

Chapter 2

Mathematical Background

This chapter collects the necessary preliminaries to develop the theory of GLT sequences. The reader who knows about measure/integration theory, general topology and matrix analysis is supposed to be familiar with most of the material presented herein. We will try, however, to be as much self-contained as possible, by proving some of the results that may not be so popular, and by providing precise bibliographic references for the results we do not prove.

2.1 Notation and Terminology

For the reader's convenience, we report in this section some of the most common notations and terminologies that will be used throughout this book. Together with the index at the end, this section can be used as a reference whenever an unknown notation/terminology is encountered.

- The cardinality of a set S is denoted by $\#S$.
- If S is any subset of a topological space, the closure of S is denoted by \bar{S} .
- $\mathbb{R}^{m \times n}$ (resp., $\mathbb{C}^{m \times n}$) is the space of real (resp., complex) $m \times n$ matrices.
- O_m and I_m denote, respectively, the $m \times m$ zero matrix and the $m \times m$ identity matrix. Sometimes, when the size m can be inferred from the context, O and I are used instead of O_m and I_m .
- If \mathbf{x} is a vector and X is a matrix, \mathbf{x}^T and \mathbf{x}^* (resp., X^T and X^*) are the transpose and the conjugate transpose of \mathbf{x} (resp., X).
- We use the abbreviations HPD, HPSD, SPD, SPSD for ‘Hermitian Positive Definite’, ‘Hermitian Positive SemiDefinite’, ‘Symmetric Positive Definite’, ‘Symmetric Positive SemiDefinite’.
- If $X, Y \in \mathbb{C}^{m \times m}$, the notation $X \geq Y$ (resp., $X > Y$) means that X, Y are Hermitian and $X - Y$ is HPSD (resp., HPD).
- If $X, Y \in \mathbb{C}^{m \times m}$, we denote by $X \circ Y$ the componentwise (or Hadamard) product of X and Y : $(X \circ Y)_{ij} = x_{ij}y_{ij}$, $i, j = 1, \dots, m$.

- If $X \in \mathbb{C}^{m \times m}$, we denote by X^\dagger the Moore–Penrose pseudoinverse of X . For more on the Moore–Penrose pseudoinverse, see Sect. 2.4.2.
- If $X \in \mathbb{C}^{m \times m}$, we denote by $\Lambda(X)$ the spectrum of X and by $\rho(X)$ the spectral radius of X , i.e., $\rho(X) = \max_{\lambda \in \Lambda(X)} |\lambda|$.
- If $X \in \mathbb{C}^{m \times m}$, we denote by $\lambda_j(X)$, $j = 1, \dots, m$, the eigenvalues of X . If the eigenvalues are real, their maximum and minimum are also denoted by $\lambda_{\max}(X)$ and $\lambda_{\min}(X)$.
- If $X \in \mathbb{C}^{m \times m}$, we denote by $\sigma_j(X)$, $j = 1, \dots, m$, the singular values of X . The maximum and minimum singular values are also denoted by $\sigma_{\max}(X)$ and $\sigma_{\min}(X)$.
- If $1 \leq p \leq \infty$, the symbol $|\cdot|_p$ denotes both the p -norm of vectors and the associated operator norm for matrices:

$$|\mathbf{x}|_p = \begin{cases} (\sum_{i=1}^m |x_i|^p)^{1/p}, & \text{if } 1 \leq p < \infty, \\ \max_{i=1, \dots, m} |x_i|, & \text{if } p = \infty, \end{cases} \quad \mathbf{x} \in \mathbb{C}^m,$$

$$|X|_p = \max_{\substack{\mathbf{x} \in \mathbb{C}^m \\ \mathbf{x} \neq \mathbf{0}}} \frac{|X\mathbf{x}|_p}{|\mathbf{x}|_p}, \quad X \in \mathbb{C}^{m \times m}.$$

The 2-norm $|\cdot|_2$ is also known as the spectral (or Euclidean) norm and it will be preferably denoted by $\|\cdot\|$. For more on p -norms, see Sect. 2.4.1.

- Given $X \in \mathbb{C}^{m \times m}$ and $1 \leq p \leq \infty$, $\|X\|_p$ denotes the Schatten p -norm of X , which is defined as the p -norm of the vector $(\sigma_1(X), \dots, \sigma_m(X))$ formed by the singular values of X . The Schatten 1-norm is also known under the names of trace-norm and nuclear norm. For more on Schatten p -norms, see Sect. 2.4.3.
- $\Re(X)$ and $\Im(X)$ are, respectively, the real and the imaginary part of the square matrix X :

$$\Re(X) = \frac{X + X^*}{2}, \quad \Im(X) = \frac{X - X^*}{2i},$$

where i is the imaginary unit ($i^2 = -1$). Note that $\Re(X)$, $\Im(X)$ are Hermitian and $X = \Re(X) + i \Im(X)$ for all square matrices X .

- If $z \in \mathbb{C}$ and $\varepsilon > 0$, we denote by $D(z, \varepsilon)$ the disk with center z and radius ε , i.e., $D(z, \varepsilon) = \{w \in \mathbb{C} : |w - z| < \varepsilon\}$. If $S \subseteq \mathbb{C}$ and $\varepsilon > 0$, we denote by $D(S, \varepsilon)$ the ε -expansion of S , which is defined as $D(S, \varepsilon) = \bigcup_{z \in S} D(z, \varepsilon)$.
- The symbol ‘something $\xrightarrow{t \rightarrow \tau}$ something else’ means that ‘something’ tends to ‘something else’ as $t \rightarrow \tau$.
- Given two sequences $\{\zeta_n\}_n$ and $\{\xi_n\}_n$, with $\zeta_n \geq 0$ and $\xi_n > 0$ for all n , the notation $\zeta_n = O(\xi_n)$ means that there exists a constant C , independent of n , such that $\zeta_n \leq C\xi_n$ for all n ; and the notation $\zeta_n = o(\xi_n)$ means that $\zeta_n/\xi_n \rightarrow 0$ as $n \rightarrow \infty$.
- $C_c(\mathbb{C})$ (resp., $C_c(\mathbb{R})$) is the space of complex-valued continuous functions defined on \mathbb{C} (resp., \mathbb{R}) with bounded support. Moreover, for $m \in \mathbb{N} \cup \{\infty\}$, $C_c^m(\mathbb{R}) = C_c(\mathbb{R}) \cap C^m(\mathbb{R})$, where $C^m(\mathbb{R})$ is the space of functions $F : \mathbb{R} \rightarrow \mathbb{C}$ such that the real and imaginary parts $\Re(F)$, $\Im(F)$ are of class C^m over \mathbb{R} in the classical sense.
- If $w_i : D_i \rightarrow \mathbb{C}$, $i = 1, \dots, d$, are arbitrary functions, $w_1 \otimes \dots \otimes w_d : D_1 \times \dots \times D_d \rightarrow \mathbb{C}$ is the tensor-product function:

$$(w_1 \otimes \cdots \otimes w_d)(\xi_1, \dots, \xi_d) = w_1(\xi_1) \cdots w_d(\xi_d)$$

for all $(\xi_1, \dots, \xi_d) \in D_1 \times \cdots \times D_d$.

- If $f : D \rightarrow E$ and $g : E \rightarrow F$ are arbitrary functions, the composite function $g \circ f$ is preferably denoted by $g(f)$.
- If $g : D \rightarrow \mathbb{C}$, we set $\|g\|_\infty = \sup_{\xi \in D} |g(\xi)|$. If we need/want to specify the domain D , we write $\|g\|_{\infty, D}$ instead of $\|g\|_\infty$. Clearly, $\|g\|_\infty < \infty$ if and only if g is bounded over its domain.
- If $g : D \rightarrow \mathbb{C}$ is continuous over D , with $D \subseteq \mathbb{C}^k$ for some k , we denote by $\omega_g(\cdot)$ the modulus of continuity of g ,

$$\omega_g(\delta) = \sup_{\substack{\mathbf{x}, \mathbf{y} \in D \\ \|\mathbf{x} - \mathbf{y}\| \leq \delta}} |g(\mathbf{x}) - g(\mathbf{y})|, \quad \delta > 0.$$

If we need/want to specify D , we will say that $\omega_g(\cdot)$ is the modulus of continuity of g over D .

- χ_E is the characteristic (or indicator) function of the set E ,

$$\chi_E(\xi) = \begin{cases} 1, & \text{if } \xi \in E, \\ 0, & \text{otherwise.} \end{cases}$$

- μ_k denotes the Lebesgue measure in \mathbb{R}^k . Throughout this book, unless otherwise stated, all the terminology coming from measure theory (such as ‘measurable set’, ‘measurable function’, ‘almost everywhere (a.e.)’, etc.) is always referred to the Lebesgue measure.
- If D is any measurable subset of some \mathbb{R}^k , we set

$$\begin{aligned} \mathfrak{M}_D &= \{f : D \rightarrow \mathbb{C} : f \text{ is measurable}\}, \\ L^p(D) &= \left\{f \in \mathfrak{M}_D : \int_D |f|^p < \infty\right\}, \quad 1 \leq p < \infty, \\ L^\infty(D) &= \{f \in \mathfrak{M}_D : \text{ess sup}_D |f| < \infty\}. \end{aligned}$$

If D is the special domain $[0, 1] \times [-\pi, \pi]$, we preferably use the notation \mathfrak{M} instead of \mathfrak{M}_D :

$$\mathfrak{M} = \{\kappa : [0, 1] \times [-\pi, \pi] \rightarrow \mathbb{C} : \kappa \text{ is measurable}\}.$$

If $f \in L^p(D)$ and the domain D is clear from the context, we write $\|f\|_{L^p}$ instead of $\|f\|_{L^p(D)}$ to indicate the L^p -norm of f , which is defined as

$$\|f\|_{L^p} = \begin{cases} (\int_D |f|^p)^{1/p}, & \text{if } 1 \leq p < \infty, \\ \text{ess sup}_D |f|, & \text{if } p = \infty. \end{cases}$$

For more on L^p spaces, see Sect. 2.2.2.

- We use a notation borrowed from probability theory to indicate sets. For example, if $f, g : D \subseteq \mathbb{R}^k \rightarrow \mathbb{C}$, then

$$\begin{aligned} \{f \neq 1\} &= \{\mathbf{x} \in D : f(\mathbf{x}) \neq 1\}, \\ \{f \in D(z, \varepsilon)\} &= \{\mathbf{x} \in D : f(\mathbf{x}) \in D(z, \varepsilon)\}, \\ \{0 \leq f \leq 1, g > 2\} &= \{\mathbf{x} \in D : 0 \leq f(\mathbf{x}) \leq 1, g(\mathbf{x}) > 2\}, \\ \mu_k\{f > 0, g < 0\} &\text{ is the measure of the set } \{\mathbf{x} \in D : f(\mathbf{x}) > 0, g(\mathbf{x}) < 0\}, \\ \chi_{\{f=0\}} &\text{ is the characteristic function of the set where } f \text{ vanishes,} \\ &\dots \end{aligned}$$

- A functional ϕ is any function defined on some vector space (such as, for example, $C_c(\mathbb{C})$ or $C_c(\mathbb{R})$) and taking values in \mathbb{C} .
- If \mathbb{K} is either \mathbb{R} or \mathbb{C} and $g : D \subset \mathbb{R}^k \rightarrow \mathbb{K}$ is a measurable function defined on a set D with $0 < \mu_k(D) < \infty$, we denote by ϕ_g the functional

$$\phi_g : C_c(\mathbb{K}) \rightarrow \mathbb{C}, \quad \phi_g(F) = \frac{1}{\mu_k(D)} \int_D F(g(\mathbf{x})) d\mathbf{x}.$$

- A matrix-sequence (or sequence of matrices) is any sequence of the form $\{A_n\}_n$, where $A_n \in \mathbb{C}^{n \times n}$ and n varies in some infinite subset of \mathbb{N} .
- We denote by \mathcal{E} the space of all matrix-sequences,

$$\mathcal{E} = \{\{A_n\}_n : \{A_n\}_n \text{ is a matrix-sequence}\}.$$

2.2 Preliminaries on Measure and Integration Theory

In this section we collect the necessary background material about measure and integration theory. Reference textbooks on this subject are, for example, [20, 95, 97]. We will anyway provide precise citations alongside each result we will not prove.

2.2.1 Essential Range

Given a measurable function $f : D \subset \mathbb{R}^k \rightarrow \mathbb{C}$, the essential range of f is denoted by $\mathcal{ER}(f)$ and is defined as the set of points $z \in \mathbb{C}$ such that, for every $\varepsilon > 0$, the measure of the set $\{f \in D(z, \varepsilon)\}$ is positive. In formulas,

$$\mathcal{ER}(f) = \{z \in \mathbb{C} : \mu_k\{f \in D(z, \varepsilon)\} > 0 \text{ for all } \varepsilon > 0\}.$$

Basic properties of the essential range are collected in the next lemma.

Lemma 2.1 *Let $f : D \subset \mathbb{R}^k \rightarrow \mathbb{C}$ be measurable. Then $\mathcal{ER}(f)$ is closed and $f \in \mathcal{ER}(f)$ a.e.*

Proof We show that the complement of $\mathcal{ER}(f)$ is open. If $z \in \mathbb{C} \setminus \mathcal{ER}(f)$ then $\mu_k\{f \in D(z, \varepsilon)\} = 0$ for some $\varepsilon > 0$. Each point $w \in D(z, \varepsilon)$ has a neighborhood $D(w, \delta)$ such that $D(w, \delta) \subseteq D(z, \varepsilon)$ and, consequently, $\mu_k\{f \in D(w, \delta)\} = 0$. We conclude that $D(z, \varepsilon) \subseteq \mathbb{C} \setminus \mathcal{ER}(f)$, hence $\mathbb{C} \setminus \mathcal{ER}(f)$ is open.

