphi(y)| = |\varphi'(\tau)| |x - y| < (\sin 1) |x - y| \quad (\text{B.27})$$

for each $x, y \in A$. Since \mathbb{R} is complete and A is closed in \mathbb{R} , the Banach fixed point theorem yields the existence of a unique fixed point $x_* \in [0, 1]$ and $\lim \varphi^n(x_0) = x_*$ for each $x_0 \in [0, 1]$.

B.6 Unit Balls in Normed Spaces

The goal of this section is to prove that a normed vector space is finite-dimensional if, and only if, its closed unit ball is compact (see Th. B.29). In preparation, we show that finite-dimensional normed vector spaces are always closed:

Theorem B.28. *Let $(X, \|\cdot\|)$ be a normed vector space over \mathbb{K} . If $U \subseteq X$ is a subspace such that $\dim U = n \in \mathbb{N}$, then U is closed.*

Proof. Let (b_1, \dots, b_n) be a basis of U . Then

$$A : U \longrightarrow \mathbb{K}^n, \quad A \left(\sum_{k=1}^n \alpha_k b_k \right) := (\alpha_1, \dots, \alpha_n), \quad (\text{B.28})$$

defines a linear isomorphism (cf. Th. A.24). We define a norm on \mathbb{K}^n by letting

$$\|\cdot\| : \mathbb{K}^n \longrightarrow \mathbb{R}_0^+, \quad \|z\| := \|A^{-1}(z)\|. \quad (\text{B.29})$$

Indeed, (B.29) defines a norm: $\|0\| = \|A^{-1}(0)\| = \|0\| = 0$; if $z \in \mathbb{K}^n$ and $\|z\| = \|A^{-1}(z)\| = 0$, then $A^{-1}(z) = 0$, i.e. $z = 0$, showing $\|\cdot\|$ to be positive definite. Moreover

$$\forall_{z \in \mathbb{K}^n} \quad \forall_{\lambda \in \mathbb{K}} \quad \|\lambda z\| = \|A^{-1}(\lambda z)\| = |\lambda| \|A^{-1}(z)\| = |\lambda| \|z\|, \quad (\text{B.30})$$

showing $\|\cdot\|$ to be homogeneous of degree 1. Finally,

$$\forall_{z, w \in \mathbb{K}^n} \quad \|z + w\| = \|A^{-1}(z + w)\| \leq \|A^{-1}(z)\| + \|A^{-1}(w)\| = \|z\| + \|w\|, \quad (\text{B.31})$$

showing the triangle inequality to hold for $\|\cdot\|$.

Let $(u^k)_{k \in \mathbb{N}}$ be a sequence in U such that $\lim_{k \rightarrow \infty} u^k = x \in X$. Then $(u^k)_{k \in \mathbb{N}}$ is a Cauchy sequence and, as A is norm-preserving in consequence of (B.29), $(Au^k)_{k \in \mathbb{N}}$ is a Cauchy sequence in \mathbb{K}^n . Since \mathbb{K}^n is complete, there is $z \in \mathbb{K}^n$ such that $\lim_{k \rightarrow \infty} Au^k = z$ and $\lim_{k \rightarrow \infty} u^k = A^{-1}z$, showing $x = A^{-1}z \in U$, i.e. U is closed. ■

Theorem B.29. *A normed vector space $(X, \|\cdot\|)$ over \mathbb{K} is finite-dimensional if, and only if, its closed unit ball $\overline{B}_1(0)$ is compact.*

Proof. Let X be finite-dimensional. If (b_1, \dots, b_n) denotes a basis of X , then (B.28) defines a linear isomorphism $A : X \longrightarrow \mathbb{K}^n$. If we define a norm on \mathbb{K}^n via (B.29), then A^{-1} becomes norm-preserving and, in particular, continuous. Then $\overline{B}_1(0)$ in X must be compact as the continuous image (under A^{-1}) of $\overline{B}_1(0)$ in \mathbb{K}^n .

Conversely, let X be infinite-dimensional. To show that $\overline{B}_1(0)$ is not compact, we construct, via recursion, a sequence $(x^k)_{k \in \mathbb{N}}$ in $\overline{B}_1(0)$ (actually in the sphere $S_1(0)$) that does not have a convergent subsequence: Fix $n \in \mathbb{N}$ and assume (x^1, \dots, x^n) to be already constructed such that

$$\forall_{k \in \{1, \dots, n\}} \quad \|x^k\| = 1, \quad (\text{B.32a})$$

$$\forall_{\substack{k, l \in \{1, \dots, n\}, \\ k \neq l}} \quad \|x^k - x^l\| \geq \frac{1}{2}. \quad (\text{B.32b})$$

Let $U := \text{span}\{x^1, \dots, x^n\}$. Since X is infinite-dimensional, we have $U \neq X$. Let $x \in X \setminus U$. Since U is closed by Th. B.28, it is

$$d := \inf \{\|x - u\| : u \in U\} > 0. \quad (\text{B.33})$$

Moreover, there exists $u_0 \in U$ such that $\|x - u_0\| \leq 2d$. Set

$$x^{n+1} := \frac{x - u_0}{\|x - u_0\|}. \quad (\text{B.34})$$

Then $\|x^{n+1}\| = 1$ and, for each $u \in U$ is $\|x - u_0\|u + u_0 \in U$, implying

$$\|u - x^{n+1}\| = \frac{\| \|x - u_0\|u - x + u_0 \|}{\|x - u_0\|} \geq \frac{d}{\|x - u_0\|} \geq \frac{1}{2}. \quad (\text{B.35})$$

Thus, (B.32) holds with n replaced by $n + 1$, where (B.32b) means that $(x^k)_{k \in \mathbb{N}}$ can not have a convergent subsequence. \blacksquare

C Differential Calculus in \mathbb{R}^n

C.1 Proof of the Chain Rule

Proof of Th. 2.28. As usual, we first consider the case $\mathbb{K} = \mathbb{R}$. Since f is differentiable at ξ and g is differentiable at $f(\xi)$, according to Lem. 2.21, there are functions $r_f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $r_g : \mathbb{R}^m \rightarrow \mathbb{R}^p$ satisfying

$$r_f(h) = f(\xi + h) - f(\xi) - Df(\xi)(h), \quad (\text{C.1a})$$

$$r_g(h) = g(f(\xi) + h) - g(f(\xi)) - Dg(f(\xi))(h) \quad (\text{C.1b})$$

for each h such that $\|h\|_2$ is sufficiently small, as well as

$$\lim_{h \rightarrow 0} \frac{r_f(h)}{\|h\|_2} = 0, \quad \lim_{h \rightarrow 0} \frac{r_g(h)}{\|h\|_2} = 0. \quad (\text{C.2})$$

Defining $r_{g \circ f} : \mathbb{R}^n \rightarrow \mathbb{R}^p$ by

$$r_{g \circ f}(h) := \begin{cases} (g \circ f)(\xi + h) - (g \circ f)(\xi) - \left(Dg(f(\xi)) \circ Df(\xi) \right)(h) & \text{for } \xi + h \in G_f, \\ 0 & \text{otherwise,} \end{cases} \quad (\text{C.3})$$

it remains to show

$$\lim_{h \rightarrow 0} \frac{r_{g \circ f}(h)}{\|h\|_2} = 0. \quad (\text{C.4})$$

For each $h \in \mathbb{R}^n$ with $\|h\|_2$ sufficiently small, we use (C.1) to compute

$$\begin{aligned} (g \circ f)(\xi + h) &= g\left(f(\xi) + Df(\xi)(h) + r_f(h)\right) \\ &= g(f(\xi)) + Dg(f(\xi))\left(Df(\xi)(h) + r_f(h)\right) + r_g\left(Df(\xi)(h) + r_f(h)\right), \end{aligned} \quad (\text{C.5a})$$

implying

$$r_{g \circ f}(h) = Dg(f(\xi))(r_f(h)) + r_g(Df(\xi)(h) + r_f(h)). \quad (\text{C.5b})$$

