# Mathematical Introduction to Deep Learning

Joshua Lee Padgett

July 29, 2021

# Preface

These lecture notes are far away from being complete and remain under construction. In particular, these lecture notes do not yet contain a suitable comparison of the presented material with existing results, arguments, and notions in the literature. This will be the subject of a future version of these lecture notes.

Joshua Lee Padgett

# Contents

# Contents

# Contents

# Contents

# Contents

# Chapter 1

# Introduction

## 1.1 Introductory comments on supervised learning

Very roughly speaking, the field *deep learning* can be divided into three subfields, deep *supervised learning*, deep *unsupervised learning*, and deep *reinforcement learning*. Algorithms in deep supervised learning seem often to be most accessible for a mathematical analysis. In the following we briefly sketch in a special situation some ideas of deep supervised learning.

Let $d, M \in \mathbb{N} = \{1, 2, 3, \dots\}$, $\mathcal{E} \in C(\mathbb{R}^d, \mathbb{R})$, $x_1, x_2, \dots, x_{M+1} \in \mathbb{R}^d$, $y_1, y_2, \dots, y_M \in \mathbb{R}$ satisfy for all $m \in \{1, 2, \dots, M\}$ that

$$y_m = \mathcal{E}(x_m). \tag{1.1}$$

In the framework described in the previous sentence we think of $M \in \mathbb{N}$ as the number of available input-output data pairs, we think of $d \in \mathbb{N}$ as the dimension of the input data, we think of $\mathcal{E} \colon \mathbb{R}^d \to \mathbb{R}$ as an unknown function which relates input and output data through (1.1), we think of $x_1, x_2, \dots, x_{M+1} \in \mathbb{R}^d$ as the available known input data, and we think of $y_1, y_2, \dots, y_M \in \mathbb{R}$ as the available known output data. The key question in the context of supervised learning is then that one intends to approximately compute the output $\mathcal{E}(x_{M+1})$ of the $(M+1)$-th input data $x_{M+1}$ without using explicit knowledge of the function $\mathcal{E} \colon \mathbb{R}^d \to \mathbb{R}$ but instead by using the knowledge of the $M$ input-output data pairs $(x_1, y_1) = (x_1, \mathcal{E}(x_1)), (x_2, y_2) = (x_2, \mathcal{E}(x_2)), \dots, (x_M, y_M) = (x_M, \mathcal{E}(x_M)) \in \mathbb{R}^d \times \mathbb{R}$. To accomplish this, one considers the optimization problem of approximately computing global minima of the function $\Phi \colon C(\mathbb{R}^d, \mathbb{R}) \to [0, \infty)$ which satisfies for all $\phi \in C(\mathbb{R}^d, \mathbb{R})$ that

$$\Phi(\phi) = \sum_{m=1}^{M} |\phi(x_m) - y_m|^2. \tag{1.2}$$

Observe that (1.1) ensures that $\Phi(\mathcal{E}) = 0$ and, in particular, we have that the unknown function $\mathcal{E} \colon \mathbb{R}^d \to \mathbb{R}$ in (1.1) above is a global minimizer of the function $\Phi \colon C(\mathbb{R}^d, \mathbb{R}) \to [0, \infty)$. The optimization problem of approximately computing minima of the function $\Phi$ is not suitable for discrete numerical computations on a computer as the function $\Phi$ is defined on the infinite dimensional Banach space $C(\mathbb{R}^d, \mathbb{R})$. To overcome this we introduce a spatially discretized version of this optimization problem. More specifically, let $\mathfrak{d} \in \mathbb{N}$, let $\psi = (\psi_\theta)_{\theta \in \mathbb{R}^{\mathfrak{d}}} \colon \mathbb{R}^{\mathfrak{d}} \to C(\mathbb{R}^d, \mathbb{R})$ be a function, and let $\Psi \colon \mathbb{R}^{\mathfrak{d}} \to [0, \infty)$ satisfy $\Psi = \Phi \circ \psi$. We think of the set

$$\{\psi_\theta \colon \theta \in \mathbb{R}^{\mathfrak{d}}\} \subseteq C(\mathbb{R}^d, \mathbb{R}) \tag{1.3}$$

as a parametrized set of functions which we employ to approximate the infinite dimensional Banach space $C(\mathbb{R}^d, \mathbb{R})$ and we think of the function $\mathbb{R}^{\mathfrak{d}} \ni \theta \mapsto \psi_\theta \in C(\mathbb{R}^d, \mathbb{R})$ as the parametrization function corresponding to this set. Taking the set in (1.3) and its parametrization function $\mathbb{R}^{\mathfrak{d}} \ni \theta \mapsto \psi_\theta \in C(\mathbb{R}^d, \mathbb{R})$ into account, we then intend to approximately compute minima of the function $\Phi$ restricted to the set $\{\psi_\theta \colon \theta \in \mathbb{R}^{\mathfrak{d}}\}$, that is, we consider the optimization problem of approximately computing minima of the function

$$\{\psi_\theta \colon \theta \in \mathbb{R}^{\mathfrak{d}}\} \ni \phi \mapsto \Phi(\phi) = \left[\sum_{m=1}^{M} |\phi(x_m) - y_m|^2\right] \in [0, \infty). \tag{1.4}$$

Employing the parametrization function $\mathbb{R}^{\mathfrak{d}} \ni \theta \mapsto \psi_\theta \in C(\mathbb{R}^d, \mathbb{R})$ one can also reformulate this optimization problem as the optimization problem of approximately computing minima of the function

$$\mathbb{R}^{\mathfrak{d}} \ni \theta \mapsto \Psi(\theta) = \Phi(\psi_\theta) = \left[\sum_{m=1}^{M} |\psi_\theta(x_m) - y_m|^2\right] \in [0, \infty) \tag{1.5}$$

and this optimization is now accessible for discrete numerical computations. In the context of deep supervised learning algorithms, one would choose the parametrization function $\mathbb{R}^{\mathfrak{d}} \ni \theta \mapsto \psi_\theta \in C(\mathbb{R}^d, \mathbb{R})$ as deep neural network parametrizations and one would then apply a stochastic gradient descent optimization algorithm to the optimization problem in (1.5) to approximately compute minima of (1.5).