To prove that $f \in \mathcal{ER}(f)$ a.e., let

$$\mathcal{B} = \left\{ D\left(q, \frac{1}{m}\right) : q = a + ib, \ a, b \in \mathbb{Q}, \ m \in \mathbb{N} \right\}.$$

\mathcal{B} is a topological basis of \mathbb{C} , i.e., for each open set $U \subseteq \mathbb{C}$ and each $u \in U$, there exists an element of \mathcal{B} which contains u and is contained in U . Since $\mathbb{C} \setminus \mathcal{ER}(f)$ is open and every $z \in \mathbb{C} \setminus \mathcal{ER}(f)$ has a neighborhood $D(z, \varepsilon)$ such that $\mu_k\{f \in D(z, \varepsilon)\} = 0$ (by definition of $\mathcal{ER}(f)$), for each $z \in \mathbb{C} \setminus \mathcal{ER}(f)$ there exists an element of \mathcal{B} , say $D_z = D(q_z, \frac{1}{m_z})$, such that $z \in D_z \subseteq \mathbb{C} \setminus \mathcal{ER}(f)$ and $\mu_k\{f \in D_z\} = 0$. Let \mathcal{C} be the subset of \mathcal{B} given by $\mathcal{C} = \{D_z : z \in \mathbb{C} \setminus \mathcal{ER}(f)\}$. Since \mathcal{B} is countable, \mathcal{C} is countable as well, say $\mathcal{C} = \{C_\ell : \ell = 1, 2, \dots\}$, and we have

$$\begin{aligned} \mu_k\{f \notin \mathcal{ER}(f)\} &= \mu_k\left(\bigcup_{z \in \mathbb{C} \setminus \mathcal{ER}(f)} \{f = z\}\right) \leq \mu_k\left(\bigcup_{z \in \mathbb{C} \setminus \mathcal{ER}(f)} \{f \in D_z\}\right) \\ &= \mu_k\left(\bigcup_{\ell=1}^{\infty} \{f \in C_\ell\}\right) \leq \sum_{\ell=1}^{\infty} \mu_k\{f \in C_\ell\} = 0, \end{aligned}$$

which completes the proof. \square

If $f : D \subseteq \mathbb{R}^k \rightarrow \mathbb{C}$ is any measurable function which is real a.e., the essential infimum [supremum] of f is defined as the infimum [supremum] of the essential range of f . Note that this makes sense, because $f \in \mathbb{R}$ a.e. and, consequently, $\mathcal{ER}(f) \subseteq \mathbb{R}$. The essential infimum of f is denoted by $\text{ess inf}_D f$ or $\text{ess inf}_{\mathbf{x} \in D} f(\mathbf{x})$. Likewise, the essential supremum of f is denoted by $\text{ess sup}_D f$ or $\text{ess sup}_{\mathbf{x} \in D} f(\mathbf{x})$. By definition,

$$\text{ess inf}_D f = \inf \mathcal{ER}(f),$$

$$\text{ess sup}_D f = \sup \mathcal{ER}(f).$$

Note that an equivalent definition of $\text{ess inf}_D f$ and $\text{ess sup}_D f$ is the following:

$$\text{ess inf}_D f = \inf \{\alpha \in \mathbb{R} : \mu_k\{f > \alpha\} > 0\},$$

$$\text{ess sup}_D f = \sup \{\beta \in \mathbb{R} : \mu_k\{f < \beta\} > 0\}.$$

Exercise 2.1 Suppose $f : D \subseteq \mathbb{R}^k \rightarrow \mathbb{C}$ is continuous and D is contained in the closure of its interior. Prove that $\mathcal{ER}(f)$ coincides with the closure of the image of f , that is, $\mathcal{ER}(f) = \overline{f(D)}$, where $f(D) = \{f(\mathbf{x}) : \mathbf{x} \in D\}$.

2.2.2 L^p Spaces

Let $D \subseteq \mathbb{R}^k$ be any measurable set. Let \mathfrak{M}_D be the space of complex-valued measurable functions defined on D ,

$$\mathfrak{M}_D = \{f : D \rightarrow \mathbb{C} : f \text{ is measurable}\},$$

and consider the following spaces:

$$\begin{aligned} L^p(D) &= \left\{ f \in \mathfrak{M}_D : \int_D |f|^p < \infty \right\}, \quad 1 \leq p < \infty, \\ L^\infty(D) &= \{f \in \mathfrak{M}_D : \text{ess sup}_D |f| < \infty\}. \end{aligned}$$

If we identify two functions $f, g \in L^p(D)$ whenever $f = g$ a.e., and if we set

$$\|f\|_{L^p} = \begin{cases} (\int_D |f|^p)^{1/p}, & \text{if } 1 \leq p < \infty, \\ \text{ess sup}_D |f|, & \text{if } p = \infty, \end{cases}$$

then $\|\cdot\|_{L^p}$ is a norm on $L^p(D)$ for all $p \in [1, \infty]$, the so-called L^p -norm; see, e.g., [97, Chap. 3]. If the domain D is not clear from the context, we will write $\|\cdot\|_{L^p(D)}$ instead of $\|\cdot\|_{L^p}$. The fact that $\|\cdot\|_{L^1}$ is a norm on $L^1(D)$ immediately implies the following ‘vanishing property’:

$$\int_D |f| = 0 \implies f = 0 \text{ a.e.} \quad (2.1)$$

Given $1 \leq p, q \leq \infty$, we say that p, q are conjugate exponents if $\frac{1}{p} + \frac{1}{q} = 1$ (it is understood that $\frac{1}{\infty} = 0$). By Hölder’s inequality [97, Theorem 3.8], if $f \in L^p(D)$ and $g \in L^q(D)$, with $1 \leq p, q \leq \infty$ conjugate exponents, then $fg \in L^1(D)$ and

$$\|fg\|_{L^1} \leq \|f\|_{L^p} \|g\|_{L^q}. \quad (2.2)$$

As a consequence, if $1 \leq r < p \leq \infty$, $f \in L^p(D)$ and $\mu_k(D) < \infty$, then $f \in L^r(D)$. This is clear for $p = \infty$, while for $p < \infty$ it follows from Hölder’s inequality and the observation that $\frac{p}{r}, \frac{p}{p-r}$ are conjugate exponents, $|f|^r \in L^{p/r}(D)$ and $1 \in L^{p/(p-r)}(D)$ (because $\mu_k(D) < \infty$). For $p = q = 2$, Hölder’s inequality (2.2) is also known as the Cauchy–Schwarz inequality.

Another important inequality is Jensen’s inequality [97, Theorem 3.3]. In combination with [97, Theorem 1.29], Jensen’s inequality implies the following result, which will be used in this book: if $f \in L^p(D)$ with $1 \leq p < \infty$ and $g \in L^1(D)$ is such that $fg \in L^1(D)$, $g \geq 0$ and $\int_D g = 1$, then

$$\left(\int_D |f|g \right)^p \leq \int_D |f|^p g. \quad (2.3)$$

Let $C_c(D)$ be the space of continuous functions $f : D \rightarrow \mathbb{C}$ such that the support $\text{supp}(f) = \overline{\{f \neq 0\}}$ is compact. The space $C_c(D)$ is dense in $L^p(D)$ for all $1 \leq p < \infty$, so for each $f \in L^p(D)$ there is a sequence $\{f_m\}_m \subset C_c(D)$ such that $f_m \rightarrow f$ in $L^p(D)$, i.e., $\|f_m - f\|_{L^p} \rightarrow 0$. For a proof of this result, see [97, Theorem 3.14]. Another density result that will be of interest herein is stated in the next lemma.

Lemma 2.2 *Let $D = [a_1, b_1] \times \cdots \times [a_k, b_k]$ be a hyperrectangle in \mathbb{R}^k and let \mathcal{P}_D be the space generated by the trigonometric monomials*

$$\left\{ e^{i \left(\frac{2\pi}{b_1 - a_1} j_1 y_1 + \cdots + \frac{2\pi}{b_k - a_k} j_k y_k \right)} : (j_1, \dots, j_k) \in \mathbb{Z}^k \right\},$$

that is, the set of all finite linear combinations of such monomials (we call it the space of scaled k -variate trigonometric polynomials). Then \mathcal{P}_D is dense in $L^1(D)$, so for each $f \in L^1(D)$ there is a sequence $\{f_m\}_m \subset \mathcal{P}_D$ such that $\|f_m - f\|_{L^1} \rightarrow 0$.

Proof Let us first consider the univariate case $k = 1$. In this case, $D = [a, b]$ is an interval and $\mathcal{P}_{[a,b]}$ is the space of scaled trigonometric polynomials,

$$\mathcal{P}_{[a,b]} = \left\{ \sum_{j=-N}^N \alpha_j e^{i \frac{2\pi}{b-a} j y} : \alpha_{-N}, \dots, \alpha_N \in \mathbb{C}, N \in \mathbb{N} \right\}.$$

Let $f \in L^1([a, b])$ and $\varepsilon > 0$. Since the space of continuous functions $C([a, b])$ is dense in $L^1([a, b])$, there exists $f_\varepsilon \in C([a, b])$ such that

$$\|f - f_\varepsilon\|_{L^1} \leq \varepsilon.$$

Now, for any $g \in L^2([a, b])$, the Fourier series of g , namely

$$\sum_{j=-\infty}^{\infty} g_j e^{i \frac{2\pi}{b-a} j y}, \quad g_j = \frac{1}{b-a} \int_a^b g(y) e^{-i \frac{2\pi}{b-a} j y} dy,$$

converges to g in $L^2([a, b])$, i.e.,

$$\lim_{N \rightarrow \infty} \left\| \sum_{j=-N}^N g_j e^{i \frac{2\pi}{b-a} j y} - g(y) \right\|_{L^2} = 0;$$

see, e.g., [97, pp. 91–92]. It follows that the set of trigonometric polynomials $\mathcal{P}_{[a,b]}$ is dense in $L^2([a, b])$. In particular, there exists a trigonometric polynomial $p_\varepsilon \in \mathcal{P}_{[a,b]}$ such that

$$\|f_\varepsilon - p_\varepsilon\|_{L^2} \leq \varepsilon.$$

Thus, by Hölder's inequality,

$$\begin{aligned}
\|f - p_\varepsilon\|_{L^1} &\leq \|f - f_\varepsilon\|_{L^1} + \|f_\varepsilon - p_\varepsilon\|_{L^1} \\
&\leq \|f - f_\varepsilon\|_{L^1} + \|1\|_{L^2} \|f_\varepsilon - p_\varepsilon\|_{L^2} \\
&\leq \varepsilon + (b - a)\varepsilon.
\end{aligned}$$

This concludes the proof for the univariate case. The extension to the k -variate case is conceptually identical: we use the density of the space of continuous functions $C(D)$ in $L^1(D)$ with respect to the L^1 -norm, and the density of \mathcal{P}_D in $L^2(D)$ with respect to the L^2 -norm, which is a consequence of the L^2 -convergence of the k -variate Fourier series of any function in $L^2(D)$ to the function itself. \square

2.2.3 Convergence in Measure, a.e., in L^p

Let $f_m, f : D \subseteq \mathbb{R}^k \rightarrow \mathbb{C}$ be measurable functions. We say that $f_m \rightarrow f$ in measure if, for every $\varepsilon > 0$,

$$\lim_{m \rightarrow \infty} \mu_k\{|f_m - f| > \varepsilon\} = 0.$$

We say that $f_m \rightarrow f$ a.e. if

$$\mu_k\{f_m \not\rightarrow f\} = 0.$$

Important results about convergence in measure, a.e. and in L^p are reported in the next lemmas.

Lemma 2.3 *Let $f_m, g_m, f, g : D \subseteq \mathbb{R}^k \rightarrow \mathbb{C}$ be measurable functions.*

- *If $f_m \rightarrow f$ in measure, then $|f_m| \rightarrow |f|$ in measure.*
- *If $f_m \rightarrow f$ in measure and $g_m \rightarrow g$ in measure, then $\alpha f_m + \beta g_m \rightarrow \alpha f + \beta g$ in measure for all $\alpha, \beta \in \mathbb{C}$.*
- *If $f_m \rightarrow f$ in measure, $g_m \rightarrow g$ in measure, and $\mu_k(D) < \infty$, then $f_m g_m \rightarrow f g$ in measure.*

For the proof of Lemma 2.3, see [20, Corollary 2.2.6] or [60, Lemma 2.3].

Lemma 2.4 *Let $f_m, f : D \subseteq \mathbb{R}^k \rightarrow \mathbb{C}$ be measurable functions and assume that $\mu_k(D) < \infty$.*

- *If $f_m \rightarrow f$ a.e., then $f_m \rightarrow f$ in measure.*
- *If $f_m \rightarrow f$ in measure, then there is a subsequence $\{f_{m_i}\}_i$ such that $f_{m_i} \rightarrow f$ a.e.*
- *If $f_m, f \in L^p(D)$ for some $1 \leq p \leq \infty$ and $f_m \rightarrow f$ in $L^p(D)$, then $f_m \rightarrow f$ in measure.*

Lemma 2.4 is stated in [97, p. 74]. The proof of the first two statements can be found in [95, pp. 100–101] or [20, Theorems 2.2.3 and 2.2.5], while the third statement is a straightforward corollary of Chebyshev's inequality: for any $f : D \subseteq \mathbb{R}^k \rightarrow \mathbb{C}$ belonging to $L^p(D)$ and for any $\varepsilon > 0$,

$$\mu_k\{|f| > \varepsilon\} = \int_D \chi_{\{|f| > \varepsilon\}} \leq \int_D \frac{|f|^p}{\varepsilon^p} = \frac{\|f\|_{L^p}^p}{\varepsilon^p}. \quad (2.4)$$

As a consequence of Lemma 2.4, if $f_m, f \in L^p(D)$ with $\mu_k(D) < \infty$ and $f_m \rightarrow f$ in $L^p(D)$, then there is a subsequence $\{f_{m_i}\}_i$ such that $f_{m_i} \rightarrow f$ a.e.

A fundamental result of measure theory, which will be used several times also in this book, is Lebesgue's dominated convergence theorem [97, Theorem 1.34]. We report the corresponding statement for the reader's convenience.

Theorem 2.1 (dominated convergence theorem) *Let $f_m, f : D \subseteq \mathbb{R}^k \rightarrow \mathbb{C}$ be measurable functions. Suppose that $f_m \rightarrow f$ a.e. and that there exists $g \in L^1(D)$ such that $|f_m| \leq g$ over D for all m . Then $f_m, f \in L^1(D)$, $f_m \rightarrow f$ in $L^1(D)$ and*

$$\int_D f_m \rightarrow \int_D f.$$

Another important result of measure theory is Lusin's theorem [97, Theorem 2.24]. One of its consequences will be used in this book and is proved here.

Theorem 2.2 *Let $f : D \subseteq \mathbb{R}^k \rightarrow \mathbb{C}$ be a measurable function defined on a set D with $0 < \mu_k(D) < \infty$. Then, there exists a sequence of functions $\{f_m\}_m \subset C_c(D)$ such that $\sup_D |f_m| \leq \text{ess sup}_D |f|$ for all m and $f_m \rightarrow f$ a.e.*

Proof Let $\tilde{f} = f \chi_{\{|f| \leq \text{ess sup}_D |f|\}}$, so that $\tilde{f} = f$ a.e. by Lemma 2.1 and $|\tilde{f}| \leq \text{ess sup}_D |f|$ over D . By Lusin's theorem, for every $m \in \mathbb{N}$ there exists $f_m \in C_c(D)$ such that

$$\mu_k\{f_m \neq \tilde{f}\} < \frac{1}{m}, \quad \sup_D |f_m| \leq \sup_D |\tilde{f}| \leq \text{ess sup}_D |f|.$$

It is then clear that $f_m \rightarrow \tilde{f}$ in measure, so $f_m \rightarrow f$ in measure as well, because $\tilde{f} = f$ a.e. In view of Lemma 2.4, passing to a subsequence of $\{f_m\}_m$ (if necessary), we may assume that $f_m \rightarrow f$ a.e. \square

We conclude this section with a series of technical lemmas that we collect here in order to simplify the presentation of future chapters. Let \mathbb{K} be either \mathbb{R} or \mathbb{C} and let $g : D \subset \mathbb{R}^k \rightarrow \mathbb{K}$ be a measurable function defined on a set D with $0 < \mu_k(D) < \infty$. Consider the functional

$$\phi_g : C_c(\mathbb{K}) \rightarrow \mathbb{C}, \quad \phi_g(F) = \frac{1}{\mu_k(D)} \int_D F(g(\mathbf{x})) d\mathbf{x}. \quad (2.5)$$

ϕ_g is a continuous linear functional on the normed vector space $(C_c(\mathbb{K}), \|\cdot\|_\infty)$, and

$$\|\phi_g\| = \sup \frac{|\phi_g(F)|}{\|F\|_\infty} \leq 1,$$

where the supremum is taken over all $F \in C_c(\mathbb{K})$ which are not identically 0. Indeed, the linearity is obvious and the continuity, as well as the bound $\|\phi_g\| \leq 1$, follows

from the observation that $|\phi_g(F)| \leq \|F\|_\infty$ for all $F \in C_c(\mathbb{K})$. If g is constant, say $g = \gamma$ a.e., then $\phi_g = \phi_\gamma$ is the evaluation functional at γ , i.e., $\phi_\gamma(F) = F(\gamma)$ for every $F \in C_c(\mathbb{K})$.