From Th. 1.67, we know that $Dg(f(\xi))$ is Lipschitz continuous with some Lipschitz constant $L_g \in \mathbb{R}_0^+$. Thus, for each $0 \neq h \in \mathbb{R}^n$,

$$0 \leq \frac{\|Dg(f(\xi))(r_f(h))\|_2}{\|h\|_2} \leq \frac{L_g \|r_f(h)\|_2}{\|h\|_2}, \quad (\text{C.6a})$$

implying

$$\lim_{h \rightarrow 0} \frac{\|Dg(f(\xi))(r_f(h))\|_2}{\|h\|_2} = 0 \quad (\text{C.6b})$$

due to (C.2). Thus, to prove (C.4), it merely remains to show

$$\lim_{h \rightarrow 0} \frac{\|r_g(Df(\xi)(h) + r_f(h))\|_2}{\|h\|_2} = 0. \quad (\text{C.7})$$

To that end, we rewrite, for $Df(\xi)(h) + r_f(h) \neq 0$,

$$\frac{\|r_g(Df(\xi)(h) + r_f(h))\|_2}{\|h\|_2} = \frac{\|Df(\xi)(h) + r_f(h)\|_2 \|r_g(Df(\xi)(h) + r_f(h))\|_2}{\|h\|_2 \|Df(\xi)(h) + r_f(h)\|_2}. \quad (\text{C.8a})$$

Next, note

$$\lim_{h \rightarrow 0} \|Df(\xi)(h) + r_f(h)\|_2 = 0 \quad \stackrel{(\text{C.2})}{\Rightarrow} \quad \lim_{h \rightarrow 0} \frac{\|r_g(Df(\xi)(h) + r_f(h))\|_2}{\|Df(\xi)(h) + r_f(h)\|_2} = 0. \quad (\text{C.8b})$$

Once again, from Th. 1.67, we know that $Df(\xi)$ is Lipschitz continuous with some Lipschitz constant $L_f \in \mathbb{R}_0^+$, implying

$$\frac{\|Df(\xi)(h) + r_f(h)\|_2}{\|h\|_2} \leq \frac{\|Df(\xi)(h)\|_2 + \|r_f(h)\|_2}{\|h\|_2} \leq L_f + 1 \quad (\text{C.8c})$$

for $0 \neq \|h\|_2$ sufficiently small. Combining (C.8a) – (C.8c) proves (C.7) and, thus, (C.4). Together with (C.3) and Lem. 2.21, this shows that $g \circ f$ is differentiable at ξ with $D(g \circ f)(\xi) = Dg(f(\xi)) \circ Df(\xi)$.

In the case $\mathbb{K} = \mathbb{C}$, we can apply the case $\mathbb{K} = \mathbb{R}$ to obtain the differentiability of $\text{Re}(g \circ f) = (\text{Re } g) \circ f$ and of $\text{Im}(g \circ f) = (\text{Im } g) \circ f$ at ξ , and, in consequence, the differentiability of $g \circ f$ at ξ . Moreover, to verify the chain rule, we use the chain rule of the case $\mathbb{K} = \mathbb{R}$ to compute

$$\begin{aligned} D(g \circ f)(\xi) &= D \text{Re}(g \circ f)(\xi) + i D \text{Im}(g \circ f)(\xi) \\ &= D((\text{Re } g) \circ f)(\xi) + i D((\text{Im } g) \circ f)(\xi) \\ &= D \text{Re } g(f(\xi)) \circ Df(\xi) + i D \text{Im } g(f(\xi)) \circ Df(\xi) \\ &= Dg(f(\xi)) \circ Df(\xi), \end{aligned} \quad (\text{C.9})$$

thereby completing the proof. ■

C.2 Bounded Derivatives Imply Lipschitz Continuity

First, we provide an \mathbb{R}^m -valued variant of Th. 2.35:

Theorem C.1. *Let $m, n \in \mathbb{N}$, let $G \subseteq \mathbb{R}^n$ be open, and let $f : G \rightarrow \mathbb{R}^m$ be differentiable. Suppose there exists $M \in \mathbb{R}_0^+$ such that $|\partial_k f_l(\xi)| \leq M$ for each $k \in \{1, \dots, n\}$, each $l \in \{1, \dots, m\}$, and each $\xi \in G$. If G is convex, then f is Lipschitz continuous with Lipschitz constant $L := mM$ with respect to the 1-norms on \mathbb{R}^n and \mathbb{R}^m and with Lipschitz constant cL , $c > 0$, with respect to arbitrary norms on \mathbb{R}^n and \mathbb{R}^m .*

Proof. According to Th. 2.35, each f_l is M -Lipschitz with respect to the 1-norm on \mathbb{R}^n . Thus, we obtain, for each $x, y \in G$,

$$\|f(y) - f(x)\|_1 = \sum_{l=1}^m |f_l(y) - f_l(x)| \leq mM\|y - x\|_1, \quad (\text{C.10})$$

showing that, with respect to the 1-norms on \mathbb{R}^n and \mathbb{R}^m , f is Lipschitz continuous with Lipschitz constant mM . Since all norms on \mathbb{R}^n and \mathbb{R}^m are equivalent, we also get that f is Lipschitz continuous with Lipschitz constant cL , $c > 0$, with respect to all other norms on \mathbb{R}^n and \mathbb{R}^m . ■

It is sometimes useful if the bound on the derivatives is the same as the resulting Lipschitz constant (which, for $m > 1$, is not the case in the above Th. C.1). The following Th. C.3 provides a variant, where the constants are the same, formulated for functions $f : I \rightarrow \mathbb{R}^n$, defined on open intervals $I \subseteq \mathbb{R}$, and making use of the Euclidean norm $\|\cdot\|_2$ on \mathbb{R}^n . We will start with some auxiliary results regarding the Euclidean norm and the Euclidean inner product:

Proposition C.2. *Let $I \subseteq \mathbb{R}$ be an open interval, and let $g, h : I \rightarrow \mathbb{R}^n$ be differentiable, $n \in \mathbb{N}$.*

(a) *The function*

$$f : I \rightarrow \mathbb{R}, \quad f(x) := g(x) \bullet h(x) = \sum_{j=1}^n g_j(x)h_j(x), \quad (\text{C.11})$$

is differentiable and

$$f' : I \rightarrow \mathbb{R}, \quad f'(x) = g'(x) \bullet h(x) + h(x) \bullet h'(x). \quad (\text{C.12})$$

(b) *The function*

$$\alpha : I \rightarrow \mathbb{R}, \quad \alpha(x) := \|g(x)\|_2 = \sqrt{g(x) \bullet g(x)}, \quad (\text{C.13})$$

is differentiable at each $x \in I$ such that $g(x) \neq 0$. Moreover,

$$\forall_{\substack{x \in I, \\ g(x) \neq 0}} \quad \alpha'(x) = \frac{g(x) \bullet g'(x)}{\alpha(x)} = \frac{g(x) \bullet g'(x)}{\|g(x)\|_2}. \quad (\text{C.14})$$

Proof. (a) is immediate from the product rule.

(b) is an easy consequence of (a), as (a) implies α to be differentiable at each $x \in I$ such that $g(x) \neq 0$, and

$$\forall_{\substack{x \in I, \\ g(x) \neq 0}} \quad \alpha'(x) = \frac{2g(x) \bullet g'(x)}{2\sqrt{g(x) \bullet g(x)}} = \frac{g(x) \bullet g'(x)}{\alpha(x)}, \quad (\text{C.15})$$

completing the proof. ■

Theorem C.3. *Let $a, b \in \mathbb{R}$ with $a < b$ and let $f :]a, b[\rightarrow \mathbb{R}^n$ be differentiable with uniformly bounded derivative, i.e. with*

$$\exists_{M \in \mathbb{R}_0^+} \quad \forall_{x \in]a, b[} \quad \|f'(x)\|_2 = \sqrt{\sum_{j=1}^n |f'_j(x)|^2} \leq M. \quad (\text{C.16})$$

Then f is M -Lipschitz, i.e.

$$\forall_{x_1, x_2 \in]a, b[} \quad \|f(x_1) - f(x_2)\|_2 \leq M |x_1 - x_2|. \quad (\text{C.17})$$

Proof. For $x_1 = x_2$, there is nothing to prove. Thus, assume $x_1 \neq x_2$ and define the auxiliary function

$$g : [0, 1] \rightarrow \mathbb{R}^n, \quad g(t) := f(x_1 + t(x_2 - x_1)) - f(x_1). \quad (\text{C.18})$$