# Chapter 2

# Basics on artificial neural networks (ANNs)

In this chapter we present two approaches on how artificial neural networks (ANNs) can be described in a rigorous mathematical way.

## 2.1 Vectorized description of ANNs

### 2.1.1 Affine functions

**Definition 2.1.1** (Affine functions). *Let $\mathfrak{d}, m, n \in \mathbb{N}$, $s \in \mathbb{N}_0$, $\theta = (\theta_1, \theta_2, \ldots, \theta_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ satisfy $\mathfrak{d} \geq s + mn + m$. Then we denote by $\mathcal{A}_{m,n}^{\theta,s} \colon \mathbb{R}^n \to \mathbb{R}^m$ the function which satisfies for all $x = (x_1, x_2, \ldots, x_n) \in \mathbb{R}^n$ that*

$$
\begin{aligned}
\mathcal{A}_{m,n}^{\theta,s}(x) &= \begin{pmatrix} \theta_{s+1} & \theta_{s+2} & \cdots & \theta_{s+n} \\ \theta_{s+n+1} & \theta_{s+n+2} & \cdots & \theta_{s+2n} \\ \theta_{s+2n+1} & \theta_{s+2n+2} & \cdots & \theta_{s+3n} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_{s+(m-1)n+1} & \theta_{s+(m-1)n+2} & \cdots & \theta_{s+mn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} + \begin{pmatrix} \theta_{s+mn+1} \\ \theta_{s+mn+2} \\ \theta_{s+mn+3} \\ \vdots \\ \theta_{s+mn+m} \end{pmatrix} \\
&= \left( \left[ \textstyle\sum_{k=1}^n x_k \theta_{s+k} \right] + \theta_{s+mn+1}, \left[ \textstyle\sum_{k=1}^n x_k \theta_{s+n+k} \right] + \theta_{s+mn+2}, \right. \\
&\qquad \left. \ldots, \left[ \sum_{k=1}^n x_k \theta_{s+(m-1)n+k} \right] + \theta_{s+mn+m} \right)
\end{aligned}
\tag{2.1}
$$

*and we call $\mathcal{A}_{m,n}^{\theta,s}$ the affine function from $\mathbb{R}^n$ to $\mathbb{R}^m$ associated to $(\theta, s)$.*

### 2.1.2 Vectorized description of ANNs

**Definition 2.1.2** (Vectorized description of ANNs). *Let $\mathfrak{d}, L \in \mathbb{N}$, $l_0, l_1, \ldots, l_L \in \mathbb{N}$, $\theta \in \mathbb{R}^{\mathfrak{d}}$ satisfy*

$$
\mathfrak{d} \geq \sum_{k=1}^L l_k(l_{k-1} + 1)
\tag{2.2}
$$

*and let $\Psi_k \colon \mathbb{R}^{l_k} \to \mathbb{R}^{l_k}$, $k \in \{1, 2, \ldots, L\}$, be functions. Then we denote by $\mathcal{N}^{\theta,l_0}_{\Psi_1,\Psi_2,\ldots,\Psi_L} \colon \mathbb{R}^{l_0} \to \mathbb{R}^{l_L}$ the function which satisfies for all $x \in \mathbb{R}^{l_0}$ that*

$$
\left(\mathcal{N}^{\theta,l_0}_{\Psi_1,\Psi_2,\ldots,\Psi_L}\right)(x) = \left(\Psi_L \circ \mathcal{A}^{\theta,\sum_{k=1}^{L-1} l_k(l_{k-1}+1)}_{l_L,l_{L-1}} \circ \Psi_{L-1} \circ \mathcal{A}^{\theta,\sum_{k=1}^{L-2} l_k(l_{k-1}+1)}_{l_{L-1},l_{L-2}} \circ \ldots \right.
$$
$$
\left. \ldots \circ \Psi_2 \circ \mathcal{A}^{\theta,l_1(l_0+1)}_{l_2,l_1} \circ \Psi_1 \circ \mathcal{A}^{\theta,0}_{l_1,l_0}\right)(x) \quad (2.3)
$$

*(cf. Definition 2.1.1) and we call $\mathcal{N}^{\theta,l_0}_{\Psi_1,\Psi_2,\ldots,\Psi_L}$ the realization of the fully connected feedforward artificial neural network associated to $\theta$ with $L+1$ layers and with dimensions $(l_0, l_1, \ldots, l_L)$ and activation functions $(\Psi_1, \Psi_2, \ldots, \Psi_L)$.*

### 2.1.3 Weights and biases of ANNs

**Remark 2.1.3.** *Let $L \in \{2, 3, \ldots\}$, $v_0, v_1, \ldots, v_{L-1} \in \mathbb{N}_0$, $l_0, l_1, \ldots, l_L$, $\mathfrak{d} \in \mathbb{N}$, $\theta = (\theta_1, \theta_2, \ldots, \theta_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ satisfy for all $k \in \{0, 1, \ldots, L-1\}$ that*

$$
\mathfrak{d} \geq \sum_{i=1}^{L} l_i(l_{i-1} + 1) \qquad and \qquad v_k = \sum_{i=1}^{k} l_i(l_{i-1} + 1), \qquad (2.4)
$$