Lemma 2.5 *Let \mathbb{K} be either \mathbb{R} or \mathbb{C} , and let $g_m, g : D \subset \mathbb{R}^k \rightarrow \mathbb{K}$ be measurable functions defined on a set D with $0 < \mu_k(D) < \infty$. If $g_m \rightarrow g$ in measure, then $F(g_m) \rightarrow F(g)$ in $L^1(D)$ for all $F \in C_c(\mathbb{K})$ and $\phi_{g_m} \rightarrow \phi_g$ pointwise over $C_c(\mathbb{K})$.*

Proof Assume that $g_m \rightarrow g$ in measure. We show that $F(g_m) \rightarrow F(g)$ in $L^1(D)$ for all $F \in C_c(\mathbb{K})$; this immediately implies that $\phi_{g_m} \rightarrow \phi_g$ pointwise over $C_c(\mathbb{K})$, because

$$|\phi_{g_m}(F) - \phi_g(F)| \leq \frac{1}{\mu_k(D)} \|F(g_m) - F(g)\|_{L^1}.$$

For every $F \in C_c(\mathbb{K})$, every m and every $\varepsilon > 0$,

$$\begin{aligned} \|F(g_m) - F(g)\|_{L^1} &= \int_D |F(g_m(\mathbf{x})) - F(g(\mathbf{x}))| d\mathbf{x} \\ &= \int_{\{|g_m - g| > \varepsilon\}} |F(g_m(\mathbf{x})) - F(g(\mathbf{x}))| d\mathbf{x} + \int_{\{|g_m - g| \leq \varepsilon\}} |F(g_m(\mathbf{x})) - F(g(\mathbf{x}))| d\mathbf{x} \\ &\leq 2\|F\|_\infty \mu_k\{|g_m - g| > \varepsilon\} + \omega_F(\varepsilon), \end{aligned} \quad (2.6)$$

where ω_F is the modulus of continuity of F , i.e.,

$$\omega_F(\varepsilon) = \sup_{\substack{y, z \in \mathbb{K} \\ |y - z| \leq \varepsilon}} |F(y) - F(z)|.$$

Since $g_m \rightarrow g$ in measure by assumption and F is uniformly continuous by the Heine–Cantor theorem [96, Theorem 4.19], we have

$$\lim_{m \rightarrow \infty} \mu_k\{|g_m - g| > \varepsilon\} = \lim_{\varepsilon \rightarrow 0} \omega_F(\varepsilon) = 0.$$

Therefore, passing first to the $\limsup_{m \rightarrow \infty}$ and then to the $\lim_{\varepsilon \rightarrow 0}$ in (2.6), we conclude that $F(g_m) \rightarrow F(g)$ in $L^1(D)$. \square

Lemma 2.5 admits the following converse.

Lemma 2.6 *Let \mathbb{K} be either \mathbb{R} or \mathbb{C} , and let $g_m, g : D \subset \mathbb{R}^k \rightarrow \mathbb{K}$ be measurable functions defined on a set D with $0 < \mu_k(D) < \infty$. If $\phi_{g_m - g} \rightarrow \phi_0$ pointwise over $C_c(\mathbb{K})$, then $g_m \rightarrow g$ in measure.*

Proof We first recall that ϕ_0 is the evaluation functional at 0. Hence, by hypothesis, for all $F \in C_c(\mathbb{K})$ we have

$$\lim_{m \rightarrow \infty} \frac{1}{\mu_k(D)} \int_D F(g_m(\mathbf{x}) - g(\mathbf{x})) d\mathbf{x} = F(0). \quad (2.7)$$

Suppose by contradiction that $g_m \not\rightarrow g$ in measure. Then, there exist $\varepsilon, \delta > 0$ and a subsequence $\{g_{m_i}\}_i$ such that, for all i ,

$$\mu_k\{|g_{m_i} - g| \geq \varepsilon\} \geq \delta.$$

Take a real function $F \in C_c(\mathbb{K})$ such that $F(0) = 1 = \max_{y \in \mathbb{K}} F(y)$ and $F(y) = 0$ over $\{y \in \mathbb{K} : |y| \geq \varepsilon\}$. By the previous inequality, for all i we have

$$\begin{aligned} \frac{1}{\mu_k(D)} \int_D F(g_{m_i}(\mathbf{x}) - g(\mathbf{x})) d\mathbf{x} &= \frac{1}{\mu_k(D)} \int_{\{|g_{m_i} - g| < \varepsilon\}} F(g_{m_i}(\mathbf{x}) - g(\mathbf{x})) d\mathbf{x} \\ &\leq \frac{\mu_k\{|g_{m_i} - g| < \varepsilon\}}{\mu_k(D)} \leq \frac{\mu_k(D) - \delta}{\mu_k(D)} < 1 = F(0), \end{aligned}$$

which is a contradiction to (2.7). \square

Remark 2.1 Let ϕ_g be defined as in (2.5) and assume that $\phi_g = \phi_0$; then $g = 0$ a.e. Indeed, if $\phi_g = \phi_0$, the constant sequence $\{\phi_g\}_m$ converges pointwise to ϕ_0 over $C_c(\mathbb{K})$. By Lemma 2.6, this implies that $g \rightarrow 0$ in measure, and hence $g \rightarrow 0$ a.e. by Lemma 2.4. Thus, $g = 0$ a.e.

We recall that a trigonometric polynomial is any finite linear combinations of the so-called Fourier frequencies $e^{ik\theta}$, $k \in \mathbb{Z}$. More explicitly, a trigonometric polynomial is a function of the form

$$p(\theta) = \sum_{k=-N}^N a_k e^{ik\theta}, \quad a_{-N}, \dots, a_N \in \mathbb{C}, \quad N \in \mathbb{N}.$$

Lemma 2.7 *Let $f \in L^1([-\pi, \pi])$. Then, there exists a sequence of trigonometric polynomials $\{p_m\}_m$ such that $\|p_m\|_\infty \leq \text{ess sup}_{[-\pi, \pi]} |f|$ for all m and $p_m \rightarrow f$ a.e. and in $L^1([-\pi, \pi])$.*

Proof It suffices to show that, for each $\varepsilon > 0$, there exists a trigonometric polynomial p_ε such that

$$\|p_\varepsilon\|_\infty \leq \text{ess sup}_{[-\pi, \pi]} |f|, \quad \|f - p_\varepsilon\|_{L^1} \leq \varepsilon.$$

Indeed, this shows the existence of a sequence of trigonometric polynomials $\{p_m\}_m$ such that $\|p_m\|_\infty \leq \text{ess sup}_{[-\pi, \pi]} |f|$ for all m and $p_m \rightarrow f$ in $L^1([-\pi, \pi])$; in view of Lemma 2.4, passing to a subsequence of $\{p_m\}_m$ (if necessary), we may assume that $p_m \rightarrow f$ a.e.

Let $\varepsilon > 0$. By Theorem 2.2, there exists $f_\varepsilon \in C_c((-\pi, \pi))$ such that

$$\|f_\varepsilon\|_\infty \leq \text{ess sup}_{[-\pi, \pi]} |f|, \quad \|f - f_\varepsilon\|_{L^1} < \varepsilon. \quad (2.8)$$

The function f_ε is continuous on $[-\pi, \pi]$ and 2π -periodic, in the sense that $f_\varepsilon(-\pi) = f_\varepsilon(\pi)$ (indeed, we have $f_\varepsilon(-\pi) = f_\varepsilon(\pi) = 0$). We can therefore follow the nice

construction in [97, pp. 89–91] to obtain a trigonometric polynomial p_ε such that

$$\|p_\varepsilon\|_\infty \leq \|f_\varepsilon\|_\infty, \quad \|f_\varepsilon - p_\varepsilon\|_\infty < \varepsilon. \quad (2.9)$$

By combining (2.8)–(2.9), we arrive at

$$\|p_\varepsilon\|_\infty \leq \operatorname{ess\,sup}_{[-\pi, \pi]} |f|, \quad \|f - p_\varepsilon\|_{L^1} \leq \varepsilon(1 + 2\pi),$$

which proves the thesis. \square

Lemma 2.8 *Let $\kappa : [0, 1] \times [-\pi, \pi] \rightarrow \mathbb{C}$ be a measurable function. Then, there exists a sequence $\{\kappa_m\}_m$ such that $\kappa_m : [0, 1] \times [-\pi, \pi] \rightarrow \mathbb{C}$ is a function of the form*

$$\kappa_m(x, \theta) = \sum_{j=-N_m}^{N_m} a_j^{(m)}(x) e^{ij\theta}, \quad a_j^{(m)} \in C^\infty([0, 1]), \quad N_m \in \mathbb{N}, \quad (2.10)$$

and $\kappa_m \rightarrow \kappa$ a.e.

Proof The function $\tilde{\kappa}_m = \kappa \chi_{\{|\kappa| \leq 1/m\}}$ belongs to $L^\infty([0, 1] \times [-\pi, \pi])$ and converges to κ in measure. Indeed, $\tilde{\kappa}_m \rightarrow \kappa$ pointwise over $[0, 1] \times [-\pi, \pi]$, and the pointwise (a.e.) convergence on a set of finite measure implies the convergence in measure (Lemma 2.4). By Lemma 2.2, the space generated by the trigonometric monomials

$$\{e^{2\pi i \ell x} e^{ij\theta} : \ell, j \in \mathbb{Z}\}$$

is dense in $L^1([0, 1] \times [-\pi, \pi])$, so we can choose a function κ_m belonging to this space such that $\|\kappa_m - \tilde{\kappa}_m\|_{L^1} \leq 1/m$. Note that κ_m is a function of the form (2.10). Moreover, for each $\varepsilon > 0$, using Chebyshev's inequality (2.4) we obtain

$$\begin{aligned} \mu_2\{|\kappa_m - \kappa| > \varepsilon\} &\leq \mu_k(\{|\kappa_m - \tilde{\kappa}_m| > \varepsilon/2\} \cup \{|\tilde{\kappa}_m - \kappa| > \varepsilon/2\}) \\ &\leq \mu_2\{|\kappa_m - \tilde{\kappa}_m| > \varepsilon/2\} + \mu_2\{|\tilde{\kappa}_m - \kappa| > \varepsilon/2\} \\ &\leq \frac{\|\kappa_m - \tilde{\kappa}_m\|_{L^1}}{(\varepsilon/2)} + \mu_2\{|\tilde{\kappa}_m - \kappa| > \varepsilon/2\}, \end{aligned}$$

which converges to 0 as $m \rightarrow \infty$. Hence, $\kappa_m \rightarrow \kappa$ in measure. Since the convergence in measure on a set of finite measure implies the existence of a subsequence that converges a.e. (Lemma 2.4), passing to a subsequence of $\{\kappa_m\}_m$ (if necessary) we may assume that $\kappa_m \rightarrow \kappa$ a.e. \square

2.2.4 Riemann-Integrable Functions

A function $a : [0, 1] \rightarrow \mathbb{C}$ is said to be Riemann-integrable if its real and imaginary parts $\Re(a), \Im(a) : [0, 1] \rightarrow \mathbb{R}$ are Riemann-integrable in the classical sense. Recall that any Riemann-integrable function is *bounded* by definition. We report below a list of properties possessed by Riemann-integrable functions that will be used in this book, either explicitly or implicitly.

- If $\alpha, \beta \in \mathbb{C}$ and $a, b : [0, 1] \rightarrow \mathbb{C}$ are Riemann-integrable, then $\alpha a + \beta b$ is Riemann-integrable.
- If $a, b : [0, 1] \rightarrow \mathbb{C}$ are Riemann-integrable, then ab is Riemann-integrable.
- If $a : [0, 1] \rightarrow \mathbb{C}$ is Riemann-integrable and $F : \mathbb{C} \rightarrow \mathbb{C}$ is continuous, then $F(a) : [0, 1] \rightarrow \mathbb{C}$ is Riemann-integrable.
- If $a : [0, 1] \rightarrow \mathbb{C}$ is Riemann-integrable, then a belongs to $L^\infty([0, 1])$ and its Lebesgue and Riemann integrals over $[0, 1]$ coincide.
- If $a : [0, 1] \rightarrow \mathbb{C}$ is bounded, then a is Riemann-integrable if and only if a is continuous a.e.

Note that the last two properties imply the first three. The proof of the second-to-last property can be found in [95, pp. 73–74] or [20, Theorem 2.10.1]. The last property is Lebesgue's characterization theorem of Riemann-integrable functions [95, p. 104]. A further property of Riemann-integrable functions that will be used in this book is stated and proved in the next lemma.

Lemma 2.9 *Let $a : [0, 1] \rightarrow \mathbb{R}$ be Riemann-integrable. For each $n \in \mathbb{N}$, consider the partition of $(0, 1]$ given by the intervals*

$$I_{i,n} = \left(\frac{i-1}{n}, \frac{i}{n} \right], \quad i = 1, \dots, n,$$

and let

$$a_{i,n} \in \left[\inf_{x \in I_{i,n}} a(x), \sup_{x \in I_{i,n}} a(x) \right], \quad i = 1, \dots, n.$$

Then

$$\sum_{i=1}^n a_{i,n} \chi_{I_{i,n}} \rightarrow a \text{ a.e. in } [0, 1] \quad (2.11)$$

and

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n a_{i,n} = \int_0^1 a(x) dx. \quad (2.12)$$

Proof Fix $\varepsilon > 0$ and let $x \in (0, 1]$ be a continuity point of a . Then there is a $\delta > 0$ such that $|a(y) - a(x)| \leq \varepsilon$ whenever $y \in [0, 1]$ and $|y - x| \leq \delta$. Take $n \geq 1/\delta$ and call $I_{k,n}$ the unique interval of the partition $(0, 1] = \bigcup_{i=1}^n I_{i,n}$ containing x . For $y \in I_{k,n}$, we have $y \in [0, 1]$ and $|y - x| \leq \delta$, hence $|a(y) - a(x)| \leq \varepsilon$. It follows that

$$\begin{aligned} \left| \sum_{i=1}^n a_{i,n} \chi_{I_{i,n}}(x) - a(x) \right| &= |a_{k,n} - a(x)| \\ &\leq \max \left(a(x) - \inf_{y \in I_{k,n}} a(y), \sup_{y \in I_{k,n}} a(y) - a(x) \right) \leq \varepsilon. \end{aligned}$$

As a consequence, $\sum_{i=1}^n a_{i,n} \chi_{I_{i,n}}(x) \rightarrow a(x)$ whenever x is a continuity point of a in $(0, 1]$. This implies (2.11), because a is Riemann-integrable and hence continuous a.e. in $[0, 1]$. Since

$$\left| \sum_{i=1}^n a_{i,n} \chi_{I_{i,n}} \right| \leq \|a\|_{\infty} < \infty, \quad \frac{1}{n} \sum_{i=1}^n a_{i,n} = \int_0^1 \left(\sum_{i=1}^n a_{i,n} \chi_{I_{i,n}} \right),$$

Equation (2.12) follows from (2.11) and from the dominated convergence theorem. \square

2.3 Preliminaries on General Topology

This section covers specific topics from general topology that will be of interest in this book and may not be included in standard university courses. The reader is supposed to be familiar with basic topology, in particular with the notions of topological space and metric space. Any university course covers these topics, as well as any good book on general topology; see, e.g., the classic book by Kelley [80].