According to the chain rule of Th. 2.28, g is differentiable on $]0, 1[$ and

$$\forall_{t \in]0, 1[} \quad g'(t) = (x_2 - x_1) f'(x_1 + t(x_2 - x_1)), \quad (\text{C.19})$$

implying

$$\forall_{t \in]0, 1[} \quad \|g'(t)\|_2 \leq M |x_1 - x_2|. \quad (\text{C.20})$$

We now introduce another auxiliary function, namely

$$\alpha : [0, 1] \rightarrow \mathbb{R}, \quad \alpha(t) := \|g(t)\|_2. \quad (\text{C.21})$$

Then α is continuous (as f and the norm are both continuous), satisfying $\alpha(0) = \|g(0)\|_2 = 0$ and $\alpha(1) = \|g(1)\|_2 = \|f(x_2) - f(x_1)\|_2$. If $\alpha(1) = 0$, then (C.17) is trivially true, and, thus, we proceed to assume $\alpha(1) > 0$. Then the continuity of α implies

$$s := \sup \{t \in [0, 1] : \alpha(t) = 0\} < 1, \quad \alpha(s) = 0. \quad (\text{C.22})$$

In consequence, α is positive on $]s, 1[$ and, thus, differentiable on $]s, 1[$ by Prop. C.2(b). The mean value theorem [Phi15a, Th. 9.17] implies the existence of $\sigma \in]s, 1[$ such that

$$\begin{aligned} \alpha(1) &= \alpha(1) - \alpha(s) = (1-s) \alpha'(\sigma) \stackrel{(\text{C.14})}{=} (1-s) \frac{g(\sigma) \bullet g'(\sigma)}{\alpha(\sigma)} \\ &\stackrel{(1.81)}{\leq} (1-s) \frac{\|g(\sigma)\|_2 \|g'(\sigma)\|_2}{\|g(\sigma)\|_2} \stackrel{(\text{C.20})}{\leq} (1-s) M |x_1 - x_2| \\ &\leq M |x_1 - x_2|, \end{aligned} \quad (\text{C.23})$$

which establishes the case. ■

C.3 Surjectivity of Directional Derivatives

We finish the proof of Th. 2.38 by showing that, for $n \geq 2$, the map

$$D : S_1(0) \longrightarrow [-\alpha, \alpha], \quad D(e) := \nabla f(\xi) \cdot e = \sum_{j=1}^n \epsilon_j \partial_j f(\xi), \quad \alpha = \|\nabla f(\xi)\|_2, \quad (\text{C.24})$$

is surjective (we already know from (2.65) that $D(e) \in [-\alpha, \alpha]$ for each $e \in S_1(0)$). We also recall $e_{\max} = \nabla f(\xi)/\alpha$, $e_{\min} = -e_{\max}$, $D(e_{\max}) = \alpha$, $D(e_{\min}) = -\alpha$.

The idea is to rotate e_{\max} into e_{\min} . This can be achieved using a suitable function

$$\rho : [0, \pi] \longrightarrow S_1(0) \subseteq \mathbb{R}^n, \quad \rho = (\rho_1, \dots, \rho_n).$$

We have to define ρ differently, depending on $n \geq 2$ being even or odd. To this end, let $(\epsilon_1, \dots, \epsilon_n) := e_{\max}$. If n is even, then define

$$\forall_{j \in \{1, \dots, n\}} \quad \rho_j : [0, \pi] \longrightarrow [-1, 1], \quad \rho_j(\theta) := \begin{cases} \epsilon_j \cos \theta + \epsilon_{j+1} \sin \theta & \text{if } j \text{ is odd,} \\ -\epsilon_{j-1} \sin \theta + \epsilon_j \cos \theta & \text{if } j \text{ is even;} \end{cases} \quad (\text{C.25a})$$

if n is odd (note $n \geq 3$ in this case), then define

$$\forall_{j \in \{1, \dots, n\}} \quad \rho_j : [0, \pi] \longrightarrow [-1, 1], \quad \rho_j(\theta) := \begin{cases} \epsilon_j \cos \theta + \epsilon_{j+1} \sin \theta & \text{if } j < n-2 \text{ is odd,} \\ -\epsilon_{j-1} \sin \theta + \epsilon_j \cos \theta & \text{if } j < n-2 \text{ is even,} \\ \epsilon_{n-2} \cos \theta + \sqrt{\epsilon_{n-1}^2 + \epsilon_n^2} \sin \theta & \text{if } j = n-2, \\ \epsilon_{n-1} \cos \theta - \frac{\epsilon_{n-2} \epsilon_{n-1}}{\sqrt{\epsilon_{n-1}^2 + \epsilon_n^2}} \sin \theta & \text{if } j = n-1, \\ \epsilon_n \cos \theta - \frac{\epsilon_{n-2} \epsilon_n}{\sqrt{\epsilon_{n-1}^2 + \epsilon_n^2}} \sin \theta & \text{if } j = n. \end{cases} \quad (\text{C.25b})$$

For the sake of readability, we assumed $\epsilon_{n-1} \neq 0$ or $\epsilon_n \neq 0$ in (C.25b). There is always at least one $j_0 \in \{1, \dots, n\}$ such that $\epsilon_{j_0} \neq 0$. If $j_0 \notin \{n-1, n\}$, then one merely needs to interchange the roles of j_0 and n in (C.25b).

Clearly, for every $n \geq 2$, each ρ_j is continuous, i.e. ρ is continuous.

Next, we verify that ρ , indeed, maps into $S_1(0)$ (which, in particular, implies each ρ_j maps into $[-1, 1]$): If $n \geq 2$ is even, then, for each odd $j \leq n-1$, one has

$$\begin{aligned} & (\rho_j(\theta))^2 + (\rho_{j+1}(\theta))^2 \\ &= (\epsilon_j \cos \theta + \epsilon_{j+1} \sin \theta)^2 + (-\epsilon_j \sin \theta + \epsilon_{j+1} \cos \theta)^2 \\ \forall_{\theta \in [0, \pi]} \quad &= \epsilon_j^2 \cos^2 \theta + 2\epsilon_j \epsilon_{j+1} \cos \theta \sin \theta + \epsilon_{j+1}^2 \sin^2 \theta \\ &+ \epsilon_j^2 \sin^2 \theta - 2\epsilon_j \epsilon_{j+1} \cos \theta \sin \theta + \epsilon_{j+1}^2 \cos^2 \theta \\ &= \epsilon_j^2 (\cos^2 \theta + \sin^2 \theta) + \epsilon_{j+1}^2 (\cos^2 \theta + \sin^2 \theta) = \epsilon_j^2 + \epsilon_{j+1}^2, \end{aligned} \quad (\text{C.26})$$

implying

$$\forall_{\theta \in [0, \pi]} \quad \|\rho(\theta)\|_2^2 = \sum_{j=1}^n (\rho_j(\theta))^2 = \sum_{j=1}^n \epsilon_j^2 = 1. \quad (\text{C.27})$$

If $n \geq 3$ is odd, then (C.26) still holds for each odd $j \leq n-4$. Additionally,

$$\begin{aligned}
& (\rho_{n-2}(\theta))^2 + (\rho_{n-1}(\theta))^2 + (\rho_n(\theta))^2 \\
&= \epsilon_1^2 \cos^2 \theta + 2\epsilon_1 \sqrt{\epsilon_2^2 + \epsilon_3^2} \sin \theta \cos \theta + (\epsilon_2^2 + \epsilon_3^2) \sin^2 \theta \\
&+ \epsilon_2^2 \cos^2 \theta - 2 \frac{\epsilon_1 \epsilon_2^2}{\sqrt{\epsilon_2^2 + \epsilon_3^2}} \sin \theta \cos \theta + \frac{\epsilon_1^2 \epsilon_2^2}{\epsilon_2^2 + \epsilon_3^2} \sin^2 \theta \\
&+ \epsilon_3^2 \cos^2 \theta - 2 \frac{\epsilon_1 \epsilon_3^2}{\sqrt{\epsilon_2^2 + \epsilon_3^2}} \sin \theta \cos \theta + \frac{\epsilon_1^2 \epsilon_3^2}{\epsilon_2^2 + \epsilon_3^2} \sin^2 \theta \\
\forall \theta \in [0, \pi] \quad &= (\epsilon_1^2 + \epsilon_2^2 + \epsilon_3^2) \cos^2 \theta \\
&+ \frac{2\epsilon_1 (\epsilon_2^2 + \epsilon_3^2 - \epsilon_2^2 - \epsilon_3^2)}{\sqrt{\epsilon_2^2 + \epsilon_3^2}} \sin \theta \cos \theta \\
&+ \left(\epsilon_2^2 + \epsilon_3^2 + \frac{\epsilon_1^2 (\epsilon_2^2 + \epsilon_3^2)}{\epsilon_2^2 + \epsilon_3^2} \right) \sin^2 \theta \\
&= (\epsilon_1^2 + \epsilon_2^2 + \epsilon_3^2) (\cos^2 \theta + \sin^2 \theta) = \epsilon_{n-2}^2 + \epsilon_{n-1}^2 + \epsilon_n^2,
\end{aligned} \tag{C.28}$$

i.e. (C.27) is true also for each $n \geq 3$ odd.