*let $W_k \in \mathbb{R}^{l_k \times l_{k-1}}$, $k \in \{1, 2, \ldots, L\}$, and $b_k \in \mathbb{R}^{l_k}$, $k \in \{1, 2, \ldots, L\}$, satisfy for all $k \in \{1, 2, \ldots, L\}$ that*

$$
W_k = \underbrace{\begin{pmatrix} \theta_{v_{k-1}+1} & \theta_{v_{k-1}+2} & \cdots & \theta_{v_{k-1}+l_{k-1}} \\ \theta_{v_{k-1}+l_{k-1}+1} & \theta_{v_{k-1}+l_{k-1}+2} & \cdots & \theta_{v_{k-1}+2l_{k-1}} \\ \theta_{v_{k-1}+2l_{k-1}+1} & \theta_{v_{k-1}+2l_{k-1}+2} & \cdots & \theta_{v_{k-1}+3l_{k-1}} \\ \vdots & \vdots & \vdots & \vdots \\ \theta_{v_{k-1}+(l_k-1)l_{k-1}+1} & \theta_{v_{k-1}+(l_k-1)l_{k-1}+2} & \cdots & \theta_{v_{k-1}+l_k l_{k-1}} \end{pmatrix}}_{weights} \qquad (2.5)
$$

*and*

$$
b_k = \underbrace{\left(\theta_{v_{k-1}+l_k l_{k-1}+1}, \theta_{v_{k-1}+l_k l_{k-1}+2}, \ldots, \theta_{v_{k-1}+l_k l_{k-1}+l_k}\right)}_{biases}, \qquad (2.6)
$$

*and let $\Psi_k \colon \mathbb{R}^{l_k} \to \mathbb{R}^{l_k}$, $k \in \{1, 2, \ldots, L\}$, be functions. Then*

*(i) it holds that*

$$
\mathcal{N}^{\theta,l_0}_{\Psi_1,\Psi_2,\ldots,\Psi_L} = \Psi_L \circ \mathcal{A}^{\theta,v_{L-1}}_{l_L,l_{L-1}} \circ \Psi_{L-1} \circ \mathcal{A}^{\theta,v_{L-2}}_{l_{L-1},l_{L-2}} \circ \Psi_{L-2} \circ \ldots \circ \mathcal{A}^{\theta,v_1}_{l_2,l_1} \circ \Psi_1 \circ \mathcal{A}^{\theta,v_0}_{l_1,l_0} \quad (2.7)
$$

*and*

*(ii) it holds for all $k \in \{1, 2, \ldots, L\}$, $x \in \mathbb{R}^{l_{k-1}}$ that $\mathcal{A}^{\theta,v_{k-1}}_{l_k,l_{k-1}}(x) = W_k x + b_k$*

*(cf. Definitions 2.1.1 and 2.1.2).*

### 2.1.4 Activation functions

#### 2.1.4.1 Multidimensional versions

To describe multidimensional activation functions, we frequently employ the concept of the multidimensional version of a function. This concept is the subject of the next notion.

**Definition 2.1.4** (Multidimensional versions). *Let $d \in \mathbb{N}$ and let $\psi \colon \mathbb{R} \to \mathbb{R}$ be a function. Then we denote by $\mathfrak{M}_{\psi,d} \colon \mathbb{R}^d \to \mathbb{R}^d$ the function which satisfies for all $x = (x_1, x_2, \ldots, x_d) \in \mathbb{R}^d$ that*

$$\mathfrak{M}_{\psi,d}(x) = (\psi(x_1), \psi(x_2), \ldots, \psi(x_d)). \tag{2.8}$$

*and we call $\mathfrak{M}_{\psi,d}$ the $d$-dimensional version of $\psi$.*

#### 2.1.4.2 Single hidden layer artificial neural networks

**Example 2.1.5.** *Let $\mathcal{I}, \mathcal{H} \in \mathbb{N}$, $\theta = (\theta_1, \theta_2, \ldots, \theta_{\mathcal{H}\mathcal{I}+2\mathcal{H}+1}) \in \mathbb{R}^{\mathcal{H}\mathcal{I}+2\mathcal{H}+1}$, $x = (x_1, x_2, \ldots, x_{\mathcal{I}}) \in \mathbb{R}^{\mathcal{I}}$ and let $\psi \colon \mathbb{R} \to \mathbb{R}$ be a function. Then*

$$
\begin{aligned}
\mathcal{N}_{\mathfrak{M}_{\psi,\mathcal{H}}, \mathrm{id}_{\mathbb{R}}}^{\theta, \mathcal{I}}(x) &= \left( (\mathrm{id}_{\mathbb{R}}) \circ \mathcal{A}_{1,\mathcal{H}}^{\theta, \mathcal{H}\mathcal{I}+\mathcal{H}} \circ \mathfrak{M}_{\psi,\mathcal{H}} \circ \mathcal{A}_{\mathcal{H},\mathcal{I}}^{\theta, 0} \right)(x) \\
&= \mathcal{A}_{1,\mathcal{H}}^{\theta, \mathcal{H}\mathcal{I}+\mathcal{H}} \big( \mathfrak{M}_{\psi,\mathcal{H}} \big( \mathcal{A}_{\mathcal{H},\mathcal{I}}^{\theta, 0}(x) \big) \big) \\
&= \left[ \sum_{k=1}^{\mathcal{H}} \theta_{\mathcal{H}\mathcal{I}+\mathcal{H}+k} \, \psi \left( \left[ \sum_{i=1}^{\mathcal{I}} x_i \theta_{(k-1)\mathcal{I}+i} \right] + \theta_{\mathcal{H}\mathcal{I}+k} \right) \right] + \theta_{\mathcal{H}\mathcal{I}+2\mathcal{H}+1}.
\end{aligned}
\tag{2.9}
$$

*(cf. Definitions 2.1.1, 2.1.2, and 2.1.4).*

#### 2.1.4.3 The rectifier function

In this subsection we formulate the rectifier function which is maybe the most commonly used activation function in deep learning applications (cf., for example, Le Cun, Bengio, & Hinton [21]).