2.3.1 Pseudometric Spaces

A pseudometric on a set X is a function $d : X \times X \rightarrow [0, \infty)$ such that, for all points $x, y, z \in X$,

- (i) $x = y \implies d(x, y) = 0$,
- (ii) $d(x, y) = d(y, x)$,
- (iii) $d(x, y) \leq d(x, z) + d(z, y)$.

A pseudometric space is a pair (X, d) where d is a pseudometric on X . A pseudometric is often referred to as a distance, and $d(x, y)$ is called the distance between x and y . The difference between a pseudometric and a metric space is simply that in a metric space the property (i) is replaced by the stronger version $x = y \iff d(x, y) = 0$. In other words, the distance between two points in a metric space is zero if and only if the two points coincide, whereas in a pseudometric space the distance can be zero even if the two points do not coincide. However, this is not so disturbing as it is not so hard to accept that the distance between different points can be zero. For example, if the purpose of a distance is to quantify the money that is necessary to go from a place to another, it may certainly happen that the distance between two distinct places is

zero (e.g., because the two places are so close to each other that one can cover the spatial distance between them by feet, without spending money). What is interesting is that, also on the mathematical level, the theory of pseudometric spaces does not differ significantly from the theory of metric spaces. Kelley, for instance, develops simultaneously the theory of metric and pseudometric spaces in Chap. 4 of his book [80]. Moreover, if we imagine to identify two points in a pseudometric space whenever the distance between them is zero, we are introducing on the pseudometric space an equivalence relation with respect to which the quotient space is a metric space; see Exercise 2.3. In the remainder of this section, we investigate some properties of pseudometric spaces that we shall use in this book. The reader will easily recognize that all these properties have an exact analog in the world of metric spaces.

Let (X, d) be a pseudometric space. If $x \in X$, the open disk of radius $r > 0$ around x is defined as

$$D(x, r) = \{y \in X : d(x, y) < r\}.$$

We say that $U \subseteq X$ is an open set if for every point $x \in U$ there exists $r > 0$ such that $D(x, r) \subseteq U$. It is easy to check that the collection of all open sets U is a topology on X , which is referred to as the pseudometric topology induced by d . This topology will be denoted by τ_d . We say that a sequence $\{x_m\}_m \subseteq X$ converges to a point $x \in X$ in (X, d) if $d(x_m, x) \rightarrow 0$ as $m \rightarrow \infty$. This is actually equivalent to saying that $x_m \rightarrow x$ in the topological space (X, τ_d) .

Lemma 2.10 *Assume that the pseudometric spaces (X, d) and (X, d') have the same convergent sequences, i.e., $x_m \rightarrow x$ in (X, d) if and only if $x_m \rightarrow x$ in (X, d') . Then $\tau_d = \tau_{d'}$.*

Proof To avoid confusion, in this proof we denote by $D_d(y, r)$ (resp., $D_{d'}(y, r)$) the disk of radius r around y in the space (X, d) (resp., (X, d')). Suppose by contradiction that $\tau_d \neq \tau_{d'}$, and assume for example that $U \in \tau_d$ and $U \notin \tau_{d'}$. Since $U \notin \tau_{d'}$, there exists a point $x \in U$ such that $D_{d'}(x, \varepsilon)$ is not contained in U for all $\varepsilon > 0$. For each $m \in \mathbb{N}$, choose a point $x_m \in D_{d'}(x, \frac{1}{m})$ such that $x_m \notin U$. It is clear that $x_m \rightarrow x$ in (X, d') as $d'(x_m, x) < \frac{1}{m} \rightarrow 0$ as $m \rightarrow \infty$. Nevertheless, x_m cannot converge to x in (X, d) because $x_m \notin U$ for all m and U is open in (X, d) , so there exists $\delta > 0$ such that $D_d(x, \delta) \subseteq U$ and $x_m \notin D_d(x, \delta)$ for all m . This is a contradiction to the assumption that (X, d) and (X, d') have the same convergent sequences. \square

A topology τ on X is said to be pseudometrizable if there exists a pseudometric d on X such that $\tau = \tau_d$. If τ is pseudometrizable, there are actually infinite pseudometrics d on X such that $\tau = \tau_d$. For example, by Lemma 2.10, if we fix a pseudometric d such that $\tau = \tau_d$ and if we define d' to be any other pseudometric on X for which there exist two positive constants $\alpha, \beta > 0$ such that

$$\alpha d(x, y) \leq d'(x, y) \leq \beta d(x, y), \quad \forall x, y \in X,$$

then $\tau_{d'} = \tau_d$. Any two pseudometrics d, d' on X such that $\tau_d = \tau_{d'}$ are said to be (topologically) equivalent.

Exercise 2.2 Let (X, τ_X) and (Y, τ_Y) be topological spaces and consider $X \times Y$ equipped with the product topology $\tau_X \times \tau_Y$. We recall that $\tau_X \times \tau_Y$ is the collection of all sets $U \subseteq X \times Y$ with the following property: for every $(x, y) \in U$ there exist $U_X \in \tau_X$ and $U_Y \in \tau_Y$ such that $(x, y) \in U_X \times U_Y \subseteq U$. Show that if (X, τ_X) , (Y, τ_Y) are pseudometrizable and d_X, d_Y are pseudometrics inducing τ_X, τ_Y , respectively, then $(X \times Y, \tau_X \times \tau_Y)$ is pseudometrizable and $\tau_X \times \tau_Y$ is induced by the pseudometric

$$d_{X \times Y}((x, y), (x', y')) = \max(d_X(x, x'), d_Y(y, y')).$$

In addition, show that $\tau_X \times \tau_Y$ is also induced by any of the pseudometrics

$$d_{X \times Y}^{(p)}((x, y), (x', y')) = |(d_X(x, x'), d_Y(y, y'))|_p, \quad 1 \leq p \leq \infty.$$

Exercise 2.3 Given a pseudometric space (X, d) , we identify two points $x, y \in X$ whenever $d(x, y) = 0$. In other words, we introduce in X the equivalence relation

$$x \sim y \iff d(x, y) = 0.$$

The equivalence class of x is the set of all points y such that $d(x, y) = 0$; we call it the zone of x and we denote it by \tilde{x} . Let \tilde{X} be the quotient space (the set of all zones) and define d on \tilde{X} as follows:

$$d(\tilde{x}, \tilde{y}) = d(x, y), \quad x, y \in X.$$

Show that (\tilde{X}, d) is a metric space.

2.3.2 The Topology τ_{measure} of Convergence in Measure

Let $D \subset \mathbb{R}^k$ be a measurable set with $0 < \mu_k(D) < \infty$, and let

$$\mathfrak{M}_D = \{f : D \rightarrow \mathbb{C} : f \text{ is measurable}\}. \quad (2.13)$$

We already introduced in Sect. 2.2.3 the notion of convergence in measure on the space \mathfrak{M}_D . In this section, we show that this convergence is related to a pseudometric topology τ_{measure} on \mathfrak{M}_D . We also identify a specific pseudometric d_{measure} inducing τ_{measure} , and we study some of its properties that will be of interest in Chap. 5.

For every $f, g \in \mathfrak{M}_D$, let

$$d_{\text{measure}}(f, g) = p_{\text{measure}}(f - g),$$

$$p_{\text{measure}}(f) = \inf \left\{ \frac{\mu_k\{f_R \neq 0\}}{\mu_k(D)} + \|f_N\|_{L^\infty} : f_R, f_N \in \mathfrak{M}_D, f_R + f_N = f \right\}, \quad (2.14)$$

where it is understood that $\|f_N\|_{L^\infty} = \text{ess sup}_D |f_N| = \infty$ whenever $f \notin L^\infty(D)$.

Theorem 2.3 *The function p_{measure} satisfies the following properties.*

- (i) $0 \leq p_{\text{measure}}(f) \leq 1$ for all $f \in \mathfrak{M}_D$.
- (ii) $p_{\text{measure}}(f) = 0$ if and only if $f = 0$ a.e.
- (iii) $p_{\text{measure}}(f + g) \leq p_{\text{measure}}(f) + p_{\text{measure}}(g)$ for all $f, g \in \mathfrak{M}_D$.
- (iv) For all $f \in \mathfrak{M}_D$,

$$p_{\text{measure}}(f) = \inf \left\{ \frac{\mu_k(E^c \cap \{f \neq 0\})}{\mu_k(D)} + \text{ess sup}_E |f| : E \subseteq D \text{ measurable} \right\} \quad (2.15)$$

$$= \inf \left\{ \frac{\mu_k(E^c)}{\mu_k(D)} + \text{ess sup}_E |f| : E \subseteq D \text{ measurable} \right\}. \quad (2.16)$$

In particular, the function d_{measure} is a pseudometric on \mathfrak{M}_D and $d_{\text{measure}}(f, g) = 0$ if and only if $f = g$ a.e.

Proof (i) It is clear that $p_{\text{measure}}(f) \geq 0$. Moreover, by taking $f_R = f$ and $f_N = 0$ in (2.14), we see that $p_{\text{measure}}(f) \leq 1$.

(ii) If $f = 0$ a.e. then $p_{\text{measure}}(f) = 0$. Conversely, suppose that $p_{\text{measure}}(f) = 0$. Then, for every $\varepsilon > 0$ there exist $f_{R,\varepsilon}, f_{N,\varepsilon} \in \mathfrak{M}_D$ such that $f_{R,\varepsilon} + f_{N,\varepsilon} = f$ and

$$\mu_k\{f_{R,\varepsilon} \neq 0\} < \varepsilon, \quad \|f_{N,\varepsilon}\|_{L^\infty} < \varepsilon.$$

Therefore,

$$\mu_k\{|f| \geq \varepsilon\} \leq \mu_k(\{f_{R,\varepsilon} \neq 0\} \cup \{|f_{N,\varepsilon}| \geq \varepsilon\}) = \mu_k\{f_{R,\varepsilon} \neq 0\} < \varepsilon.$$

It follows that $\mu_k\{|f| \geq \varepsilon\} = 0$ for each $\varepsilon > 0$. Since $\{f \neq 0\} = \bigcup_{m=1}^{\infty} \{|f| \geq \frac{1}{m}\}$ is the union of countably many sets of zero measure, we conclude that $\mu_k\{f \neq 0\} = 0$ and $f = 0$ a.e.

(iii) By definition of $p_{\text{measure}}(f)$ and $p_{\text{measure}}(g)$, for every $\varepsilon > 0$ there exist four functions $f_{R,\varepsilon}, f_{N,\varepsilon}, g_{R,\varepsilon}, g_{N,\varepsilon} \in \mathfrak{M}_D$ such that $f_{R,\varepsilon} + f_{N,\varepsilon} = f$, $g_{R,\varepsilon} + g_{N,\varepsilon} = g$ and

$$\frac{\mu_k\{f_{R,\varepsilon} \neq 0\}}{\mu_k(D)} + \|f_{N,\varepsilon}\|_{L^\infty} < p_{\text{measure}}(f) + \varepsilon,$$

$$\frac{\mu_k\{g_{R,\varepsilon} \neq 0\}}{\mu_k(D)} + \|g_{N,\varepsilon}\|_{L^\infty} < p_{\text{measure}}(g) + \varepsilon.$$

Hence, there exist $(f + g)_{R,\varepsilon} = f_{R,\varepsilon} + g_{R,\varepsilon}$ and $(f + g)_{N,\varepsilon} = f_{N,\varepsilon} + g_{N,\varepsilon}$ such that $(f + g)_{R,\varepsilon} + (f + g)_{N,\varepsilon} = f + g$ and

$$\begin{aligned}
p_{\text{measure}}(f + g) &\leq \frac{\mu_k\{(f + g)_{R,\varepsilon} \neq 0\}}{\mu_k(D)} + \|(f + g)_{N,\varepsilon}\|_{L^\infty} \\
&\leq \frac{\mu_k\{f_{R,\varepsilon} \neq 0\} + \mu_k\{g_{R,\varepsilon} \neq 0\}}{\mu_k(D)} + \|f_{N,\varepsilon}\|_{L^\infty} + \|g_{N,\varepsilon}\|_{L^\infty} \\
&< p_{\text{measure}}(f) + p_{\text{measure}}(g) + 2\varepsilon.
\end{aligned}$$

It follows that $p_{\text{measure}}(f + g) \leq p_{\text{measure}}(f) + p_{\text{measure}}(g)$.

(iv) The equality (2.16) follows from the observation that, on the one hand, the left-hand side is clearly less than or equal to the right-hand side as $E^c \cap \{f \neq 0\} \subseteq E^c$; and, on the other hand, for each measurable $E \subseteq D$ we have

$$\frac{\mu_k(E^c \cap \{f = 0\})}{\mu_k(D)} + \text{ess sup}_E |f| = \frac{\mu_k(\hat{E}^c)}{\mu_k(D)} + \text{ess sup}_{\hat{E}} |f|,$$

where $\hat{E} = E \cup \{f = 0\}$. We prove the equality (2.15). Let $q(f)$ be the right-hand side of (2.15). For each measurable $E \subseteq D$, setting $f_R = f \chi_{E^c}$ and $f_N = f \chi_E$, we have $f_R + f_N = f$ and

$$\frac{\mu_k\{f_R \neq 0\}}{\mu_k(D)} + \|f_N\|_{L^\infty} = \frac{\mu_k(E^c \cap \{f \neq 0\})}{\mu_k(D)} + \text{ess sup}_E |f|.$$

Hence $q(f) \geq p_{\text{measure}}(f)$. Conversely, for each $f_R, f_N \in \mathfrak{M}_D$ such that $f_R + f_N = f$, setting $E = \{f_R = 0\}$ and noting that $f_N = f$ on E , we have

$$\frac{\mu_k\{f_R \neq 0\}}{\mu_k(D)} + \|f_N\|_{L^\infty} = \frac{\mu_k(E^c)}{\mu_k(D)} + \|f_N\|_{L^\infty} \geq \frac{\mu_k(E^c \cap \{f \neq 0\})}{\mu_k(D)} + \text{ess sup}_E |f|.$$

Hence $p_{\text{measure}}(f) \geq q(f)$. □

Since d_{measure} is a pseudometric on \mathfrak{M}_D , it induces on \mathfrak{M}_D a topology, which we denote by τ_{measure} . The next theorem shows that the notion of convergence associated with τ_{measure} is precisely the notion of convergence in measure; that is, a sequence $\{f_m\}_m \subseteq \mathfrak{M}_D$ converges in measure to $f \in \mathfrak{M}_D$ if and only if $d_{\text{measure}}(f_m, f) \rightarrow 0$.

Theorem 2.4 *Let $f_m, f : D \rightarrow \mathbb{C}$ be measurable functions. Then, the following conditions are equivalent.*

- (i) $f_m \rightarrow f$ in measure.
- (ii) $p_{\text{measure}}(f_m - f) \rightarrow 0$ as $m \rightarrow \infty$.

Proof (i) \implies (ii). Assume that $f_m \rightarrow f$ in measure and fix $\varepsilon > 0$. Then, there exists $M(\varepsilon)$ such that, for $m \geq M(\varepsilon)$,

$$\mu_k\{|f_m - f| \geq \varepsilon\} < \varepsilon.$$

Setting $E_{m,\varepsilon} = \{|f_m - f| < \varepsilon\}$ and using Theorem 2.3(iv), for $m \geq M(\varepsilon)$ we obtain

$$p_{\text{measure}}(f_m - f) \leq \frac{\mu_k(E_{m,\varepsilon}^c)}{\mu_k(D)} + \text{ess sup}_{E_{m,\varepsilon}} |f_m - f| < \frac{\varepsilon}{\mu_k(D)} + \varepsilon.$$

It follows that $p_{\text{measure}}(f_m - f) \rightarrow 0$.