Clearly, D is also continuous and, thus, so is $D \circ \rho : [0, \pi] \rightarrow [-\alpha, \alpha]$. Moreover, as $\sin(0) = \sin(\pi) = 0$, $\cos(0) = 1$, $\cos(\pi) = -1$, we obtain

$$\forall_{n \geq 2} \quad \forall_{j \in \{1, \dots, n\}} \quad \left(\rho_j(0) = \epsilon_j \quad \wedge \quad \rho_j(\pi) = -\epsilon_j \right), \tag{C.29}$$

implying

$$\forall_{n \geq 2} \quad \left(\rho(0) = e_{\max} \quad \wedge \quad (D \circ \rho)(0) = \alpha \quad \wedge \quad \rho(\pi) = e_{\min} \quad \wedge \quad (D \circ \rho)(\pi) = -\alpha \right). \tag{C.30}$$

The continuity of $D \circ \rho$ and the intermediate value theorem [Phi15a, Th. 7.57] imply $D \circ \rho$ to be surjective, i.e. D must be surjective as well.

C.4 Implicit Function Theorem

We start with a preparatory proposition.

Proposition C.4. *Let $\|\cdot\|$ be some norm on \mathbb{R}^n , $n \in \mathbb{N}$. Moreover, let $a \in \mathbb{R}^n$, $r > 0$, and let $f : B_r(a) \rightarrow \mathbb{R}^n$ be defined on the open r -ball with center a with respect to $\|\cdot\|$. If A is an invertible $n \times n$ matrix over \mathbb{R} such that*

$$\|A^{-1}f(a)\| < \frac{r}{2} \tag{C.31}$$

and such that the map

$$F : B_r(a) \rightarrow \mathbb{R}^n, \quad F(x) := x - A^{-1}f(x), \tag{C.32}$$

is Lipschitz continuous with Lipschitz constant $L = 1/2$, then f has a unique zero $\xi \in B_r(a)$. Moreover, for each $x_0 \in B_r(a)$, ξ is the limit of the sequence $(x_k)_{k \in \mathbb{N}_0}$, recursively defined by

$$\forall_{k \in \mathbb{N}_0} \quad x_{k+1} := F(x_k). \quad (\text{C.33})$$

Proof. Set

$$s_0 := \max \left\{ 2 \|A^{-1}f(a)\|, \|x_0 - a\| \right\} \stackrel{(\text{C.31})}{\in} [0, r[. \quad (\text{C.34})$$

The idea is to show that, for each $s_0 < s < r$, the Banach fixed point Th. B.26 applies to the contraction

$$F_s : \overline{B}_s(a) \longrightarrow \overline{B}_s(a), \quad F_s(x) := F(x). \quad (\text{C.35})$$

We verify that F_s , indeed, maps $\overline{B}_s(a)$ into $\overline{B}_s(a)$: If $x \in \overline{B}_s(a)$, then

$$\begin{aligned} \|F(x) - a\| &\leq \|F(x) - F(a)\| + \|F(a) - a\| \leq \frac{1}{2}\|x - a\| + \|A^{-1}f(a)\| \\ &\leq \frac{s}{2} + \frac{s_0}{2} < s, \end{aligned} \quad (\text{C.36})$$

showing $F_s(x) \in \overline{B}_s(a)$ (in particular, this shows the x_k are well-defined by (C.33)). As F is Lipschitz continuous with Lipschitz constant $L = 1/2$, so is F_s , i.e. F_s is, indeed, a contraction. As $\overline{B}_s(a)$ is closed, the Banach fixed point Th. B.26 yields that F_s has a unique fixed point ξ and, moreover, $\xi = \lim_{k \rightarrow \infty} x_k$. Since this holds for each $s \in]s_0, r[$, ξ must also be the unique fixed point of F . The proof is concluded by noting

$$\forall_{y \in B_r(a)} \quad f(y) = 0 \quad \Leftrightarrow \quad F(y) = y - A^{-1}f(y) = y, \quad (\text{C.37})$$

that means y is a zero of f if, and only if, y is a fixed point of F . ■

Remark C.5. If the map f in Prop. C.4 is differentiable with invertible derivatives $Df(x)$, and if, instead of using a constant matrix A in the definition of (C.33), one uses $(Df(x_k))^{-1}$, then the iteration defined by (C.33) is known as *Newton's method* (in n dimensions, cf. [Phi15b, Sec. 6.3]). In consequence, if $A \approx (Df(x_k))^{-1}$ in (C.33), then the defined iteration is sometimes referred to as a *simplified Newton's method*.

Notation C.6. Let $k, m, n \in \mathbb{N}$, let $G \subseteq \mathbb{R}^n \times \mathbb{R}^m$ be open, and consider a map $f : G \longrightarrow \mathbb{R}^k$. If $(\xi, \eta) \in G$ and f is differentiable at (ξ, η) , then let $D_y f(\xi, \eta)$ and $D_x f(\xi, \eta)$ denote the linear maps

$$D_y f(\xi, \eta) : \mathbb{R}^m \longrightarrow \mathbb{R}^k, \quad (D_y f(\xi, \eta))(h) := (Df(\xi, \eta))(0, h), \quad (\text{C.38a})$$

$$D_x f(\xi, \eta) : \mathbb{R}^n \longrightarrow \mathbb{R}^k, \quad (D_x f(\xi, \eta))(h) := (Df(\xi, \eta))(h, 0), \quad (\text{C.38b})$$

respectively.

Theorem C.7 (Implicit Function Theorem). *Let $m, n \in \mathbb{N}$, let $G \subseteq \mathbb{R}^n \times \mathbb{R}^m$ be open, and let $f : G \longrightarrow \mathbb{R}^m$ be continuously differentiable, i.e. $f \in C^1(G, \mathbb{R}^m)$. If $(\xi, \eta) \in G$ is such that*

$$f(\xi, \eta) = 0 \quad \text{and} \quad A := D_y f(\xi, \eta) \text{ is invertible}, \quad (\text{C.39})$$

then there exist open neighborhoods $U_\xi \subseteq \mathbb{R}^n$ of ξ and $V_\eta \subseteq \mathbb{R}^m$ of η , and a continuously differentiable map $g : U_\xi \longrightarrow V_\eta$ such that the zeros of f in $U_\xi \times V_\eta$ are given precisely by the graph of g , i.e.

$$(U_\xi \times V_\eta) \cap f^{-1}\{0\} = \{(x, g(x)) : x \in U_\xi\}, \quad (\text{C.40a})$$

which can be restated as

$$\forall_{(x,y) \in U_\xi \times V_\eta} \quad \left(f(x, y) = 0 \iff y = g(x) \right). \quad (\text{C.40b})$$

Moreover,

$$\forall_{x \in U_\xi} \quad Dg(x) = -\left(D_y f(x, g(x))\right)^{-1} D_x f(x, g(x)) \quad (\text{C.41})$$

and, if $f \in C^\alpha(G, \mathbb{R}^m)$, $\alpha \in \mathbb{N} \cup \{\infty\}$, then $g \in C^\alpha(U_\xi, \mathbb{R}^m)$.