**Definition 2.1.6** (Rectifier function). *We denote by $\mathfrak{r} \colon \mathbb{R} \to \mathbb{R}$ the function which satisfies for all $x \in \mathbb{R}$ that*

$$\mathfrak{r}(x) = \max\{x, 0\}. \tag{2.10}$$

*and we call $\mathfrak{r}$ the rectifier function.*

**Definition 2.1.7** (Multidimensional rectifier functions). *Let $d \in \mathbb{N}$. Then we denote by $\mathfrak{R}_d \colon \mathbb{R}^d \to \mathbb{R}^d$ the function given by*

$$\mathfrak{R}_d = \mathfrak{M}_{\mathfrak{r},d} \tag{2.11}$$

*(cf. Definitions 2.1.4 and 2.1.6) and we call $\mathfrak{R}_d$ the $d$-dimensional rectifier function.*

**Proposition 2.1.8** (An artificial neural network with the rectifier function as the activation function). *Let $W_1 = w_1 = 1$, $W_2 = w_2 = -1$, $b_1 = b_2 = B = 0$. Then it holds for all $x \in \mathbb{R}$ that*

$$x = W_1 \max\{w_1 x + b_1, 0\} + W_2 \max\{w_2 x + b_2, 0\} + B. \tag{2.12}$$

*Proof of Proposition 2.1.8.* Observe that for all $x \in \mathbb{R}$ it holds that

$$
\begin{aligned}
&W_1 \max\{w_1 x + b_1, 0\} + W_2 \max\{w_2 x + b_2, 0\} + B \\
&= \max\{w_1 x + b_1, 0\} - \max\{w_2 x + b_2, 0\} = \max\{x, 0\} - \max\{-x, 0\} \\
&= \max\{x, 0\} + \min\{x, 0\} = x.
\end{aligned}
\tag{2.13}
$$

The proof of Proposition 2.1.8 is thus complete. $\qquad\square$

**Exercise 2.1.1** (Real identity). *Prove or disprove the following statement: There exist* $\mathfrak{d}, L \in \mathbb{N}$, $l_1, l_2, \ldots, l_L \in \mathbb{N}$, $\theta = (\theta_1, \theta_2, \ldots, \theta_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ *with* $\mathfrak{d} \geq 2l_1 + \left[\sum_{k=2}^{L} l_k(l_{k-1} + 1)\right] + l_L + 1$ *such that for all* $x \in \mathbb{R}$ *it holds that*

$$\left(\mathcal{N}^{\theta,1}_{\mathfrak{R}_{l_1}, \mathfrak{R}_{l_2}, \ldots, \mathfrak{R}_{l_L}, \mathrm{id}_{\mathbb{R}}}\right)(x) = x \tag{2.14}$$

*(cf. Definitions 2.1.2 and 2.1.7).*

The statement of the next lemma, Lemma 2.1.9, provides a partial answer to Exercise 2.1.1. Lemma 2.1.9 follows from an application of Proposition 2.1.8 and the detailed proof of Lemma 2.1.9 is left as an exercise.

**Lemma 2.1.9** (Real identity). *Let* $\theta = (1, -1, 0, 0, 1, -1, 0) \in \mathbb{R}^7$. *Then it holds for all* $x \in \mathbb{R}$ *that*

$$\left(\mathcal{N}^{\theta,1}_{\mathfrak{R}_2, \mathrm{id}_{\mathbb{R}}}\right)(x) = x \tag{2.15}$$

*(cf. Definitions 2.1.2 and 2.1.7).*

**Exercise 2.1.2** (Absolute value). *Prove or disprove the following statement: There exist* $\mathfrak{d}, L \in \mathbb{N}$, $l_1, l_2, \ldots, l_L \in \mathbb{N}$, $\theta = (\theta_1, \theta_2, \ldots, \theta_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ *with* $\mathfrak{d} \geq 2l_1 + \left[\sum_{k=2}^{L} l_k(l_{k-1} + 1)\right] + l_L + 1$ *such that for all* $x \in \mathbb{R}$ *it holds that*

$$\left(\mathcal{N}^{\theta,1}_{\mathfrak{R}_{l_1}, \mathfrak{R}_{l_2}, \ldots, \mathfrak{R}_{l_L}, \mathrm{id}_{\mathbb{R}}}\right)(x) = |x| \tag{2.16}$$

*(cf. Definitions 2.1.2 and 2.1.7).*

**Exercise 2.1.3** (Exponential). *Prove or disprove the following statement: There exist* $\mathfrak{d}, L \in \mathbb{N}$, $l_1, l_2, \ldots, l_L \in \mathbb{N}$, $\theta = (\theta_1, \theta_2, \ldots, \theta_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ *with* $\mathfrak{d} \geq 2l_1 + \left[\sum_{k=2}^{L} l_k(l_{k-1} + 1)\right] + l_L + 1$ *such that for all* $x \in \mathbb{R}$ *it holds that*

$$\left(\mathcal{N}^{\theta,1}_{\mathfrak{R}_{l_1}, \mathfrak{R}_{l_2}, \ldots, \mathfrak{R}_{l_L}, \mathrm{id}_{\mathbb{R}}}\right)(x) = e^x \tag{2.17}$$

*(cf. Definitions 2.1.2 and 2.1.7).*

**Exercise 2.1.4** (Two-dimensional maximum). *Prove or disprove the following statement: There exist* $\mathfrak{d}, L \in \mathbb{N}$, $l_1, l_2, \ldots, l_L \in \mathbb{N}$, $\theta = (\theta_1, \theta_2, \ldots, \theta_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ *with* $\mathfrak{d} \geq 3l_1 + \left[\sum_{k=2}^{L} l_k(l_{k-1} + 1)\right] + l_L + 1$ *such that for all* $x, y \in \mathbb{R}$ *it holds that*

$$\left(\mathcal{N}^{\theta,2}_{\mathfrak{R}_{l_1}, \mathfrak{R}_{l_2}, \ldots, \mathfrak{R}_{l_L}, \mathrm{id}_{\mathbb{R}}}\right)(x, y) = \max\{x, y\} \tag{2.18}$$