(ii) \implies (i). Assume that $p_{\text{measure}}(f_m - f) \rightarrow 0$ and fix $\varepsilon > 0$. Then, there exists $M(\varepsilon)$ such that, for $m \geq M(\varepsilon)$,

$$p_{\text{measure}}(f_m - f) < \varepsilon.$$

By Theorem 2.3(iv), for $m \geq M(\varepsilon)$ there is a measurable $E_{m,\varepsilon} \subseteq D$ such that

$$\frac{\mu_k(E_{m,\varepsilon}^c)}{\mu_k(D)} + \text{ess sup}_{E_{m,\varepsilon}} |f_m - f| < \varepsilon \implies \frac{\mu_k(E_{m,\varepsilon}^c)}{\mu_k(D)} < \varepsilon, \quad \text{ess sup}_{E_{m,\varepsilon}} |f_m - f| < \varepsilon.$$

Thus,

$$\mu_k\{|f_m - f| \geq \varepsilon\} \leq \mu_k(E_{m,\varepsilon}^c) < \varepsilon \mu_k(D).$$

We conclude that $f_m \rightarrow f$ in measure. \square

Remark 2.2 (equivalent pseudometrics inducing τ_{measure}) Other pseudometrics on \mathfrak{M}_D , which are topologically equivalent to d_{measure} and induce on \mathfrak{M}_D the topology τ_{measure} of convergence in measure, are the following (see Exercise 2.4):

$$\begin{aligned} d'_{\text{measure}}(f, g) &= p'_{\text{measure}}(f - g), \\ p'_{\text{measure}}(f) &= \int_D \min(|f|, 1), \end{aligned} \tag{2.17}$$

and

$$\begin{aligned} d''_{\text{measure}}(f, g) &= p''_{\text{measure}}(f - g), \\ p''_{\text{measure}}(f) &= \int_D \frac{|f|}{1 + |f|}. \end{aligned} \tag{2.18}$$

The pseudometrics (2.17)–(2.18) are actually much more common than our nonstandard pseudometric d_{measure} as they are usually proposed in standard textbooks; see, e.g., [95, p. 102] and [20, p. 306]. Nevertheless, our pseudometric is illuminating for the purposes of Chap. 5, and this is why we preferred it to the others.

In view of Exercise 2.4, we recall that a function $\varphi : I \rightarrow \mathbb{R}$, defined on some interval $I \subseteq \mathbb{R}$, is said to be concave if

$$\varphi(x + \lambda(y - x)) \geq \varphi(x) + \lambda(\varphi(y) - \varphi(x)) \tag{2.19}$$

for every $x, y \in I$ such that $x < y$ and every $\lambda \in (0, 1)$. Graphically, the condition is expressed as follows: if $x < u < y$, then the point $(u, \varphi(u))$ lies above or on the line connecting the points $(x, \varphi(x))$ and $(y, \varphi(y))$. In other words, φ is concave if and only if

$$\varphi(u) \geq \varphi(x) + \frac{u-x}{y-x}(\varphi(y) - \varphi(x)) \iff \frac{\varphi(y) - \varphi(x)}{y-x} \leq \frac{\varphi(u) - \varphi(x)}{u-x} \quad (2.20)$$

for every $x, u, y \in I$ such that $x < u < y$; see also Fig. 2.1, in which

$$\xi_{x,y}(u) = \varphi(x) + \frac{u-x}{y-x}(\varphi(y) - \varphi(x)). \quad (2.21)$$

A simple characterization of concave functions is the following: $\varphi : I \rightarrow \mathbb{R}$ is concave if and only if the incremental ratios of φ satisfy the property

$$\frac{\varphi(y) - \varphi(u)}{y-u} \leq \frac{\varphi(u) - \varphi(x)}{u-x} \quad (2.22)$$

for every $x, u, y \in I$ such that $x < u < y$. Indeed, defining $\xi_{x,y}(u)$ as in (2.21) and taking into account that the points $(x, \varphi(x))$, $(u, \xi_{x,y}(u))$, $(y, \varphi(y))$ lie on the same line (see Fig. 2.1), we can prove that (2.22) is equivalent to (2.20) as follows:

$$\begin{aligned} \frac{\varphi(y) - \varphi(u)}{y-u} &\leq \frac{\varphi(u) - \varphi(x)}{u-x} \\ \iff \varphi(u) &\geq \varphi(x) + \frac{u-x}{y-u}(\varphi(y) - \varphi(u)) \\ \iff \varphi(u) &\geq \xi_{x,y}(u) + \frac{u-x}{y-u}(\varphi(y) - \varphi(u)) - \frac{u-x}{y-x}(\varphi(y) - \varphi(x)) \\ \iff \varphi(u) &\geq \xi_{x,y}(u) + \frac{u-x}{y-u}(\varphi(y) - \varphi(u)) - \frac{u-x}{y-u}(\varphi(y) - \xi_{x,y}(u)) \\ \iff \varphi(u) &\geq \xi_{x,y}(u). \end{aligned}$$

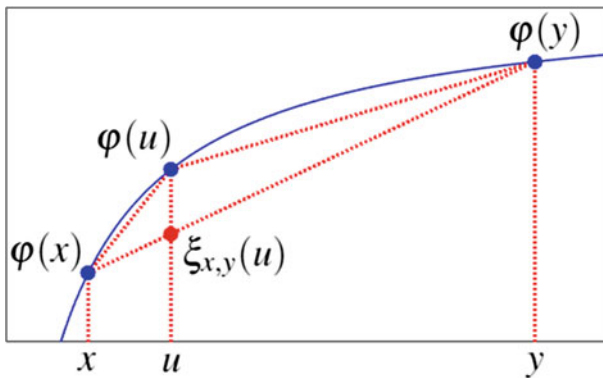


Fig. 2.1 Illustration of the properties of a concave function

It can be shown that a concave function $\varphi : I \rightarrow \mathbb{R}$ is continuous on the interior of I ; see [97, Theorem 3.2] and take into account that a function φ is concave if and only if $-\varphi$ is convex according to [97, Definition 3.1].

Exercise 2.4 Let $\varphi : [0, \infty) \rightarrow \mathbb{R}$ be a concave function such that $\varphi(0) = 0$.

(a) Show that φ is subadditive, i.e.,

$$\varphi(x + y) \leq \varphi(x) + \varphi(y), \quad \forall x, y \in [0, \infty).$$

(b) Suppose that $\varphi(x)$ does not diverge to $-\infty$ as $x \rightarrow \infty$. Show that φ is non-decreasing (and hence nonnegative due to the condition $\varphi(0) = 0$).

(c) Suppose that φ is bounded, continuous at 0 (hence continuous on the whole $[0, \infty)$ by the result mentioned above) and positive on $(0, \infty)$ (hence non-decreasing on the whole $[0, \infty)$ by (b)). Let \mathfrak{M}_D be as in (2.13) and for every $f, g \in \mathfrak{M}_D$ let

$$\begin{aligned} d_{\text{measure}}^\varphi(f, g) &= p_{\text{measure}}^\varphi(f - g), \\ p_{\text{measure}}^\varphi(f) &= \int_D \varphi(|f|). \end{aligned} \tag{2.23}$$

Show that $d_{\text{measure}}^\varphi$ is a pseudometric on \mathfrak{M}_D inducing the topology τ_{measure} of convergence in measure.

2.4 Preliminaries on Matrix Analysis

In this section we collect the necessary background material about matrix analysis. Reference textbooks on this subject are, for example, [13, 17, 69]. The reader, however, is not required to know everything about these books, which contain much more than is needed here. As in Sect. 2.2, we will provide precise citations next to each result we will state without a proof.

2.4.1 p -norms

Given $1 \leq p \leq \infty$ and a vector $\mathbf{x} \in \mathbb{C}^m$, we denote by $|\mathbf{x}|_p$ the p -norm of \mathbf{x} , i.e.,

$$|\mathbf{x}|_p = \begin{cases} (\sum_{i=1}^m |x_i|^p)^{1/p}, & \text{if } 1 \leq p < \infty, \\ \max_{i=1, \dots, m} |x_i|, & \text{if } p = \infty. \end{cases}$$

The p -norm of a matrix $X \in \mathbb{C}^{m \times m}$ is simply the operator norm of X regarded as an application from $(\mathbb{C}^m, |\cdot|_p)$ into itself. In formulas,

$$|X|_p = \max_{\substack{\mathbf{x} \in \mathbb{C}^m \\ \mathbf{x} \neq \mathbf{0}}} \frac{|X\mathbf{x}|_p}{|\mathbf{x}|_p} = \max_{\substack{\mathbf{x} \in \mathbb{C}^m \\ |\mathbf{x}|_p=1}} |X\mathbf{x}|_p.$$

The 2-norm of both vectors and matrices is also known as the spectral (or Euclidean) norm and will be preferably denoted by $\|\cdot\|$. Just like any other operator norm, the p -norm $|X|_p$ satisfies the following inequalities:

$$\rho(X) \leq |X|_p \quad (2.24)$$

and

$$|X\mathbf{x}|_p \leq |X|_p |\mathbf{x}|_p \quad (2.25)$$

for all $\mathbf{x} \in \mathbb{C}^m$. Moreover, $|X|_p$ is the smallest constant for which (2.25) holds for all $\mathbf{x} \in \mathbb{C}^m$. Therefore, observing that

$$|XY\mathbf{x}|_p \leq |X|_p |Y\mathbf{x}|_p \leq |X|_p |Y|_p |\mathbf{x}|_p,$$

one immediately obtains the submultiplicative property:

$$|XY|_p \leq |X|_p |Y|_p, \quad X, Y \in \mathbb{C}^{m \times m}. \quad (2.26)$$

The most important among the p -norms are undoubtedly the 1-norm, the 2-norm and the ∞ -norm. For each of these norms, a special formula for the computation of $|X|_p$ is available:

$$|X|_1 = \max_{j=1,\dots,m} \sum_{i=1}^m |x_{ij}|, \quad (2.27)$$

$$\|X\| = \sqrt{\rho(X^*X)} = \sqrt{\lambda_{\max}(X^*X)}, \quad (2.28)$$

$$|X|_\infty = \max_{i=1,\dots,m} \sum_{j=1}^m |x_{ij}|; \quad (2.29)$$

see, e.g., [17, Theorem 3.9] or [69, pp. 72–73]. Formulas (2.27) and (2.29) show that $|X|_1$ and $|X|_\infty$ are particularly easy to compute as they admit explicit expressions in terms of the components of X . More precisely, $|X|_1$ is the maximum among the 1-norms of the columns of X , and $|X|_\infty$ is the maximum among the 1-norms of the rows of X . Formula (2.28) shows that the 2-norm is unitarily invariant, that is,

$$\|X\| = \|UXV\| \quad (2.30)$$

for all $X \in \mathbb{C}^{m \times m}$ and all unitary matrices $U, V \in \mathbb{C}^{m \times m}$. Indeed,

$$\begin{aligned}\|UXV\| &= \sqrt{\rho((UXV)^*(UXV))} = \sqrt{\rho(V^*X^*U^*UXV)} = \sqrt{\rho(V^*X^*XV)} \\ &= \sqrt{\rho(X^*X)} = \|X\|,\end{aligned}$$

where the second-to-last equality is due to the fact that V^*X^*XV is similar to X^*X , so $\Lambda(V^*X^*XV) = \Lambda(X^*X)$. An important inequality involving the norms (2.27)–(2.29) is the following:

$$\|X\| \leq \sqrt{|X|_1 |X|_\infty}, \quad X \in \mathbb{C}^{m \times m}. \quad (2.31)$$

The proof of (2.31) is simple. Let $\mathbf{x} \neq \mathbf{0}$ be such that $X^*X\mathbf{x} = \lambda_{\max}(X^*X)\mathbf{x}$. Passing to the norms and using (2.25)–(2.28), we obtain

$$\|X\|^2 |\mathbf{x}|_\infty = \lambda_{\max}(X^*X) |\mathbf{x}|_\infty = |X^*X\mathbf{x}|_\infty \leq |X^*|_\infty |X|_\infty |\mathbf{x}|_\infty = |X|_1 |X|_\infty |\mathbf{x}|_\infty,$$

which yields (2.31). In view of (2.27) and (2.29), the inequality (2.31) is particularly useful to estimate the spectral norm of a matrix when we have upper bounds for its components.

A matrix $X \in \mathbb{C}^{m \times m}$ such that $XX^* = X^*X$ is said to be normal. If X is Hermitian ($X^* = X$) or skew-Hermitian ($X^* = -X$), then X is normal. If X is normal, then X is unitarily diagonalizable, meaning that there exist a unitary matrix U and a diagonal matrix D (whose diagonal entries are the eigenvalues of X) such that $X = UDU^*$; see, e.g., [17, Theorem 2.28] or [69, Corollary 7.1.4]. This result in combination with (2.30) and (2.24) implies that

$$\|X\| = \rho(X) \leq |X|_p, \quad 1 \leq p \leq \infty, \quad X \in \mathbb{C}^{m \times m} \text{ normal}. \quad (2.32)$$

2.4.2 Singular Value Decomposition

The fundamental theorem about the singular value decomposition of a matrix $X \in \mathbb{C}^{m \times m}$ is formally stated here. The proof can be found in [17, Theorem 7.8] or [69, Theorem 2.4.1].

Theorem 2.5 (singular value decomposition) *Let $X \in \mathbb{C}^{m \times m}$. Then, there exist two unitary matrices $U, V \in \mathbb{C}^{m \times m}$ and a diagonal matrix $\Sigma \in \mathbb{R}^{m \times m}$ with diagonal entries $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_m \geq 0$ such that $X = U\Sigma V^*$.*

If $X \in \mathbb{C}^{m \times m}$, any decomposition of X of the form $X = U\Sigma V^*$, in which $U, V \in \mathbb{C}^{m \times m}$ are unitary and $\Sigma \in \mathbb{R}^{m \times m}$ is diagonal with diagonal entries $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_m \geq 0$, is referred to as a singular value decomposition (SVD) of X . Some important properties of the SVD are listed below.

- If $X = U\Sigma V^*$ and $X = \tilde{U}\tilde{\Sigma}\tilde{V}^*$ are two SVDs of X , then $X^*X = V\Sigma^2V^*$ and $X^*X = \tilde{V}\tilde{\Sigma}^2\tilde{V}^*$. Hence, the diagonal entries of both Σ and $\tilde{\Sigma}$ are the square roots of the eigenvalues of X^*X . Since the diagonal entries of Σ and $\tilde{\Sigma}$ are sorted in non-increasing order, we conclude that $\Sigma = \tilde{\Sigma}$. In conclusion:

- the diagonal matrix Σ is always the same in any SVD of X ;
- the diagonal entries of Σ are referred to as the singular values of X ; they are the square roots of the eigenvalues of X^*X and are denoted by

$$\sigma_{\max}(X) = \sigma_1(X) \geq \sigma_2(X) \geq \dots \geq \sigma_m(X) = \sigma_{\min}(X).$$

- For any matrix $X \in \mathbb{C}^{m \times m}$, we have

$$\text{rank}(X) = \#\{i \in \{1, \dots, m\} : \sigma_i(X) \neq 0\}, \quad (2.33)$$

$$\|X\| = \sigma_{\max}(X) \geq |x_{ij}|, \quad i, j = 1, \dots, m. \quad (2.34)$$

Indeed, let $X = U\Sigma V^*$ be an SVD of X . Since U, V are invertible, it is clear that $\text{rank}(X) = \text{rank}(\Sigma)$ and so (2.33) holds. Moreover, the unitary invariance of $\|\cdot\|$ yields $\|X\| = \|\Sigma\| = \sigma_{\max}(X)$, which gives the equality in (2.34). To prove the inequality in (2.34), note that

$$|x_{ij}| = |(U\Sigma V^*)_{ij}| = \left| \sum_{\ell=1}^m \sigma_{\ell}(X) u_{i\ell} \overline{v_{j\ell}} \right| \leq \sigma_{\max}(X) \sum_{\ell=1}^m |u_{i\ell}| |v_{j\ell}| \leq \|X\|,$$

where in the last inequality we used the equation in (2.34), the Cauchy–Schwarz inequality for vectors of \mathbb{C}^m , and the fact that the 2-norms of the rows and columns of U and V are equal to 1.