Proof. Fix some arbitrary norms on \mathbb{R}^n and on the set $\mathcal{M}(m, \mathbb{R})$ of real $m \times m$ matrices (for readability's sake, we will denote both norms by $\|\cdot\|$). On \mathbb{R}^m , we will use the 1-norm $\|\cdot\|_1$ to apply Th. C.1. According to the hypothesis, A is invertible. Thus, $\det(A) > 0$. Since the map $B \mapsto \det(B)$ is continuous (cf. Ex. 1.66(a)), and the map $D_y f : G \longrightarrow \mathcal{M}(m, \mathbb{R})$ is continuous due to the assumed continuous differentiability of f , the set

$$G_0 := \{(x, y) \in G : \det(D_y f(x, y)) > 0\} \subseteq G \quad (\text{C.42})$$

is an open neighborhood of (ξ, η) . Next, we consider the map

$$F : G_0 \longrightarrow \mathbb{R}^m, \quad F(x, y) := y - A^{-1}f(x, y). \quad (\text{C.43})$$

Then F is continuously differentiable with

$$D_y F : G_0 \longrightarrow \mathcal{M}(m, \mathbb{R}), \quad D_y F(x, y) = \text{Id} - A^{-1}D_y f(x, y), \quad (\text{C.44})$$

being continuous as well. Thus, since $D_y F(\xi, \eta) = \text{Id} - A^{-1}A = 0$, there exists $r > 0$ such that the open r -balls $B_r(\xi) \subseteq \mathbb{R}^n$ and $B_r(\eta) \subseteq \mathbb{R}^m$ satisfy

$$(\xi, \eta) \in B_r(\xi) \times B_r(\eta) \subseteq \left\{ (x, y) \in G_0 : \forall_{k,l=1,\dots,m} |\partial_{y_k} F_l(x, y)| < \frac{1}{2m} \right\} \subseteq G_0. \quad (\text{C.45})$$

As we assume f to be continuous with $f(\xi, \eta) = 0$, there exists $s \in]0, r]$ such that

$$B_s(\xi) \subseteq \left\{ x \in \mathbb{R}^n : \|A^{-1}f(x, \eta)\|_1 < \frac{r}{2} \right\} \subseteq B_r(\xi). \quad (\text{C.46})$$

To construct the map $g : B_s(\xi) \longrightarrow B_r(\eta)$, we fix $x \in B_s(\xi)$ and apply Prop. C.4 to the function

$$f_x : B_r(\eta) \longrightarrow \mathbb{R}^m, \quad f_x(y) := f(x, y). \quad (\text{C.47})$$

To verify that the hypotheses of Prop. C.4 are satisfied, we observe $\|A^{-1}f_x(\eta)\|_1 < \frac{r}{2}$ holds due to $x \in B_s(\xi)$ and (C.46), the map $F_x : B_r(\eta) \longrightarrow \mathbb{R}^m$, $F_x(y) := y - A^{-1}f_x(y) = F(x, y)$, is Lipschitz continuous with Lipschitz constant $L = m \frac{1}{2m} = \frac{1}{2}$ due

to (C.45) and Th. C.1. Thus, according to Prop. C.4, the function f_x has a unique zero $g(x)$ in $B_r(\eta)$, which defines the function g .

Note that, in the above argument, for each $0 < \rho < r$, one can choose $s(\rho) < s$ such that (C.46) holds with s replaced by $s(\rho)$ and r replaced by ρ , then showing that g maps $B_{s(\rho)}$ into B_ρ . We now choose some arbitrary $\rho \in]0, r[$ and set

$$U_\xi := B_{s(\rho)}(\xi), \quad V_\eta := B_\rho(\eta) \quad (\text{C.48})$$

for the desired neighborhoods of the theorem. We verify g to be continuous on U_ξ : Let $x \in U_\xi$ and let $(x_k)_{k \in \mathbb{N}}$ be a sequence in U_ξ with $\lim_{k \rightarrow \infty} x_k = x$. We have to show $\lim_{k \rightarrow \infty} g(x_k) = g(x)$. If $\lim_{k \rightarrow \infty} g(x_k) = g(x)$ does not hold, then, without loss of generality, we may assume that there exists $\epsilon > 0$ such that $\|g(x_k) - g(x)\| > \epsilon$ for each $k \in \mathbb{N}$ (after having replaced $(x_k)_{k \in \mathbb{N}}$ with a suitable subsequence). Moreover, we may replace $(x_k)_{k \in \mathbb{N}}$ with another subsequence such that there exists $y \in \overline{B}_\rho(\eta) \subseteq B_r(\eta)$ satisfying $y = \lim_{k \rightarrow \infty} g(x_k)$ (this is due to the Bolzano-Weierstrass Th. 1.16(b), as $g(x_k) \in B_\rho(\eta)$ for each $k \in \mathbb{N}$). Then the continuity of f implies

$$f(x, y) = \lim_{k \rightarrow \infty} f(x_k, g(x_k)) = 0, \quad (\text{C.49})$$

showing $g(x) = y = \lim_{k \rightarrow \infty} g(x_k)$ (due to (C.40b) – here we need $y \in B_r(\eta)$, which was the reason for choosing $\rho < r$), which is in contradiction to the choice of the x_k , and proves the continuity of g .

Next, we show that g is differentiable at each $x \in U_\xi$, where the derivative is given by (C.41): To this end, let $x \in U_\xi$ and note the existence of $(D_y f(x, g(x)))^{-1}$ due to $(x, g(x)) \in U_\xi \times V_\eta \subseteq G_0$. According to Def. 2.19, we have to show

$$\lim_{h \rightarrow 0} \frac{g(x+h) - g(x) + (D_y f(x, g(x)))^{-1} D_x f(x, g(x)) h}{\|h\|} = 0. \quad (\text{C.50})$$

Let $0 \neq h \in \mathbb{R}^n$ be sufficiently small such that $x+h \in U_\xi$. Using the notation of the mean value Th. 2.32, for each $l \in \{1, \dots, m\}$, there exist $x_{h,l} \in S_{x,x+h}$ and $y_{h,l} \in S_{g(x),g(x+h)}$ such that

$$\begin{aligned} 0 &= f_l(x+h, g(x+h)) - f_l(x, g(x)) \\ &= f_l(x+h, g(x+h)) - f_l(x, g(x+h)) + f_l(x, g(x+h)) - f_l(x, g(x)) \\ &= D_x f_l(x_{h,l}, g(x+h))(h) + D_y f_l(x, y_{h,l})(g(x+h) - g(x)). \end{aligned} \quad (\text{C.51})$$

Note that the two derivatives occurring in (C.51) have the form of gradients, which, according to our usual convention, we can interpret as row vectors. Joining m row vectors into a matrix, we can write the m equations of (C.51) in matrix form as

$$0 = X_h h + Y_h (g(x+h) - g(x)), \quad (\text{C.52})$$

where

$$X_h := \begin{pmatrix} D_x f_1(x_{h,1}, g(x+h)) \\ \vdots \\ D_x f_m(x_{h,m}, g(x+h)) \end{pmatrix}, \quad Y_h := \begin{pmatrix} D_y f_1(x, y_{h,1}) \\ \vdots \\ D_y f_m(x, y_{h,m}) \end{pmatrix}. \quad (\text{C.53})$$

As we already know g to be continuous, $h \rightarrow 0$ implies $g(x+h) \rightarrow g(x)$. Thus, since $y_{h,l} \in S_{g(x),g(x+h)}$, $h \rightarrow 0$ implies $y_{h,l} \rightarrow g(x)$ for each $l \in \{1, \dots, m\}$, and, as all partials of f are continuous as well, $Y_h \rightarrow D_y f(x, g(x))$. Since the maps $B \mapsto \det(B)$ and $B \mapsto \|B^{-1}\|$ are continuous (cf. Ex. 1.53 and Ex. 1.66(a),(b)), $h \rightarrow 0$ implies $\det(Y_h) \rightarrow \det(D_y f(x, g(x))) \neq 0$ and Y_h is invertible for sufficiently small h with $(Y_h)^{-1} \rightarrow (D_y f(x, g(x)))^{-1}$. For such sufficiently small h , we can rewrite (C.52) as

$$g(x+h) - g(x) = -(Y_h)^{-1} X_h h. \quad (\text{C.54})$$

Also, since $x_{h,l} \in S_{x,x+h}$, $h \rightarrow 0$ implies $x_{h,l} \rightarrow x$ and, then, the continuity of g together with the continuity of the partials of f implies $X_h \rightarrow D_x f(x, g(x))$. Thus, we can finish the proof of (C.41) by noting

$$\begin{aligned} & \lim_{h \rightarrow 0} \frac{\|g(x+h) - g(x) + (D_y f(x, g(x)))^{-1} D_x f(x, g(x)) h\|_1}{\|h\|} \\ &= \lim_{h \rightarrow 0} \frac{\|-(Y_h)^{-1} X_h h + (D_y f(x, g(x)))^{-1} D_x f(x, g(x)) h\|_1}{\|h\|} = 0. \end{aligned} \quad (\text{C.55})$$