*(cf. Definitions 2.1.2 and 2.1.7).*

**Exercise 2.1.5** (Real identity with two hidden layers). *Prove or disprove the following statement: There exist* $\mathfrak{d}, l_1, l_2 \in \mathbb{N}$, $\theta = (\theta_1, \theta_2, \ldots, \theta_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ *with* $\mathfrak{d} \geq 2l_1 + l_1 l_2 + 2l_2 + 1$ *such that for all* $x \in \mathbb{R}$ *it holds that*

$$\left(\mathcal{N}^{\theta,1}_{\mathfrak{R}_{l_1}, \mathfrak{R}_{l_2}, \mathrm{id}_{\mathbb{R}}}\right)(x) = x \tag{2.19}$$

*(cf. Definitions 2.1.2 and 2.1.7).*

The statement of the next lemma, Lemma 2.1.10, provides a partial answer to Exercise 2.1.5. The proof of Lemma 2.1.10 is left as an exercise.

**Lemma 2.1.10** (Real identity with two hidden layers)**.** *Let* $\theta = (1, -1, 0, 0, 1, -1, -1,$ $1, 0, 0, 1, -1, 0) \in \mathbb{R}^{13}$*. Then it holds for all* $x \in \mathbb{R}$ *that*

$$\left(\mathcal{N}^{\theta,1}_{\mathfrak{R}_2,\mathfrak{R}_2,\mathrm{id}_{\mathbb{R}}}\right)(x) = x \tag{2.20}$$

*(cf. Definitions 2.1.2 and 2.1.7).*

**Exercise 2.1.6** (Three-dimensional maximum)**.** *Prove or disprove the following statement: There exist* $\mathfrak{d}, L \in \mathbb{N}$*,* $l_1, l_2, \ldots, l_L \in \mathbb{N}$*,* $\theta = (\theta_1, \theta_2, \ldots, \theta_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ *with* $\mathfrak{d} \geq$ $4l_1 + \left[\sum_{k=2}^{L} l_k(l_{k-1} + 1)\right] + l_L + 1$ *such that for all* $x, y, z \in \mathbb{R}$ *it holds that*

$$\left(\mathcal{N}^{\theta,3}_{\mathfrak{R}_{l_1},\mathfrak{R}_{l_2},\ldots,\mathfrak{R}_{l_L},\mathrm{id}_{\mathbb{R}}}\right)(x,y,z) = \max\{x,y,z\} \tag{2.21}$$

*(cf. Definition 2.1.2 and Definition 2.1.7).*

**Exercise 2.1.7** (Multidimensional maxima)**.** *Prove or disprove the following statement: For every* $k \in \mathbb{N}$ *there exist* $\mathfrak{d}, L \in \mathbb{N}$*,* $l_1, l_2, \ldots, l_L \in \mathbb{N}$*,* $\theta = (\theta_1, \theta_2, \ldots, \theta_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ *with* $\mathfrak{d} \geq (k+1)l_1 + \left[\sum_{k=2}^{L} l_k(l_{k-1} + 1)\right] + l_L + 1$ *such that for all* $x_1, x_2, \ldots, x_k \in \mathbb{R}$ *it holds that*

$$\left(\mathcal{N}^{\theta,k}_{\mathfrak{R}_{l_1},\mathfrak{R}_{l_2},\ldots,\mathfrak{R}_{l_L},\mathrm{id}_{\mathbb{R}}}\right)(x_1,x_2,\ldots,x_k) = \max\{x_1,x_2,\ldots,x_k\} \tag{2.22}$$

*(cf. Definitions 2.1.2 and 2.1.7).*

**Exercise 2.1.8** (Hat function)**.** *Prove or disprove the following statement: There exist* $\mathfrak{d}, l \in \mathbb{N}$*,* $\theta = (\theta_1, \theta_2, \ldots, \theta_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ *with* $\mathfrak{d} \geq 3l + 1$ *such that for all* $x \in \mathbb{R}$ *it holds that*

$$\left(\mathcal{N}^{\theta,1}_{\mathfrak{R}_l,\mathrm{id}_{\mathbb{R}}}\right)(x) = \begin{cases} 1 & : x \leq 2 \\ x-1 & : 2 < x \leq 3 \\ 5-x & : 3 < x \leq 4 \\ 1 & : x > 4 \end{cases} \tag{2.23}$$

*(cf. Definition 2.1.2 and Definition 2.1.7).*

**Exercise 2.1.9.** *Prove or disprove the following statement: There exist* $\mathfrak{d}, l \in \mathbb{N}$*,* $\theta = (\theta_1, \theta_2, \ldots, \theta_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ *with* $\mathfrak{d} \geq 3l + 1$ *such that for all* $x \in \mathbb{R}$ *it holds that*

$$\left(\mathcal{N}^{\theta,1}_{\mathfrak{R}_l,\mathrm{id}_{\mathbb{R}}}\right)(x) = \begin{cases} -2 & : x \leq 1 \\ 2x-4 & : 1 < x \leq 3 \\ 2 & : x > 3 \end{cases} \tag{2.24}$$

*(cf. Definition 2.1.2 and Definition 2.1.7).*

**Exercise 2.1.10.** *Prove or disprove the following statement: There exist* $\mathfrak{d}, l \in \mathbb{N}$*,* $\theta = (\theta_1, \theta_2, \ldots, \theta_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ *with* $\mathfrak{d} \geq 3l + 1$ *such that for all* $x \in [0,1]$ *it holds that*

$$\left(\mathcal{N}^{\theta,1}_{\mathfrak{R}_l,\mathrm{id}_{\mathbb{R}}}\right)(x) = x^2 \tag{2.25}$$