- If $X \in \mathbb{C}^{m \times m}$ is normal, then we already noted in Sect. 2.4.1 that $X = UDU^*$ for some unitary matrix U and some diagonal matrix $D = \text{diag}_{i=1, \dots, m} \lambda_i(X)$, containing the eigenvalues of X as diagonal entries. By permuting the columns of U (if necessary), we may assume that $|\lambda_1(X)| \geq |\lambda_2(X)| \geq \dots \geq |\lambda_m(X)|$. Multiply each eigenvalue $\lambda_i(X)$ by the phase factor $e^{i\theta_i}$ such that $e^{i\theta_i} \lambda_i(X) = |\lambda_i(X)|$ and set $\Sigma = \text{diag}_{i=1, \dots, m} |\lambda_i(X)|$ and $V = e^{i\theta_1} \dots e^{i\theta_m} U$. Then $X = U\Sigma V^*$ is an SVD of X , from which we see that the singular values of X are $|\lambda_i(X)|$, $i = 1, \dots, m$. In conclusion, the singular values of a normal matrix coincide with the moduli of the eigenvalues. In view of (2.34), this also provides another proof of the equation $\|X\| = \rho(X)$ for normal matrices X .

A crucial result about the SVD is the following approximation theorem, which is sometimes referred to as the Eckart–Young theorem. It says that the matrix X_s obtained by truncating the SVD of X at the s th singular value is the closest to X in spectral norm among all the matrices with rank bounded by s ; moreover, the minimum distance $\|X - X_s\|$ equals $\sigma_{s+1}(X)$. The proof of the Eckart–Young theorem can be found, e.g., in [17, Theorem 7.13] or [69, Theorem 2.4.8]. In the corresponding statement below, we denote by $\mathbf{u}_1, \dots, \mathbf{u}_m$ and $\mathbf{v}_1, \dots, \mathbf{v}_m$ the columns of U and V , respectively.

Theorem 2.6 (the Eckart–Young theorem) *Let $X \in \mathbb{C}^{m \times m}$ and $1 \leq s \leq m - 1$. Let $X = U\Sigma V^* = \sum_{i=1}^m \sigma_i(X) \mathbf{u}_i \mathbf{v}_i^*$ be an SVD of X and set $X_s = U\Sigma_s V^* = \sum_{i=1}^s \sigma_i(X) \mathbf{u}_i \mathbf{v}_i^*$, where $\Sigma_s = \text{diag}(\sigma_1(X), \dots, \sigma_s(X), 0, \dots, 0)$. Then*

$$\min\{\|X - Y\| : \text{rank}(Y) \leq s\} = \|X - X_s\| = \sigma_{s+1}(X).$$

We conclude this section about the SVD by talking about the Moore–Penrose pseudoinverse. If $X \in \mathbb{C}^{m \times m}$ and $X = U\Sigma V^*$ is an SVD of X , we define

$$X^\dagger = V\Sigma^\dagger U^*, \quad (2.35)$$

where $\Sigma^\dagger = \text{diag}(1/\sigma_1(X), \dots, 1/\sigma_r(X), 0, \dots, 0)$ is the diagonal matrix obtained from $\Sigma = \text{diag}(\sigma_1(X), \dots, \sigma_r(X), 0, \dots, 0)$ by inverting the nonzero singular values of X ($r = \text{rank}(X)$). The matrix X^\dagger is called the Moore–Penrose pseudoinverse of X . The Moore–Penrose pseudoinverse is well-defined, in the sense that its definition (2.35) is independent of the considered SVD of X . Indeed, regardless of the considered SVD, the matrix X^\dagger defined by (2.35) is the unique matrix of $\mathbb{C}^{m \times m}$ that associates to any $\mathbf{y} \in \mathbb{C}^m$ the solution of the least squares problem

$$\min_{\mathbf{x} \in \mathbb{C}^m} \|X\mathbf{x} - \mathbf{y}\|,$$

i.e., the unique minimum norm vector \mathbf{x}^\dagger such that $\min_{\mathbf{x} \in \mathbb{C}^m} \|X\mathbf{x} - \mathbf{y}\| = \|X\mathbf{x}^\dagger - \mathbf{y}\|$. For more details, we refer the reader to [17, Chap. 7 (especially Theorems 7.1, 7.15, and p. 457)] or [69, Sect. 5.5 (especially Subsections 5.5.1–5.5.2)].

2.4.3 Schatten p -norms

Given $1 \leq p \leq \infty$ and a matrix $X \in \mathbb{C}^{m \times m}$, we denote by $\|X\|_p$ the Schatten p -norm of X , which is defined as the p -norm of the vector $(\sigma_1(X), \dots, \sigma_m(X))$ formed by the singular values of X . Note that the Schatten p -norms are unitarily invariant, i.e., $\|PXQ\|_p = \|X\|_p$ for all $p \in [1, \infty]$, all $X \in \mathbb{C}^{m \times m}$ and all unitary matrices $P, Q \in \mathbb{C}^{m \times m}$. This follows from the fact that X and UXV have the same singular values (see Theorem 2.5). The Schatten p -norms, along with all unitarily invariant norms, are deeply studied in Chap. IV of Bhatia's book [13].

If $1 \leq p, q \leq \infty$ are conjugate exponents, the following Hölder-type inequality holds for the Schatten norms (see [13, Problem III.6.2 and Corollary IV.2.6]):

$$\|XY\|_1 \leq \|X\|_p \|Y\|_q, \quad X, Y \in \mathbb{C}^{m \times m}. \quad (2.36)$$

An analogous inequality actually holds for all unitarily invariant norms, as shown in [13, Corollary IV.2.6]. If $1 \leq p, q \leq \infty$ are conjugate exponents, then

$$\|X\|_1 \leq (\text{rank}(X))^{1/q} \|X\|_p \leq m^{1/q} \|X\|_p, \quad X \in \mathbb{C}^{m \times m}, \quad (2.37)$$

where it is understood that $1/\infty = 0$. To prove (2.37), let $X = U\Sigma V^*$ be an SVD of X and let J be the matrix obtained from the identity I by setting to 0 all the diagonal

entries corresponding to indices exceeding $\text{rank}(X)$. Since $\text{rank}(X)$ is the number of nonzero singular values of X , we have $\text{rank}(J) = \text{rank}(X)$ and $J\Sigma = \Sigma$. Hence, by (2.36),

$$\|X\|_1 = \|\Sigma\|_1 = \|J\Sigma\|_1 \leq \|J\|_q \|\Sigma\|_p = (\text{rank}(X))^{1/q} \|X\|_p.$$

In the next lemma we provide a variational characterization of Schatten p -norms. A completely analogous characterization actually holds for all unitarily invariant norms, as proved in [113, Theorem 2.1]. Throughout this book, we use the abbreviations HPD and HPSD for ‘Hermitian positive definite’ and ‘Hermitian positive semidefinite’, respectively.

Lemma 2.11 *If $1 \leq p \leq \infty$ and $X \in \mathbb{C}^{m \times m}$, then*

$$\|X\|_p = \sup |(\mathbf{u}_1^* X \mathbf{v}_1, \dots, \mathbf{u}_m^* X \mathbf{v}_m)|_p, \quad (2.38)$$

where the supremum is taken over all pairs of orthonormal bases $\{\mathbf{u}_i\}_{i=1}^m$, $\{\mathbf{v}_i\}_{i=1}^m$ of \mathbb{C}^m . If moreover X is HPSD, then

$$\|X\|_p = \sup |(\mathbf{u}_1^* X \mathbf{u}_1, \dots, \mathbf{u}_m^* X \mathbf{u}_m)|_p, \quad (2.39)$$

where the supremum is taken over all orthonormal bases $\{\mathbf{u}_i\}_{i=1}^m$ of \mathbb{C}^m .

Proof Let $X = U\Sigma V^*$ be an SVD of X , so that $U^*XV = \text{diag}(\sigma_1(X), \dots, \sigma_m(X))$. If $\{\mathbf{u}_i\}_{i=1}^m$ and $\{\mathbf{v}_i\}_{i=1}^m$ are, respectively, the columns of U and V , then

$$\|X\|_p = |(\sigma_1(X), \dots, \sigma_m(X))|_p = |(\mathbf{u}_1^* X \mathbf{v}_1, \dots, \mathbf{u}_m^* X \mathbf{v}_m)|_p.$$

Hence, \leq holds in (2.38). On the other hand, suppose that U (with columns $\{\mathbf{u}_i\}_{i=1}^m$) and V (with columns $\{\mathbf{v}_i\}_{i=1}^m$) are any two unitary matrices. If $P_i = \mathbf{e}_i \mathbf{e}_i^*$ is the orthogonal projection onto the subspace of \mathbb{C}^m generated by \mathbf{e}_i (the i th vector of the canonical basis), then the singular values of $\sum_{i=1}^m P_i U^* X V P_i = \sum_{i=1}^m (\mathbf{u}_i^* X \mathbf{v}_i) P_i$ are $|\mathbf{u}_i^* X \mathbf{v}_i|$, $i = 1, \dots, m$. Thus, from the pinching inequality [13, Formula (IV.52)], we obtain

$$\begin{aligned} |(\mathbf{u}_1^* X \mathbf{v}_1, \dots, \mathbf{u}_m^* X \mathbf{v}_m)|_p &= |(|\mathbf{u}_1^* X \mathbf{v}_1|, \dots, |\mathbf{u}_m^* X \mathbf{v}_m|)|_p = \left\| \sum_{i=1}^m P_i U^* X V P_i \right\|_p \\ &\leq \|U^* X V\|_p = \|X\|_p. \end{aligned}$$

Since $\{\mathbf{u}_i\}_{i=1}^m$ and $\{\mathbf{v}_i\}_{i=1}^m$ are arbitrary orthonormal bases, we infer that also \geq holds in (2.38), and this completes the proof of (2.38).

To prove (2.39), we first note that \geq certainly holds in (2.39) by (2.38). The proof of \leq is the same as before; it suffices to observe that, since X is HPSD, we have $\lambda_i(X) = \sigma_i(X)$ for all $i = 1, \dots, m$, and, moreover, we can take an SVD of X of the form $X = U\Sigma U^*$, with $\Sigma = \text{diag}(\lambda_1(X), \dots, \lambda_m(X))$. \square

The most important among the Schatten p -norms are the Schatten 1-norm, the Schatten 2-norm and the Schatten ∞ -norm. The Schatten ∞ -norm $\|X\|_\infty$ is equal to $\sigma_{\max}(X)$ by definition, and since $\sigma_{\max}(X) = \|X\|$ we have $\|X\|_\infty = \|X\|$. Moreover, it is not difficult to see that, in the case $p = \infty$, Eqs. (2.38)–(2.39) yield

$$\sigma_{\max}(X) = \|X\| = \|X\|_\infty = \sup_{\|u\|=\|v\|=1} |u^* X v|, \quad X \in \mathbb{C}^{m \times m}, \quad (2.40)$$

$$\lambda_{\max}(X) = \|X\| = \|X\|_\infty = \sup_{\|u\|=1} u^* X u, \quad X \in \mathbb{C}^{m \times m} \text{ HPSPD}. \quad (2.41)$$

The Schatten 2-norm $\|X\|_2$ is also known as the Frobenius norm and admits an explicit expression in terms of the components of X , namely

$$\|X\|_2 = \left(\sum_{i,j=1}^m |x_{ij}|^2 \right)^{1/2}, \quad X \in \mathbb{C}^{m \times m}.$$

Indeed, recalling that the squares of the singular values of X coincide with the eigenvalues of X^*X , we have

$$\|X\|_2^2 = \sum_{i=1}^m (\sigma_i(X))^2 = \text{trace}(X^*X) = \sum_{i,j=1}^m |x_{ij}|^2.$$

The Schatten 1-norm $\|X\|_1$ is also known under the names of trace-norm and nuclear norm. From (2.37) we immediately obtain the following important trace-norm inequality:

$$\|X\|_1 \leq \text{rank}(X) \|X\| \leq m \|X\|, \quad X \in \mathbb{C}^{m \times m}. \quad (2.42)$$

Note that (2.42) is actually a direct consequence of the equation $\sigma_{\max}(X) = \|X\|$ and the definition $\|X\|_1 = \sum_{i=1}^m \sigma_i(X) = \sum_{i=1}^{\text{rank}(X)} \sigma_i(X)$. Other interesting trace-norm inequalities, which provide an upper and lower bound for the trace-norm in terms of the components, are the following:

$$|\text{trace}(X)| \leq \|X\|_1 \leq \sum_{i,j=1}^m |x_{ij}|, \quad X \in \mathbb{C}^{m \times m}. \quad (2.43)$$

The proof is simple. Let $X = U \Sigma V^*$ be an SVD of X . Then, setting $Q = V U^*$, the matrix Q is unitary and we have

$$\begin{aligned} \|X\|_1 &= \text{trace}(\Sigma) = \text{trace}(U^* X V) = \text{trace}(X Q) \\ &\leq \sum_{i=1}^m \sum_{k=1}^m |x_{ik} q_{ki}| \leq \sum_{i=1}^m \max_{k=1, \dots, m} |q_{ki}| \sum_{k=1}^m |x_{ik}| \leq \sum_{i=1}^m \sum_{k=1}^m |x_{ik}|. \end{aligned}$$

Moreover, using the Cauchy–Schwarz inequality for vectors of \mathbb{C}^m and the fact that the 2-norms of the rows and columns of U and V are equal to 1, we get

$$\begin{aligned} |\text{trace}(X)| &= \left| \sum_{i=1}^m (U \Sigma V^*)_{ii} \right| = \left| \sum_{i=1}^m \sum_{k=1}^m \sigma_k(X) u_{ik} \bar{v}_{ik} \right| = \left| \sum_{k=1}^m \sigma_k(X) \sum_{i=1}^m u_{ik} \bar{v}_{ik} \right| \\ &\leq \sum_{k=1}^m \sigma_k(X) \sum_{i=1}^m |u_{ik}| |v_{ik}| \leq \sum_{k=1}^m \sigma_k(X) = \|X\|_1. \end{aligned}$$

It is worth noting that the left inequality in (2.43) is a special case of Weyl’s majorization theorem [13, Theorem II.3.6]. This theorem actually implies a stronger inequality, i.e.,

$$\sum_{j=1}^m |\lambda_j(X)| \leq \|X\|_1, \quad X \in \mathbb{C}^{m \times m}. \quad (2.44)$$

2.4.4 Singular Value and Eigenvalue Inequalities

We begin by reporting the minimax principles for singular values and eigenvalues. We use the notation $V \subseteq_{\text{sp.}} \mathbb{C}^m$ to indicate that V is a subspace of \mathbb{C}^m .