It remains to prove that g is C^α if f is C^α , $\alpha \in \mathbb{N} \cup \{\infty\}$. To this end, for $\alpha \in \mathbb{N}$, we will show by induction on $\beta = 1, \dots, \alpha$ that each partial derivative of g at $x \in U_\xi$ of order β is a rational function of partials of f of order $\leq \beta$, all taken at $(x, g(x))$, and of partials of g of order $\leq \beta - 1$, all taken at x (in particular, the denominator of this rational function does not have any zeros in U_ξ): For $\beta = 1$, the claim follows from (C.41): The entries of $D_x f(x, g(x))$ are polynomials of first partials of f taken at $(x, g(x))$; the entries of $(D_y f(x, g(x)))^{-1}$ are, according to Th. A.51(c), rational functions, where both the numerator and the denominator polynomial are polynomials of first partials of f taken at $(x, g(x))$ (in particular, the entries of the right-hand side of (C.41) do not involve any first partials of g). For the induction step, let $1 < \beta \leq \alpha$. By induction, we know the partials of g of order $\beta - 1$ are rational functions of partials of f of order $\leq \beta - 1$, all taken at $(x, g(x))$, and of partials of g of order $\leq \beta - 2$, all taken at x . Taking the derivative of partials of g of order $\leq \beta - 2$ evaluated at x , yields partials of g of order $\leq \beta - 1$ still evaluated at x ; according to the chain rule of Th. 2.28, taking the derivative of partials of f of order $\leq \beta - 1$ evaluated at $(x, g(x))$, yields polynomials of partials of f of order $\leq \beta$ evaluated at $(x, g(x))$ and of first partials of g evaluated at x . Thus, applying the product and the quotient rule establishes the case. ■

Theorem C.8 (Inverse Function Theorem). *Let $n \in \mathbb{N}$, let $G \subseteq \mathbb{R}^n$ be open, and let $f : G \rightarrow \mathbb{R}^n$ be continuously differentiable, i.e. $f \in C^1(G, \mathbb{R}^n)$. If $\xi \in G$ is such that*

$$Df(\xi) \text{ is invertible}, \quad (\text{C.56})$$

then there exists an open neighborhood $U \subseteq G$ of ξ such that $V := f(U)$ is open and the restriction $f : U \rightarrow V$ is bijective with continuously differentiable inverse function $f^{-1} : V \rightarrow U$. Moreover,

$$\forall_{y \in V} \quad D(f^{-1})(y) = \left(Df(f^{-1}(y)) \right)^{-1} \quad (\text{C.57})$$

and, if $f \in C^\alpha(U, \mathbb{R}^n)$, $\alpha \in \mathbb{N} \cup \{\infty\}$, then $f^{-1} \in C^\alpha(V, \mathbb{R}^n)$.

Proof. The idea is to apply the implicit function Th. C.7 to the continuously differentiable map

$$F : G \times \mathbb{R}^n \longrightarrow \mathbb{R}^n, \quad F(x, y) := f(x) - y. \quad (\text{C.58})$$

Here, as compared to Th. C.7, the roles of the variables x and y are switched. Letting $\eta := f(\xi)$, we have

$$F(\xi, \eta) = f(\xi) - \eta = 0, \quad \text{and} \quad D_x F(\xi, \eta) = Df(\xi) \text{ is invertible.} \quad (\text{C.59})$$

Thus, Th. C.7 applies and yields an open neighborhood $\tilde{U} \subseteq G$ of ξ , an open neighborhood $V \subseteq \mathbb{R}^n$ of η , and a C^1 map $g : V \longrightarrow \tilde{U}$ such that

$$\forall_{(x,y) \in \tilde{U} \times V} \quad \left(F(x, y) = f(x) - y = 0 \quad \Leftrightarrow \quad x = g(y) \right). \quad (\text{C.60})$$

If we let $U := g(V)$, then $U \subseteq \tilde{U}$ is a neighborhood of $\xi = g(f(\xi))$, and (C.60) implies that $f : U \longrightarrow V$ and $g : V \longrightarrow U$ are inverse to each other, in particular, they are both bijective with $f^{-1} = g$. To verify that U is open, consider the (still continuous) map $f : \tilde{U} \longrightarrow \mathbb{R}^n$ and observe $U = f^{-1}(V)$. As V is open, Th. 1.54(ii) implies the existence of $O \subseteq \mathbb{R}^n$ open with $U = O \cap \tilde{U}$. Since both O and \tilde{U} are open, U must be open as well.

Using (C.41), we obtain, for each $y \in V$,

$$\begin{aligned} Dg(y) &= -\left(D_x F(g(y), y)\right)^{-1} D_y F(g(y), y) \\ &= -\left(Df(g(y))\right)^{-1} (-\text{Id}) = \left(Df(g(y))\right)^{-1}, \end{aligned} \quad (\text{C.61})$$

proving (C.57). Finally, if f is C^α on U , then F is C^α on $U \times \mathbb{R}^n$, such that Th. C.7 implies $g = f^{-1}$ to be C^α as well. \blacksquare

Corollary C.9. *Let $n \in \mathbb{N}$, let $G \subseteq \mathbb{R}^n$ be open, and let $f : G \longrightarrow \mathbb{R}^n$ be continuously differentiable, i.e. $f \in C^1(G, \mathbb{R}^n)$. If $Df(x)$ is invertible for each $x \in G$, then f maps open sets to open sets, i.e. if $O \subseteq G$ is open, then $f(O)$ is open as well.*

Proof. Let $O \subseteq G$ be open. We have to show that each point $\eta \in f(O)$ is an interior point of $f(O)$. To this end, let $\eta \in f(O)$ and let $\xi \in O$ be such that $f(\xi) = \eta$. Since $Df(\xi)$ is invertible by hypothesis, we can apply the inverse function Th. C.8 to the restriction of f to O , obtaining open neighborhoods $U \subseteq O$ of ξ and $V \subseteq f(O)$ of η such that $f : U \longrightarrow V$ is bijective. In particular, η is an interior point of $f(O)$, proving $f(O)$ to be open. \blacksquare

D Riemann Integral for \mathbb{C} -Valued Functions

D.1 Riemann Integrability

Notation D.1. Let $I := [a, b] \subseteq \mathbb{R}^n$ be an interval, $a, b \in \mathbb{R}^n$, $a < b$. By $\mathcal{R}(I, \mathbb{R}) := \mathcal{R}(I)$ we denote the set of all Riemann integrable functions $f : I \longrightarrow \mathbb{R}$ (cf. Def. 4.5(b)).

Definition D.2. Let $I := [a, b] \subseteq \mathbb{R}^n$ be an interval, $a, b \in \mathbb{R}^n$, $a < b$. We call a function $f : I \rightarrow \mathbb{C}$ *Riemann integrable* if, and only if, both $\operatorname{Re} f$ and $\operatorname{Im} f$ are Riemann integrable. The set of all Riemann integrable functions $f : I \rightarrow \mathbb{C}$ is denoted by $\mathcal{R}(I, \mathbb{C})$. If $f \in \mathcal{R}(I, \mathbb{C})$, then

$$\int_I f := \left(\int_I \operatorname{Re} f, \int_I \operatorname{Im} f \right) = \int_I \operatorname{Re} f + i \int_I \operatorname{Im} f \in \mathbb{C} \quad (\text{D.1})$$

is called the Riemann integral of f over I .