*(cf. Definition 2.1.2 and Definition 2.1.7).*

**2.1.4.4 Clipping functions**

**Definition 2.1.11** (Clipping function)**.** *Let $u \in [-\infty, \infty)$, $v \in (u, \infty]$. Then we denote by $\mathfrak{c}_{u,v} \colon \mathbb{R} \to \mathbb{R}$ the function which satisfies for all $x \in \mathbb{R}$ that*

$$\mathfrak{c}_{u,v}(x) = \max\{u, \min\{x, v\}\}. \tag{2.26}$$

*and we call $\mathfrak{c}_{u,v}$ the $(u, v)$-clipping function.*

**Definition 2.1.12** (Multidimensional clipping functions)**.** *Let $d \in \mathbb{N}$, $u \in [-\infty, \infty)$, $v \in (u, \infty]$. Then we denote by $\mathfrak{C}_{u,v,d} \colon \mathbb{R}^d \to \mathbb{R}^d$ the function given by*

$$\mathfrak{C}_{u,v,d} = \mathfrak{M}_{\mathfrak{c}_{u,v},d} \tag{2.27}$$

*(cf. Definitions 2.1.4 and 2.1.11) and we call $\mathfrak{C}_{u,v,d}$ the $d$-dimensional $(u, v)$-clipping function.*

**2.1.4.5 The softplus function**

**Definition 2.1.13** (Softplus function)**.** *We denote by $\mathfrak{s} \colon \mathbb{R} \to \mathbb{R}$ the function which satisfies for all $x \in \mathbb{R}$ that*

$$\mathfrak{s}(x) = \ln(1 + \exp(x)) \tag{2.28}$$

*and we call $\mathfrak{s}$ the softplus function.*

The next result, Lemma 2.1.14 below, presents a few elementary properties of the softplus function.

**Lemma 2.1.14** (Properties of the softplus function)**.** *It holds*

  *(i) for all $x \in [0, \infty)$ that $x \le \mathfrak{s}(x) \le x + 1$,*

  *(ii) that $\lim_{x \to -\infty} \mathfrak{s}(x) = 0$,*

  *(iii) that $\lim_{x \to \infty} \mathfrak{s}(x) = \infty$, and*

  *(iv) that $\mathfrak{s}(0) = \ln(2)$*

*(cf. Definition 2.1.13).*

*Proof of Lemma 2.1.14.* Observe that the fact that $2 \le \exp(1)$ ensures that for all $x \in [0, \infty)$ it holds that

$$\begin{aligned} x = \ln(\exp(x)) &\le \ln(1 + \exp(x)) = \ln(\exp(0) + \exp(x)) \\ &\le \ln(\exp(x) + \exp(x)) = \ln(2\exp(x)) \le \ln(\exp(1)\exp(x)) \\ &= \ln(\exp(x + 1)) = x + 1. \end{aligned} \tag{2.29}$$

The proof of Lemma 2.1.14 is thus complete. □

Note that Lemma 2.1.14 ensures that $\mathfrak{s}(0) = \ln(2) = 0.693\ldots$ (cf. Definition 2.1.13). In the next step we introduce the multidimensional version of the softplus function (cf. Definitions 2.1.4 and 2.1.13 above).

**Definition 2.1.15** (Multidimensional softplus functions)**.** *Let $d \in \mathbb{N}$. Then we denote by $\mathfrak{S}_d \colon \mathbb{R}^d \to \mathbb{R}^d$ the function given by*

$$\mathfrak{S}_d = \mathfrak{M}_{\mathfrak{s},d} \tag{2.30}$$

*(cf. Definitions 2.1.4 and 2.1.13) and we call $\mathfrak{S}_d$ the $d$-dimensional softplus function.*

**2.1.4.6 The standard logistic function**

**Definition 2.1.16** (Standard logistic function). *We denote by* $\mathfrak{l}\colon \mathbb{R} \to \mathbb{R}$ *the function which satisfies for all* $x \in \mathbb{R}$ *that*

$$\mathfrak{l}(x) = \frac{1}{1+\exp(-x)} = \frac{\exp(x)}{\exp(x)+1} \tag{2.31}$$

*and we call* $\mathfrak{l}$ *the standard logistic function.*

**Definition 2.1.17** (Multidimensional standard logistic functions). *Let* $d \in \mathbb{N}$*. Then we denote by* $\mathfrak{L}_d\colon \mathbb{R}^d \to \mathbb{R}^d$ *the function given by*

$$\mathfrak{L}_d = \mathfrak{M}_{\mathfrak{l},d} \tag{2.32}$$

*(cf. Definitions 2.1.4 and 2.1.16) and we call* $\mathfrak{L}_d$ *the d-dimensional standard logistic function.*

**2.1.4.7 Derivative of the standard logistic function**

**Proposition 2.1.18** (Logistic differential equation). *It holds that* $\mathfrak{l}\colon \mathbb{R} \to \mathbb{R}$ *is infinitely often differentiable and it holds for all* $x \in \mathbb{R}$ *that*

$$\mathfrak{l}(0) = {}^1\!/{}_2, \qquad \mathfrak{l}'(x) = \mathfrak{l}(x)(1-\mathfrak{l}(x)) = \mathfrak{l}(x) - [\mathfrak{l}(x)]^2, \qquad and \tag{2.33}$$

$$\mathfrak{l}''(x) = \mathfrak{l}(x)(1-\mathfrak{l}(x))(1-2\,\mathfrak{l}(x)) = 2[\mathfrak{l}(x)]^3 - 3[\mathfrak{l}(x)]^2 + \mathfrak{l}(x) \tag{2.34}$$