Theorem 2.7 (minimax principle for singular values) *Let $X \in \mathbb{C}^{m \times m}$ and let $\sigma_1(X) \geq \dots \geq \sigma_m(X)$ be the singular values of X sorted, as always, in non-increasing order. Then,*

$$\sigma_j(X) = \max_{\substack{V \subseteq_{\text{sp.}} \mathbb{C}^m \\ \dim V = j}} \min_{\substack{\mathbf{x} \in V \\ \|\mathbf{x}\|=1}} \|X\mathbf{x}\| = \min_{\substack{V \subseteq_{\text{sp.}} \mathbb{C}^m \\ \dim V = m-j+1}} \max_{\substack{\mathbf{x} \in V \\ \|\mathbf{x}\|=1}} \|X\mathbf{x}\|, \quad j = 1, \dots, m.$$

In particular,

$$\sigma_{\max}(X) = \max_{\substack{\mathbf{x} \in \mathbb{C}^m \\ \|\mathbf{x}\|=1}} \|X\mathbf{x}\|, \quad \sigma_{\min}(X) = \min_{\substack{\mathbf{x} \in \mathbb{C}^m \\ \|\mathbf{x}\|=1}} \|X\mathbf{x}\|.$$

Theorem 2.8 (minimax principle for eigenvalues) *Let $X \in \mathbb{C}^{m \times m}$ be Hermitian and let $\lambda_1(X) \geq \dots \geq \lambda_m(X)$ be the eigenvalues of X sorted in non-increasing order. Then,*

$$\lambda_j(X) = \max_{\substack{V \subseteq_{\text{sp.}} \mathbb{C}^m \\ \dim V = j}} \min_{\substack{\mathbf{x} \in V \\ \|\mathbf{x}\|=1}} (\mathbf{x}^* X \mathbf{x}) = \min_{\substack{V \subseteq_{\text{sp.}} \mathbb{C}^m \\ \dim V = m-j+1}} \max_{\substack{\mathbf{x} \in V \\ \|\mathbf{x}\|=1}} (\mathbf{x}^* X \mathbf{x}), \quad j = 1, \dots, m.$$

In particular,

$$\lambda_{\max}(X) = \max_{\substack{\mathbf{x} \in \mathbb{C}^m \\ \|\mathbf{x}\|=1}} (\mathbf{x}^* X \mathbf{x}), \quad \lambda_{\min}(X) = \min_{\substack{\mathbf{x} \in \mathbb{C}^m \\ \|\mathbf{x}\|=1}} (\mathbf{x}^* X \mathbf{x}).$$

Theorem 2.7 follows from Theorem 2.8 applied to the Hermitian matrix X^*X , whose eigenvalues are the squares of the singular values of X . Theorem 2.8 is proved, e.g., in [17, Theorem 6.7], [69, Theorem 8.1.2] and [13, Corollary III.1.2]. As a consequence of Theorem 2.8, for all matrices $X \in \mathbb{C}^{m \times m}$ we have the following localization of the spectrum:

$$\Lambda(X) \subseteq [\lambda_{\min}(\Re(X)), \lambda_{\max}(\Re(X))] \times [\lambda_{\min}(\Im(X)), \lambda_{\max}(\Im(X))] \subset \mathbb{C}. \quad (2.45)$$

Indeed, if λ is an eigenvalue of X and \mathbf{x} is a corresponding eigenvector with $\|\mathbf{x}\| = 1$, then, by Theorem 2.8 applied to $\Re(X)$ and $\Im(X)$,

$$\begin{aligned} \lambda &= \mathbf{x}^* X \mathbf{x} = \mathbf{x}^* \Re(X) \mathbf{x} + i \mathbf{x}^* \Im(X) \mathbf{x} \\ &\in [\lambda_{\min}(\Re(X)), \lambda_{\max}(\Re(X))] \times [\lambda_{\min}(\Im(X)), \lambda_{\max}(\Im(X))]. \end{aligned}$$

In the next theorems, we provide some important perturbation and interlacing theorems for singular values and eigenvalues.

Theorem 2.9 (perturbation theorem for singular values) *Let $X, Y \in \mathbb{C}^{m \times m}$ and let $\sigma_1(X) \geq \dots \geq \sigma_m(X)$ and $\sigma_1(Y) \geq \dots \geq \sigma_m(Y)$ be their respective singular values sorted, as always, in non-increasing order. Then,*

$$|\sigma_j(X) - \sigma_j(Y)| \leq \|X - Y\|, \quad j = 1, \dots, m.$$

Theorem 2.10 (perturbation theorem for eigenvalues) *Let $X, Y \in \mathbb{C}^{m \times m}$ be Hermitian and let $\lambda_1(X) \geq \dots \geq \lambda_m(X)$ and $\lambda_1(Y) \geq \dots \geq \lambda_m(Y)$ be their respective eigenvalues sorted in non-increasing order. Then,*

$$|\lambda_j(X) - \lambda_j(Y)| \leq \|X - Y\|, \quad j = 1, \dots, m.$$

Theorem 2.9 (resp., 2.10) is a corollary of the minimax principle for singular values (resp., eigenvalues). For example, to prove Theorem 2.10 one simply observes that

$$|\mathbf{x}^*(X - Y)\mathbf{x}| \leq \|\mathbf{x}\| \|(X - Y)\mathbf{x}\| \leq \|\mathbf{x}\| \|X - Y\| \|\mathbf{x}\|,$$

hence $|\mathbf{x}^*(X - Y)\mathbf{x}| \leq \|X - Y\|$ for all vectors \mathbf{x} such that $\|\mathbf{x}\| = 1$, and

$$\begin{aligned} \lambda_j(X) &= \max_{\substack{V \subseteq \mathbb{C}^m \\ \dim V = j}} \min_{\substack{\mathbf{x} \in V \\ \|\mathbf{x}\| = 1}} (\mathbf{x}^* X \mathbf{x}) = \max_{\substack{V \subseteq \mathbb{C}^m \\ \dim V = j}} \min_{\substack{\mathbf{x} \in V \\ \|\mathbf{x}\| = 1}} (\mathbf{x}^*(X - Y)\mathbf{x} + \mathbf{x}^* Y \mathbf{x}) \\ &\leq \|X - Y\| + \max_{\substack{V \subseteq \mathbb{C}^m \\ \dim V = j}} \min_{\substack{\mathbf{x} \in V \\ \|\mathbf{x}\| = 1}} (\mathbf{x}^* Y \mathbf{x}) = \|X - Y\| + \lambda_j(Y). \end{aligned}$$

Theorem 2.9 is proved likewise. Theorem 2.10 is known as Weyl's perturbation theorem [13, Corollary III.2.6]. We refer the reader to [13, Problem II.6.13] for a general perturbation theorem for singular values, which extends Theorem 2.9.

To simplify the statement of Theorems 2.11–2.12, we here adopt the following convention: for each matrix $X \in \mathbb{C}^{m \times m}$ with singular values $\sigma_1(X) \geq \dots \geq \sigma_m(X)$, let $\sigma_j(X) = +\infty$ if $j < 1$ and $\sigma_j(X) = -\infty$ if $j > m$; for each Hermitian matrix $X \in \mathbb{C}^{m \times m}$ with eigenvalues $\lambda_1(X) \geq \dots \geq \lambda_m(X)$, let $\lambda_j(X) = +\infty$ if $j < 1$ and $\lambda_j(X) = -\infty$ if $j > m$.

Theorem 2.11 (interlacing theorem for singular values) *Let $Y = X + E$, where $X, E \in \mathbb{C}^{m \times m}$ and $\text{rank}(E) \leq k$. Let $\sigma_1(X) \geq \dots \geq \sigma_m(X)$ and $\sigma_1(Y) \geq \dots \geq \sigma_m(Y)$ be the singular values of X and Y . Then,*

$$\sigma_{j-k}(X) \geq \sigma_j(Y) \geq \sigma_{j+k}(X), \quad j = 1, \dots, m.$$

Proof For every $A \in \mathbb{C}^{m \times m}$, the eigenvalues of the $2m \times 2m$ Hermitian matrix

$$\tilde{A} = \begin{bmatrix} O & A \\ A^* & O \end{bmatrix}$$

are $\sigma_j(A)$, $-\sigma_j(A)$, $j = 1, \dots, m$. Indeed, if $A = U\Sigma V^*$ is an SVD of A , then

$$\begin{aligned} \tilde{A} &= \begin{bmatrix} O & U\Sigma V^* \\ V\Sigma U^* & O \end{bmatrix} = \begin{bmatrix} U & O \\ O & V \end{bmatrix} \begin{bmatrix} O & \Sigma \\ \Sigma & O \end{bmatrix} \begin{bmatrix} U^* & O \\ O & V^* \end{bmatrix} \\ &= \begin{bmatrix} U & O \\ O & V \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} I & I \\ I & -I \end{bmatrix} \begin{bmatrix} \Sigma & O \\ O & -\Sigma \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} I & I \\ I & -I \end{bmatrix} \begin{bmatrix} U^* & O \\ O & V^* \end{bmatrix} = Q \begin{bmatrix} \Sigma & O \\ O & -\Sigma \end{bmatrix} Q^*, \end{aligned}$$

where

$$Q = \begin{bmatrix} U & O \\ O & V \end{bmatrix} \frac{1}{\sqrt{2}} \begin{bmatrix} I & I \\ I & -I \end{bmatrix}$$

is unitary, being the product of two unitary matrices. Therefore, by applying Theorem 2.12 with \tilde{Y} , \tilde{X} , \tilde{E} in place of Y , X , E , we obtain the thesis. \square

Theorem 2.12 (interlacing theorem for eigenvalues) *Let $Y = X + E$, where $X, E \in \mathbb{C}^{m \times m}$ are Hermitian. Let $\lambda_1(X) \geq \dots \geq \lambda_m(X)$ and $\lambda_1(Y) \geq \dots \geq \lambda_m(Y)$ be the eigenvalues of X and Y , and let k^+ , k^- be the number of positive and negative eigenvalues of E :*

$$k^+ = \#\{j \in \{1, \dots, m\} : \lambda_j(E) > 0\}, \quad k^- = \#\{j \in \{1, \dots, m\} : \lambda_j(E) < 0\}.$$

Then,

$$\lambda_{j-k^+}(X) \geq \lambda_j(Y) \geq \lambda_{j+k^-}(X), \quad j = 1, \dots, m.$$

In particular, if $\text{rank}(E) \leq k$ then

$$\lambda_{j-k}(X) \geq \lambda_j(Y) \geq \lambda_{j+k}(X), \quad j = 1, \dots, m.$$

Proof Throughout this proof, we adopt the convention stated before Theorem 2.11. Moreover, we will make use of the following result.

Let $B = A + \eta \mathbf{u} \mathbf{u}^*$, where $A \in \mathbb{C}^{m \times m}$ is a Hermitian matrix, $\mathbf{u} \in \mathbb{C}^m$ is a (column) vector and $\eta \geq 0$. Let $\lambda_1(A) \geq \dots \geq \lambda_m(A)$ and $\lambda_1(B) \geq \dots \geq \lambda_m(B)$ be the eigenvalues of A and B . Then,

$$\lambda_{j-1}(A) \geq \lambda_j(B) \geq \lambda_j(A), \quad j = 1, \dots, m,$$

or, equivalently,

$$\lambda_j(B) \geq \lambda_j(A) \geq \lambda_{j+1}(B), \quad j = 1, \dots, m.$$

This result, which is actually a special case of Theorem 2.12, can be found in [17, Theorem 6.12], [69, Theorem 8.1.8] and [13, Exercise III.2.4].

Let us now prove Theorem 2.12. Let $\alpha_1, \dots, \alpha_{k^+}$ and $\beta_1, \dots, \beta_{k^-}$ be, respectively, the positive and the negative eigenvalues of E . Since E is Hermitian, we can write

$$\begin{aligned} E &= Q \operatorname{diag}(\alpha_1, \dots, \alpha_{k^+}, \beta_1, \dots, \beta_{k^-}, 0, \dots, 0) Q^* \\ &= \sum_{i=1}^{k^+} \alpha_i \mathbf{u}_i \mathbf{u}_i^* + \sum_{i=1}^{k^-} \beta_i \mathbf{v}_i \mathbf{v}_i^* = E^+ + E^-, \end{aligned}$$

where:

- Q is a unitary matrix;
- $\mathbf{u}_1, \dots, \mathbf{u}_{k^+}$ are the columns $1, \dots, k^+$ of Q , which correspond to the positive eigenvalues $\alpha_1, \dots, \alpha_{k^+}$;
- $\mathbf{v}_1, \dots, \mathbf{v}_{k^-}$ are the columns $k^+ + 1, \dots, k^+ + k^-$ of Q , which correspond to the negative eigenvalues $\beta_1, \dots, \beta_{k^-}$;
- $E^+ = \sum_{i=1}^{k^+} \alpha_i \mathbf{u}_i \mathbf{u}_i^*$ and $E^- = \sum_{i=1}^{k^-} \beta_i \mathbf{v}_i \mathbf{v}_i^*$.

By repeated applications of the result quoted above, for $j = 1, \dots, m$ we obtain

$$\begin{aligned} \lambda_{j-1}(X) &\geq \lambda_j(X + \alpha_1 \mathbf{u}_1 \mathbf{u}_1^*) \geq \lambda_j(X), \\ \lambda_{j-1}(X + \alpha_1 \mathbf{u}_1 \mathbf{u}_1^*) &\geq \lambda_j(X + \alpha_1 \mathbf{u}_1 \mathbf{u}_1^* + \alpha_2 \mathbf{u}_2 \mathbf{u}_2^*) \geq \lambda_j(X + \alpha_1 \mathbf{u}_1 \mathbf{u}_1^*), \\ &\dots \\ \lambda_{j-1}\left(X + \sum_{i=1}^{k^+-1} \alpha_i \mathbf{u}_i \mathbf{u}_i^*\right) &\geq \lambda_j(X + E^+) \geq \lambda_j\left(X + \sum_{i=1}^{k^+-1} \alpha_i \mathbf{u}_i \mathbf{u}_i^*\right). \end{aligned}$$

Thus,

$$\lambda_{j-k^+}(X) \geq \lambda_j(X + E^+) \geq \lambda_j(X), \quad j = 1, \dots, m. \quad (2.46)$$

By repeated applications of the result quoted above, for $j = 1, \dots, m$ we obtain

$$\begin{aligned}
\lambda_j(X + E^+) &\geq \lambda_j(X + E^+ + \beta_1 \mathbf{v}_1 \mathbf{v}_1^*) \geq \lambda_{j+1}(X + E^+), \\
\lambda_j(X + E^+ + \beta_1 \mathbf{v}_1 \mathbf{v}_1^*) &\geq \lambda_j(X + E^+ + \beta_1 \mathbf{v}_1 \mathbf{v}_1^* + \beta_2 \mathbf{v}_2 \mathbf{v}_2^*) \geq \lambda_{j+1}(X + E^+ + \beta_1 \mathbf{v}_1 \mathbf{v}_1^*), \\
&\dots \\
\lambda_j\left(X + E^+ + \sum_{i=1}^{k^- - 1} \beta_i \mathbf{v}_i \mathbf{v}_i^*\right) &\geq \lambda_j(Y) \geq \lambda_{j+1}\left(X + E^+ + \sum_{i=1}^{k^- - 1} \beta_i \mathbf{v}_i \mathbf{v}_i^*\right).
\end{aligned}$$

Thus,

$$\lambda_j(X + E^+) \geq \lambda_j(Y) \geq \lambda_{j+k^-}(X + E^+), \quad j = 1, \dots, m. \quad (2.47)$$

By combining (2.46) and (2.47), we obtain the thesis. \square

Important inequalities involving the imaginary parts of the eigenvalues of X and the eigenvalues of $\Im(X)$ are provided in the next theorem, which is due to Ky-Fan.