Theorem D.3. Let $I := [a, b] \subseteq \mathbb{R}^n$, $a, b \in \mathbb{R}^n$, $n \in \mathbb{N}$, $a < b$, $f : I \rightarrow \mathbb{C}$. If f is continuous, then f is Riemann integrable over I .

Proof. If f is continuous, then $\operatorname{Re} f$ and $\operatorname{Im} f$ are both continuous, and, thus, the statement follows from the real-valued case of Th. 4.14. \blacksquare

Theorem D.4. Let $I := [a, b] \subseteq \mathbb{R}^n$, $a, b \in \mathbb{R}^n$, $n \in \mathbb{N}$, $a < b$.

- (a) If $f, g \in \mathcal{R}(I, \mathbb{C})$, then $\bar{f}, fg \in \mathcal{R}(I, \mathbb{C})$. If, in addition, there exists $\delta > 0$ such that $|g(x)| \geq \delta$ for each $x \in I$, then $f/g \in \mathcal{R}(I, \mathbb{C})$.
- (b) If $f \in \mathcal{R}(I, \mathbb{R})$ and $\phi : f(I) \rightarrow \mathbb{C}$ is Lipschitz continuous, then $\phi \circ f \in \mathcal{R}(I, \mathbb{C})$.
- (c) If $f \in \mathcal{R}(I, \mathbb{C})$ and $\phi : f(I) \rightarrow \mathbb{R}$ is Lipschitz continuous, then $\phi \circ f \in \mathcal{R}(I, \mathbb{R})$.

Proof. All the following proofs are completely analogous to the respective 1-dimensional case in [Phi15a, Th. G.4].

(a): Since

$$\bar{f} = (\operatorname{Re} f, -\operatorname{Im} f), \quad (\text{D.2a})$$

$$fg = (\operatorname{Re} f \operatorname{Re} g - \operatorname{Im} f \operatorname{Im} g, \operatorname{Re} f \operatorname{Im} g + \operatorname{Im} f \operatorname{Re} g), \quad (\text{D.2b})$$

$$1/g = (\operatorname{Re} g/|g|^2, -\operatorname{Im} g/|g|^2), \quad (\text{D.2c})$$

everything follows from the real-valued case of Th. 4.12(a) and of Th. 4.15(b),(c), where $|g| \geq \delta > 0$ guarantees $|g|^2 \geq \delta^2 > 0$.

(b): Assume ϕ to be L -Lipschitz, $L \geq 0$. For each $x, y \in f(I)$, one has

$$|\operatorname{Re} \phi(x) - \operatorname{Re} \phi(y)| \stackrel{[\text{Phi15a, Th. 5.11(d)}]}{\leq} |\phi(x) - \phi(y)| \leq L|x - y|, \quad (\text{D.3a})$$

$$|\operatorname{Im} \phi(x) - \operatorname{Im} \phi(y)| \stackrel{[\text{Phi15a, Th. 5.11(d)}]}{\leq} |\phi(x) - \phi(y)| \leq L|x - y|, \quad (\text{D.3b})$$

showing $\operatorname{Re} \phi$ and $\operatorname{Im} \phi$ are L -Lipschitz, such that $\operatorname{Re}(\phi \circ f)$ and $\operatorname{Im}(\phi \circ f)$ are Riemann integrable by Th. 4.15(a).

(c): Assume ϕ to be L -Lipschitz, $L \geq 0$. If $f \in \mathcal{R}(I, \mathbb{C})$, then $\operatorname{Re} f, \operatorname{Im} f \in \mathcal{R}(I, \mathbb{R})$, and, given $\epsilon > 0$, Riemann's integrability criterion of Th. 4.13 provides partitions Δ_1, Δ_2 of I such that $R(\Delta_1, \operatorname{Re} f) - r(\Delta_1, \operatorname{Re} f) < \epsilon/2L$, $R(\Delta_2, \operatorname{Im} f) - r(\Delta_2, \operatorname{Im} f) < \epsilon/2L$, where

R and r denote upper and lower Riemann sums, respectively (cf. (4.11)). Letting Δ be a joint refinement of Δ_1 and Δ_2 , we have (cf. Def. 4.8(a),(b) and Th. 4.10(a))

$$R(\Delta, \operatorname{Re} f) - r(\Delta, \operatorname{Re} f) < \epsilon/2L, \quad R(\Delta, \operatorname{Im} f) - r(\Delta, \operatorname{Im} f) < \epsilon/2L. \quad (\text{D.4})$$

Recalling that, for each $g : I \rightarrow \mathbb{R}$, it is

$$r(\Delta, g) = \sum_{p \in P(\Delta)} m_p(g) |I_p|, \quad (\text{D.5a})$$

$$R(\Delta, g) = \sum_{p \in P(\Delta)} M_p(g) |I_p|, \quad (\text{D.5b})$$

where $P(\Delta)$ is according to Def. 4.2,

$$m_p(g) := \inf\{g(x) : x \in I_p\}, \quad M_p(g) := \sup\{g(x) : x \in I_p\}, \quad (\text{D.5c})$$

we obtain, for each $\xi_p, \eta_p \in I_p$,

$$\begin{aligned} & |(\phi \circ f)(\xi_p) - (\phi \circ f)(\eta_p)| \\ & \leq L |f(\xi_p) - f(\eta_p)| \stackrel{[\text{Phi15a, Th. 5.11(d)}]}{\leq} L |\operatorname{Re} f(\xi_p) - \operatorname{Re} f(\eta_p)| + L |\operatorname{Im} f(\xi_p) - \operatorname{Im} f(\eta_p)| \\ & \leq L (M_p(\operatorname{Re} f) - m_p(\operatorname{Re} f)) + L (M_p(\operatorname{Im} f) - m_p(\operatorname{Im} f)), \end{aligned} \quad (\text{D.6})$$

and, thus,

$$\begin{aligned} R(\Delta, \phi \circ f) - r(\Delta, \phi \circ f) &= \sum_{p \in P(\Delta)} (M_p(\phi \circ f) - m_p(\phi \circ f)) |I_p| \\ &\stackrel{(\text{D.6})}{\leq} \sum_{p \in P(\Delta)} L (M_p(\operatorname{Re} f) - m_p(\operatorname{Re} f)) |I_p| + \sum_{p \in P(\Delta)} L (M_p(\operatorname{Im} f) - m_p(\operatorname{Im} f)) |I_p| \\ &= L (R(\Delta, \operatorname{Re} f) - r(\Delta, \operatorname{Re} f)) + L (R(\Delta, \operatorname{Im} f) - r(\Delta, \operatorname{Im} f)) \stackrel{(\text{D.4})}{<} \epsilon. \end{aligned} \quad (\text{D.7})$$

Thus, $\phi \circ f \in \mathcal{R}(I, \mathbb{R})$ by Th. 4.13. ■

Theorem D.5. Let $n \in \mathbb{N}$, $a, b \in \mathbb{R}^n$, $a < b$, $I := [a, b]$.