*(cf. Definition 2.1.16).*

*Proof of Proposition 2.1.18.* Observe that (2.31) ensures that for all $x \in \mathbb{R}$ it holds that

$$
\begin{aligned}
\mathfrak{l}'(x) &= \frac{\exp(-x)}{(1+\exp(-x))^2} = \mathfrak{l}(x)\left(\frac{\exp(-x)}{1+\exp(-x)}\right) \\
&= \mathfrak{l}(x)\left(\frac{1+\exp(-x)-1}{1+\exp(-x)}\right) = \mathfrak{l}(x)\left(1 - \frac{1}{1+\exp(-x)}\right) \\
&= \mathfrak{l}(x)(1-\mathfrak{l}(x)).
\end{aligned}
\tag{2.35}
$$

Hence, we obtain that for all $x \in \mathbb{R}$ it holds that

$$
\begin{aligned}
\mathfrak{l}''(x) &= \big[\mathfrak{l}(x)(1-\mathfrak{l}(x))\big]' = \mathfrak{l}'(x)(1-\mathfrak{l}(x)) + \mathfrak{l}(x)(1-\mathfrak{l}(x))' \\
&= \mathfrak{l}'(x)(1-\mathfrak{l}(x)) - \mathfrak{l}(x)\,\mathfrak{l}'(x) = \mathfrak{l}'(x)(1-2\,\mathfrak{l}(x)) \\
&= \mathfrak{l}(x)(1-\mathfrak{l}(x))(1-2\,\mathfrak{l}(x)) \\
&= \big(\mathfrak{l}(x) - [\mathfrak{l}(x)]^2\big)(1-2\,\mathfrak{l}(x)) = \mathfrak{l}(x) - [\mathfrak{l}(x)]^2 - 2[\mathfrak{l}(x)]^2 + 2[\mathfrak{l}(x)]^3 \\
&= 2[\mathfrak{l}(x)]^3 - 3[\mathfrak{l}(x)]^2 + \mathfrak{l}(x).
\end{aligned}
\tag{2.36}
$$

The proof of Proposition 2.1.18 is thus complete. □

### 2.1.4.8   Integral of the standard logistic function

**Lemma 2.1.19** (Primitive of the standard logistic function)**.** *It holds for all $x \in \mathbb{R}$ that*

$$\int_{-\infty}^{x} \mathfrak{l}(y)\, \mathrm{d}y = \int_{-\infty}^{x} \left( \frac{1}{1 + e^{-y}} \right) \mathrm{d}y = \ln(1 + \exp(x)) = \mathfrak{s}(x) \tag{2.37}$$

*(cf. Definitions 2.1.13 and 2.1.16).*

*Proof of Lemma 2.1.19.* Observe that (2.28) implies that for all $x \in \mathbb{R}$ it holds that

$$\mathfrak{s}'(x) = \left[ \frac{1}{1 + \exp(x)} \right] \exp(x) = \mathfrak{l}(x). \tag{2.38}$$

The fundamental theorem of calculus hence shows that for all $w, x \in \mathbb{R}$ with $w \leq x$ it holds that

$$\int_{w}^{x} \underbrace{\mathfrak{l}(y)}_{\geq 0}\, \mathrm{d}y = \mathfrak{s}(x) - \mathfrak{s}(w). \tag{2.39}$$

Combining this with the fact that $\lim_{w \to -\infty} \mathfrak{s}(w) = 0$ establishes (2.37). The proof of Lemma 2.1.19 is thus complete. $\qquad\square$

### 2.1.4.9   The hyperbolic tangent function

**Definition 2.1.20** (Hyperbolic tangent)**.** *We denote by* $\tanh \colon \mathbb{R} \to \mathbb{R}$ *the function which satisfies for all $x \in \mathbb{R}$ that*

$$\tanh(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)} \tag{2.40}$$

*and we call* $\tanh$ *the hyperbolic tangent.*

**Definition 2.1.21** (Multidimensional hyperbolic tangent functions)**.** *Let $d \in \mathbb{N}$. Then we denote by $\mathfrak{T}_d \colon \mathbb{R}^d \to \mathbb{R}^d$ the function given by*

$$\mathfrak{T}_d = \mathfrak{M}_{\tanh, d} \tag{2.41}$$

*(cf. Definitions 2.1.4 and 2.1.20) and we call $\mathfrak{T}_d$ the d-dimensional hyperbolic tangent.*

**Lemma 2.1.22.** *It holds for all $x \in \mathbb{R}$ that*

$$\tanh(x) = 2\, \mathfrak{l}(2x) - 1 \tag{2.42}$$

*(cf. Definitions 2.1.16 and 2.1.20).*

*Proof of Lemma 2.1.22.* Observe that (2.31) and (2.40) ensure that for all $x \in \mathbb{R}$ it holds that

$$\begin{aligned}
2\, \mathfrak{l}(2x) - 1 &= 2 \left( \frac{\exp(2x)}{\exp(2x) + 1} \right) - 1 = \frac{2 \exp(2x) - (\exp(2x) + 1)}{\exp(2x) + 1} \\
&= \frac{\exp(2x) - 1}{\exp(2x) + 1} = \frac{\exp(x)(\exp(x) - \exp(-x))}{\exp(x)(\exp(x) + \exp(-x))} \\
&= \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)} = \tanh(x).
\end{aligned} \tag{2.43}$$

The proof of Lemma 2.1.22 is thus complete. $\qquad\square$

### 2.1.4.10 The Heaviside function

**Definition 2.1.23** (Heaviside function). *We denote by $\mathfrak{h}\colon \mathbb{R} \to \mathbb{R}$ the function which satisfies for all $x \in \mathbb{R}$ that*

$$\mathfrak{h}(x) = \mathbb{1}_{[0,\infty)}(x) = \begin{cases} 1 & : x \geq 0 \\ 0 & : x < 0 \end{cases} \tag{2.44}$$

*and we call $\mathfrak{h}$ the Heaviside function (we call $\mathfrak{h}$ the Heaviside step function, we call $\mathfrak{h}$ the unit step function).*