Theorem 2.13 *Let $X \in \mathbb{C}^{m \times m}$ and label the eigenvalues of X and $\Im(X)$ so that $\Im(\lambda_1(X)) \geq \dots \geq \Im(\lambda_m(X))$ and $\lambda_1(\Im(X)) \geq \dots \geq \lambda_m(\Im(X))$. Then,*

$$\sum_{j=1}^k \Im(\lambda_j(X)) \leq \sum_{j=1}^k \lambda_j(\Im(X)) \quad (2.48)$$

for all $k = 1, \dots, m$. Moreover, for $k = m$, the equality holds in (2.48).

The proof of Theorem 2.13 can be found in [13, Proposition III.5.3]. Note that Theorem 2.13 is stated in [13] with ‘ \Re ’ instead of ‘ \Im ’, but this is not an issue because $\Re(X) = \Im(iX)$.

Lemma 2.12 *For every $X \in \mathbb{C}^{m \times m}$,*

$$\sum_{j=1}^m |\Im(\lambda_j(X))| \leq \|\Im(X)\|_1. \quad (2.49)$$

In particular, for every $X \in \mathbb{C}^{m \times m}$ and every $\varepsilon > 0$ we have

$$\#\{j \in \{1, \dots, m\} : |\Im(\lambda_j(X))| > \varepsilon\} \leq \frac{\|\Im(X)\|_1}{\varepsilon}, \quad (2.50)$$

and if $\Lambda(\Re(X))$ is contained in the interval $I \subseteq \mathbb{R}$, then

$$\#\{j \in \{1, \dots, m\} : \lambda_j(X) \notin I \times [-\varepsilon, \varepsilon]\} \leq \frac{\|\Im(X)\|_1}{\varepsilon}. \quad (2.51)$$

Proof Label the eigenvalues of X and $\Im(X)$ so that $\Im(\lambda_1(X)) \geq \dots \geq \Im(\lambda_m(X))$ and $\lambda_1(\Im(X)) \geq \dots \geq \lambda_m(\Im(X))$. By Theorem 2.13,

$$\begin{aligned}
\sum_{j: \lambda_j(\Im(X)) \geq 0} \lambda_j(\Im(X)) &= \max_{k=1, \dots, m} \sum_{j=1}^k \lambda_j(\Im(X)) \\
&\geq \max_{k=1, \dots, m} \sum_{j=1}^k \Im(\lambda_j(X)) = \sum_{j: \Im(\lambda_j(X)) \geq 0} \Im(\lambda_j(X)). \quad (2.52)
\end{aligned}$$

Again by Theorem 2.13,

$$\begin{aligned}
&\sum_{j: \lambda_j(\Im(X)) < 0} \lambda_j(\Im(X)) + \sum_{j: \lambda_j(\Im(X)) \geq 0} \lambda_j(\Im(X)) \\
&= \sum_{j: \Im(\lambda_j(X)) < 0} \Im(\lambda_j(X)) + \sum_{j: \Im(\lambda_j(X)) \geq 0} \Im(\lambda_j(X)),
\end{aligned}$$

so (2.52) implies that

$$\sum_{j: \lambda_j(\Im(X)) < 0} \lambda_j(\Im(X)) \leq \sum_{j: \Im(\lambda_j(X)) < 0} \Im(\lambda_j(X)). \quad (2.53)$$

Since $\Im(X)$ is Hermitian, its trace-norm equals the sum of the absolute values of its eigenvalues. Thus, by (2.52)–(2.53),

$$\begin{aligned}
\|\Im(X)\|_1 &= \sum_{j=1}^m |\lambda_j(\Im(X))| \\
&= \sum_{j: \lambda_j(\Im(X)) \geq 0} \lambda_j(\Im(X)) - \sum_{j: \lambda_j(\Im(X)) < 0} \lambda_j(\Im(X)) \\
&\geq \sum_{j: \Im(\lambda_j(X)) \geq 0} \Im(\lambda_j(X)) - \sum_{j: \Im(\lambda_j(X)) < 0} \Im(\lambda_j(X)) \\
&= \sum_{j=1}^m |\Im(\lambda_j(X))|.
\end{aligned}$$

This proves (2.49). The inequality (2.50) follows from (2.49) and the observation that

$$\begin{aligned}
\sum_{j=1}^m |\Im(\lambda_j(X))| &\geq \sum_{j: |\Im(\lambda_j(X))| > \varepsilon} |\Im(\lambda_j(X))| \\
&\geq \varepsilon \cdot \#\{j \in \{1, \dots, m\} : |\Im(\lambda_j(X))| > \varepsilon\}.
\end{aligned}$$

The inequality (2.51) follows from (2.50) and (2.45). \square

2.4.5 Tensor Products and Direct Sums

If X, Y are matrices of any dimension, say $X \in \mathbb{C}^{m_1 \times m_2}$ and $Y \in \mathbb{C}^{\ell_1 \times \ell_2}$, the tensor (Kronecker) product of X and Y is the $m_1 \ell_1 \times m_2 \ell_2$ matrix defined by

$$X \otimes Y = [x_{ij}Y]_{\substack{i=1,\dots,m_1 \\ j=1,\dots,m_2}} = \begin{bmatrix} x_{11}Y & \cdots & x_{1m_2}Y \\ \vdots & & \vdots \\ x_{m_11}Y & \cdots & x_{m_1m_2}Y \end{bmatrix},$$

and the direct sum of X and Y is the $(m_1 + \ell_1) \times (m_2 + \ell_2)$ matrix defined by

$$X \oplus Y = \text{diag}(X, Y) = \begin{bmatrix} X & O \\ O & Y \end{bmatrix}.$$

Tensor products and direct sums possess a lot of nice algebraic properties.

- (i) Associativity: for all matrices X, Y, Z ,

$$\begin{aligned} (X \otimes Y) \otimes Z &= X \otimes (Y \otimes Z), \\ (X \oplus Y) \oplus Z &= X \oplus (Y \oplus Z). \end{aligned}$$

We can therefore omit parentheses in expressions like $X_1 \otimes X_2 \otimes \cdots \otimes X_d$ or $X_1 \oplus X_2 \oplus \cdots \oplus X_d$.

- (ii) If X_1, X_2 can be multiplied and Y_1, Y_2 can be multiplied, then

$$\begin{aligned} (X_1 \otimes Y_1)(X_2 \otimes Y_2) &= (X_1 X_2) \otimes (Y_1 Y_2), \\ (X_1 \oplus Y_1)(X_2 \oplus Y_2) &= (X_1 X_2) \oplus (Y_1 Y_2). \end{aligned}$$

- (iii) For all matrices X, Y ,

$$\begin{aligned} (X \otimes Y)^* &= X^* \otimes Y^*, & (X \otimes Y)^T &= X^T \otimes Y^T \\ (X \oplus Y)^* &= X^* \oplus Y^*, & (X \oplus Y)^T &= X^T \oplus Y^T. \end{aligned}$$

- (iv) Bilinearity (of tensor products): for each fixed matrix X , the application

$$Y \mapsto X \otimes Y$$

is linear on $\mathbb{C}^{\ell_1 \times \ell_2}$ for all $\ell_1, \ell_2 \in \mathbb{N}$; for each fixed matrix Y , the application

$$X \mapsto X \otimes Y$$

is linear on $\mathbb{C}^{m_1 \times m_2}$ for all $m_1, m_2 \in \mathbb{N}$.

From (i)–(iv), a lot of other interesting properties follow. For example, if X, Y are invertible, then $X \otimes Y$ is invertible, with inverse $X^{-1} \otimes Y^{-1}$. If X, Y are normal (resp., Hermitian, symmetric, unitary) then $X \otimes Y$ is also normal (resp., Hermitian, symmetric, unitary). If $X \in \mathbb{C}^{m \times m}$ and $Y \in \mathbb{C}^{\ell \times \ell}$, the eigenvalues and singular values of $X \otimes Y$ are

$$\{\lambda_i(X)\lambda_j(Y) : i = 1, \dots, m, j = 1, \dots, \ell\}, \quad (2.54)$$

$$\{\sigma_i(X)\sigma_j(Y) : i = 1, \dots, m, j = 1, \dots, \ell\}; \quad (2.55)$$

and the eigenvalues and singular values of $X \oplus Y$ are

$$\{\lambda_i(X) : i = 1, \dots, m\} \cup \{\lambda_j(Y) : j = 1, \dots, \ell\}, \quad (2.56)$$

$$\{\sigma_i(X) : i = 1, \dots, m\} \cup \{\sigma_j(Y) : j = 1, \dots, \ell\}; \quad (2.57)$$

see Exercise 2.5. In particular, for all $X \in \mathbb{C}^{m \times m}$, $Y \in \mathbb{C}^{\ell \times \ell}$, and $1 \leq p \leq \infty$, we have

$$\|X \otimes Y\|_p = \|X\|_p \|Y\|_p, \quad (2.58)$$

$$\|X \oplus Y\|_p = |(\|X\|_p, \|Y\|_p)|_p = \begin{cases} (\|X\|_p^p + \|Y\|_p^p)^{1/p}, & \text{if } 1 \leq p < \infty, \\ \max(\|X\|_\infty, \|Y\|_\infty), & \text{if } p = \infty. \end{cases} \quad (2.59)$$

and

$$\text{rank}(X \otimes Y) = \text{rank}(X)\text{rank}(Y), \quad (2.60)$$

$$\text{rank}(X \oplus Y) = \text{rank}(X) + \text{rank}(Y). \quad (2.61)$$

Exercise 2.5 Let $X \in \mathbb{C}^{m \times m}$ and $Y \in \mathbb{C}^{\ell \times \ell}$. Show that:

- (a) the eigenvalues and the singular values of $X \otimes Y$ are given by (2.54)–(2.55);
- (b) the eigenvalues and the singular values of $X \oplus Y$ are given by (2.56)–(2.57).

2.4.6 Matrix Functions

Given a diagonalizable matrix $A \in \mathbb{C}^{m \times m}$, if $\lambda_1, \dots, \lambda_t$ are the *distinct* eigenvalues of A and V_1, \dots, V_t are their respective eigenspaces, we have

$$\mathbb{C}^m = \bigoplus_{i=1}^t V_i.$$

For each function $f : \Lambda(A) \rightarrow \mathbb{C}$ we define $f(A)$ as the matrix such that

$$f(A)\mathbf{v} = f(\lambda_i)\mathbf{v} \text{ for every } \mathbf{v} \in V_i \text{ and every } i = 1, \dots, t. \quad (2.62)$$

In practice, $f(A)$ is the matrix that possesses the same eigenspaces V_1, \dots, V_t as A with corresponding eigenvalues $f(\lambda_1), \dots, f(\lambda_t)$. Note that such a matrix $f(A)$ exists and is unique. To show the uniqueness, simply note that, if B is another matrix such that $B\mathbf{v} = f(\lambda_i)\mathbf{v} = f(A)\mathbf{v}$ for every $\mathbf{v} \in V_i$ and every $i = 1, \dots, t$, then B coincides with $f(A)$ on each basis of \mathbb{C}^m formed by eigenvectors of A , hence $B = f(A)$. To show the existence, fix a basis of \mathbb{C}^m formed by eigenvectors of A , define $f(A)$ on this basis in the unique possible way to meet (2.62), and extend the definition to the whole \mathbb{C}^m by linearity. It is not difficult to check that the matrix $f(A)$ defined in this way satisfies (2.62).

Now, let $A \in \mathbb{C}^{m \times m}$ be diagonalizable and let $\lambda_1, \dots, \lambda_m$ denote *all* the eigenvalues of A , i.e., all the roots of the characteristic polynomial of A , each of them counted with its multiplicity. If $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ is a basis of \mathbb{C}^m formed by eigenvectors of A ,

$$A\mathbf{v}_i = \lambda_i\mathbf{v}_i, \quad i = 1, \dots, m,$$

then for each $f : \Lambda(A) \rightarrow \mathbb{C}$ we have

$$f(A)\mathbf{v}_i = f(\lambda_i)\mathbf{v}_i, \quad i = 1, \dots, m.$$

This is a spectral decomposition of $f(A)$, which can be rewritten in matrix form as

$$f(A)V = V \begin{bmatrix} f(\lambda_1) & & & \\ & f(\lambda_2) & & \\ & & \ddots & \\ & & & f(\lambda_m) \end{bmatrix}, \quad V = \left[\begin{array}{c|c|c|c} \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_m \end{array} \right],$$

or, equivalently,

$$f(A) = V \begin{bmatrix} f(\lambda_1) & & & \\ & f(\lambda_2) & & \\ & & \ddots & \\ & & & f(\lambda_m) \end{bmatrix} V^{-1}, \quad V = \left[\begin{array}{c|c|c|c} \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_m \end{array} \right]. \quad (2.63)$$

As a straightforward consequence of (2.63), if $f(\lambda) = \lambda$ then $f(A) = A$. Moreover, if A is invertible and $f(\lambda) = \lambda^{-1}$, then $f(A) = A^{-1}$. Further properties of matrix functions are given in Exercise 2.6.

Exercise 2.6 Let A be a diagonalizable matrix. Prove the following properties.

- If $\alpha, \beta \in \mathbb{C}$ and $f, g : \Lambda(A) \rightarrow \mathbb{C}$, let $\alpha f + \beta g : \Lambda(A) \rightarrow \mathbb{C}$ be defined in the natural way as $(\alpha f + \beta g)(\lambda) = \alpha f(\lambda) + \beta g(\lambda)$. Then, $(\alpha f + \beta g)(A) = \alpha f(A) + \beta g(A)$.
- If $f, g : \Lambda(A) \rightarrow \mathbb{C}$, let $fg : \Lambda(A) \rightarrow \mathbb{C}$ be the product function $(fg)(\lambda) = f(\lambda)g(\lambda)$. Then, $(fg)(A) = f(A)g(A)$. Moreover, $f(A)g(A) = g(A)f(A)$, i.e., two functions of the same matrix always commute.

- (c) Suppose $\Lambda(A) \subset (0, \infty)$, so that the functions $\lambda^{1/2}, \lambda^{-1/2} : \Lambda(A) \rightarrow \mathbb{R}$ are well-defined and hence also the matrices $A^{1/2}, A^{-1/2}$ via definition (2.62). Then, $A^{1/2}A^{-1/2} = I$, $(A^{1/2})^2 = A$, $(A^{-1/2})^2 = A^{-1}$.
- (d) If $p(\lambda) = \sum_{j=0}^r a_j \lambda^j$ is a polynomial, then the matrix $p(A)$ obtained from definition (2.62) coincides with $\sum_{j=0}^r a_j A^j$.
- (e) $e^A = \sum_{k=0}^{\infty} \frac{A^k}{k!}$.

For more on matrix functions, including the definition of $f(A)$ in the case where A is not diagonalizable, we refer the reader to Higham's book [71].

Generalized Locally Toeplitz Sequences: Theory and
Applications

Volume I

Garoni, C.; Serra-Capizzano, S.

2017, XI, 312 p. 16 illus. in color., Hardcover

ISBN: 978-3-319-53678-1