(a) The integral is linear: More precisely, if $f, g \in \mathcal{R}(I, \mathbb{C})$ and $\lambda, \mu \in \mathbb{C}$, then $\lambda f + \mu g \in \mathcal{R}(I, \mathbb{C})$ and

$$\int_I (\lambda f + \mu g) = \lambda \int_I f + \mu \int_I g. \quad (\text{D.8})$$

(b) For each $f \in \mathcal{R}(I, \mathbb{C})$, one has $|f| \in \mathcal{R}(I, \mathbb{R})$ and

$$\left| \int_I f \right| \leq \int_I |f|. \quad (\text{D.9})$$

Proof. (a): One computes, using the real-valued case of Th. 4.12(a),

$$\begin{aligned} \int_I (\lambda f) &= \left(\int_I (\operatorname{Re} \lambda \operatorname{Re} f - \operatorname{Im} \lambda \operatorname{Im} f), \int_I (\operatorname{Re} \lambda \operatorname{Im} f + \operatorname{Im} \lambda \operatorname{Re} f) \right) \\ &= \left(\operatorname{Re} \lambda \int_I \operatorname{Re} f - \operatorname{Im} \lambda \int_I \operatorname{Im} f, \operatorname{Re} \lambda \int_I \operatorname{Im} f + \operatorname{Im} \lambda \int_I \operatorname{Re} f \right) \\ &= \lambda \int_I f \end{aligned} \quad (\text{D.10a})$$

and

$$\begin{aligned} \int_I (f + g) &= \left(\int_I \operatorname{Re}(f + g), \int_I \operatorname{Im}(f + g) \right) = \left(\int_I \operatorname{Re} f + \int_I \operatorname{Re} g, \int_I \operatorname{Im} f + \int_I \operatorname{Im} g \right) \\ &= \left(\int_I \operatorname{Re} f, \int_I \operatorname{Im} f \right) + \left(\int_I \operatorname{Re} g, \int_I \operatorname{Im} g \right) = \int_I f + \int_I g. \end{aligned} \quad (\text{D.10b})$$

(b): As the modulus is 1-Lipschitz by the inverse triangle inequality, $|f| \in \mathcal{R}(I, \mathbb{R})$ by Th. D.4(c). Let Δ be an arbitrary partition of I . Then, using the notation from the proof of Th. D.4(c) above, we obtain the following estimate of intermediate Riemann sums (cf. (4.11c)):

$$\begin{aligned} \left| (\rho(\Delta, \operatorname{Re} f), \rho(\Delta, \operatorname{Im} f)) \right| &:= \left| \left(\sum_{p \in P(\Delta)} \operatorname{Re} f(\xi_p) |I_p|, \sum_{p \in P(\Delta)} \operatorname{Im} f(\xi_p) |I_p| \right) \right| \\ &\leq \sum_{p \in P(\Delta)} \left| (\operatorname{Re} f(\xi_p), \operatorname{Im} f(\xi_p)) \right| |I_p| \\ &= \sum_{p \in P(\Delta)} |f(\xi_p)| |I_p| = \rho(\Delta, |f|). \end{aligned} \quad (\text{D.11})$$

Since the intermediate Riemann sums in (D.11) converge to the respective integrals by (4.30b), one obtains

$$\left| \int_I f \right| = \lim_{|\Delta| \rightarrow 0} \left| (\rho(\Delta, \operatorname{Re} f), \rho(\Delta, \operatorname{Im} f)) \right| \stackrel{(\text{D.11})}{\leq} \lim_{|\Delta| \rightarrow 0} \rho(\Delta, |f|) = \int_I |f|, \quad (\text{D.12})$$

proving (D.9). ■

D.2 Fubini Theorem

Definition D.6. Let $I = [a, b] \subseteq \mathbb{R}^n$ be an interval, $a, b \in \mathbb{R}^n$, $n \in \mathbb{N}$, $a < b$, and suppose $f : I \rightarrow \mathbb{C}$ is bounded (i.e. both $\operatorname{Re} f$ and $\operatorname{Im} f$ are bounded). Define

$$J_*(f, I) := J_*(\operatorname{Re} f, I) + i J_*(\operatorname{Im} f, I), \quad (\text{D.13a})$$

$$J^*(f, I) := J^*(\operatorname{Re} f, I) + i J^*(\operatorname{Im} f, I). \quad (\text{D.13b})$$

As in the \mathbb{R} -valued case, we call $J_*(f, I)$ the *lower Riemann integral* of f over I and $J^*(f, I)$ the *upper Riemann integral* of f over I .

Theorem D.7. *Let $a, b, c, d, e, f \in \mathbb{R}^n$, $n \in \mathbb{N}$, $a < b$, $c < d$, $e < f$, $I = [a, b]$, $J = [c, d]$, $K = [e, f]$. If $I = J \times K$ and $f \in \mathcal{R}(I, \mathbb{C})$, then*

$$\int_I f = \int_I f(x, y) \, d(x, y) = \int_K \int_J f(x, y) \, dx \, dy = \int_J \int_K f(x, y) \, dy \, dx. \quad (\text{D.14})$$

As in the real-valued Th. 4.16, there is a slight abuse of notation in (D.14), as it can happen that a function $x \mapsto f(x, y)$ is not Riemann integrable over J and that a function $y \mapsto f(x, y)$ is not Riemann integrable over K . As in Th. 4.16, in that case, one can choose either the lower or the upper Riemann integral for the inner integrals in (D.14). Independently of the choice, the resulting function $y \mapsto \int_J f(x, y) \, dx$ is Riemann integrable over K , $x \mapsto \int_K f(x, y) \, dy$ is Riemann integrable over J , and the validity of (D.14) is unaffected. By applying (D.14) inductively, one obtains

$$\int_I f = \int_I f(x) \, dx = \int_{a_1}^{b_1} \cdots \int_{a_n}^{b_n} f(x_1, \dots, x_n) \, dx_n \cdots dx_1, \quad (\text{D.15})$$

where, for the inner integrals, one can choose the upper or lower Riemann integral, and one can also permute their order without changing the overall value.

Proof. We show how to obtain the present \mathbb{C} -valued case from the \mathbb{R} -valued case of Th. 4.16. One computes, using lower Riemann integrals for the inner integrals,

$$\begin{aligned} \int_I f &= \int_I \operatorname{Re} f + i \int_I \operatorname{Im} f \\ &\stackrel{\text{Th. 4.16}}{=} \int_K J_*(\operatorname{Re} f(\cdot, y), J) \, dy + i \int_K J_*(\operatorname{Im} f(\cdot, y), J) \, dy \\ &\stackrel{(\text{D.13a})}{=} \int_K J_*(f(\cdot, y), J) \, dy, \end{aligned} \quad (\text{D.16})$$

proving $\int_I f = \int_K \int_J f(x, y) \, dx \, dy$ with the inner integral interpreted as lower Riemann integral. Clearly, the same calculation works if the inner integral is interpreted as upper Riemann integral, and it also still works if J and K are switched, completing the proof of (D.14). As mentioned in the statement, (D.15) follows from (D.14) by induction. ■

D.3 Change of Variables

Theorem D.8. *Let $a, b, c, d \in \mathbb{R}^n$, $n \in \mathbb{N}$, $a < b$, $c < d$, $I := [a, b]$, $J := [c, d]$, $\phi : I \rightarrow \mathbb{R}^n$, $f : J \rightarrow \mathbb{C}$. If, on the interior of I , ϕ is one-to-one, Lipschitz continuous, and has continuous first partials, $\phi(I) \subseteq J$, and $(f \circ \phi)|\det J_\phi| \in \mathcal{R}(I)$, then $f\chi_{\phi(I)} \in \mathcal{R}(J)$ and the following change of variables formula holds:*

$$\int_J f\chi_{\phi(I)} = \int_I (f \circ \phi)|\det J_\phi|. \quad (\text{D.17})$$

Proof. We show how to obtain the present \mathbb{C} -valued case from the \mathbb{R} -valued case of Th. 4.18. One computes

$$\begin{aligned}
 \int_J f \chi_{\phi(I)} &= \int_J \operatorname{Re} f \chi_{\phi(I)} + i \int_J \operatorname{Im} f \chi_{\phi(I)} \\
 &\stackrel{(4.61)}{=} \int_I ((\operatorname{Re} f) \circ \phi) |\det J_\phi| + i \int_I ((\operatorname{Im} f) \circ \phi) |\det J_\phi| \\
 &= \int_I (f \circ \phi) |\det J_\phi|,
 \end{aligned} \tag{D.18}$$

thereby establishing the case. ■

References

- [Phi15a] P. PHILIP. *Calculus I for Computer Science and Statistics Students*. Lecture Notes, Ludwig-Maximilians-Universität, Germany, 2014/2015, available in PDF format at http://www.math.lmu.de/~philip/publications/lectureNotes/calc1_forInfAndStatStudents.pdf.
- [Phi15b] P. PHILIP. *Numerical Analysis I*. Lecture Notes, Ludwig-Maximilians-Universität, Germany, 2014/2015, available in PDF format at <http://www.math.lmu.de/~philip/publications/lectureNotes/numericalAnalysis.pdf>.
- [Str08] GERNOT STROTH. *Lineare Algebra*, 2nd ed. Berliner Studienreihe zur Mathematik, Vol. 7, Heldermann Verlag, Lemgo, Germany, 2008 (German).