**Definition 2.1.24** (Multidimensional Heaviside functions). *Let $d \in \mathbb{N}$. Then we denote by $\mathfrak{H}_d\colon \mathbb{R}^d \to \mathbb{R}^d$ the function given by*

$$\mathfrak{H}_d = \mathfrak{M}_{\mathfrak{h},d} \tag{2.45}$$

*(cf. Definitions 2.1.4 and 2.1.23) and we call $\mathfrak{H}_d$ the d-dimensional Heaviside function (we call $\mathfrak{H}_d$ the d-dimensional Heaviside step function, we call $\mathfrak{H}_d$ the d-dimensional unit step function).*

### 2.1.4.11 The softmax function

**Definition 2.1.25** (The softmax function). *Let $d \in \mathbb{N}$. Then we denote by $\mathscr{S}_d = (\mathscr{S}_d^{(1)}, \mathscr{S}_d^{(2)}, \ldots, \mathscr{S}_d^{(d)})\colon \mathbb{R}^d \to \mathbb{R}^d$ the function which satisfies for all $x = (x_1, x_2, \ldots, x_d) \in \mathbb{R}^d$ that*

$$\begin{aligned} \mathscr{S}_d(x) &= \left( \mathscr{S}_d^{(1)}(x), \mathscr{S}_d^{(2)}(x), \ldots, \mathscr{S}_d^{(d)}(x) \right) \\ &= \left( \frac{\exp(x_1)}{\left(\sum_{i=1}^d \exp(x_i)\right)}, \frac{\exp(x_2)}{\left(\sum_{i=1}^d \exp(x_i)\right)}, \ldots, \frac{\exp(x_d)}{\left(\sum_{i=1}^d \exp(x_i)\right)} \right) \end{aligned} \tag{2.46}$$

*and we call $\mathscr{S}_d$ the d-dimensional softmax function.*

**Lemma 2.1.26.** *Let $d \in \mathbb{N}$. Then*

(i) *it holds for all $x \in \mathbb{R}^d$, $k \in \{1, 2, \ldots, d\}$ that $\mathscr{S}_d^{(k)}(x) \in (0, 1]$ and*

(ii) *it holds for all $x \in \mathbb{R}^d$ that*

$$\sum_{k=1}^d \mathscr{S}_d^{(k)}(x) = 1 \tag{2.47}$$

*(cf. Definition 2.1.25).*

*Proof of Lemma 2.1.26.* Observe that (2.46) demonstrates that for all $x = (x_1, x_2, \ldots, x_d) \in \mathbb{R}^d$ it holds that

$$\sum_{k=1}^d \mathscr{S}_d^{(k)}(x) = \sum_{k=1}^d \frac{\exp(x_k)}{\left(\sum_{i=1}^d \exp(x_i)\right)} = \frac{\sum_{k=1}^d \exp(x_k)}{\sum_{i=1}^d \exp(x_i)} = 1. \tag{2.48}$$

The proof of Lemma 2.1.26 is thus complete. $\qquad\square$

### 2.1.5 Rectified clipped ANNs

**Definition 2.1.27** (Rectified clipped ANNs). *Let $L, \mathfrak{d} \in \mathbb{N}$, $u \in [-\infty, \infty)$, $v \in (u, \infty]$, $\mathbf{l} = (l_0, l_1, \ldots, l_L) \in \mathbb{N}^{L+1}$, $\theta \in \mathbb{R}^{\mathfrak{d}}$ satisfy*

$$\mathfrak{d} \geq \sum_{k=1}^{L} l_k (l_{k-1} + 1). \tag{2.49}$$

*Then we denote by $\mathscr{N}_{u,v}^{\theta,\mathbf{l}} \colon \mathbb{R}^{l_0} \to \mathbb{R}^{l_L}$ the function which satisfies for all $x \in \mathbb{R}^{l_0}$ that*

$$\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(x) = \begin{cases} \left( \mathscr{N}_{\mathfrak{C}_{u,v,l_L}}^{\theta,l_0} \right)(x) & : L = 1 \\ \left( \mathscr{N}_{\mathfrak{R}_{l_1}, \mathfrak{R}_{l_2}, \ldots, \mathfrak{R}_{l_{L-1}}, \mathfrak{C}_{u,v,l_L}}^{\theta,l_0} \right)(x) & : L > 1 \end{cases} \tag{2.50}$$

*(cf. Definitions 2.1.2, 2.1.7, and 2.1.12).*

## 2.2 Structured description of ANNs

### 2.2.1 Structured description of ANNs

**Definition 2.2.1** (Structured description of ANNs). *We denote by $\mathbf{N}$ the set given by*

$$\mathbf{N} = \bigcup_{L \in \mathbb{N}} \bigcup_{l_0, l_1, \ldots, l_L \in \mathbb{N}} \left( \bigtimes_{k=1}^{L} (\mathbb{R}^{l_k \times l_{k-1}} \times \mathbb{R}^{l_k}) \right), \tag{2.51}$$

*we denote by $\mathcal{P} \colon \mathbf{N} \to \mathbb{N}$, $\mathcal{L} \colon \mathbf{N} \to \mathbb{N}$, $\mathcal{I} \colon \mathbf{N} \to \mathbb{N}$, $\mathcal{O} \colon \mathbf{N} \to \mathbb{N}$, $\mathcal{H} \colon \mathbf{N} \to \mathbb{N}_0$, $\mathcal{D} \colon \mathbf{N} \to \left( \bigcup_{L=2}^{\infty} \mathbb{N}^L \right)$, and $\mathbb{D}_n \colon \mathbf{N} \to \mathbb{N}_0$, $n \in \mathbb{N}_0$, the functions which satisfy for all $L \in \mathbb{N}$, $l_0, l_1, \ldots, l_L \in \mathbb{N}$, $\Phi \in \left( \bigtimes_{k=1}^{L} (\mathbb{R}^{l_k \times l_{k-1}} \times \mathbb{R}^{l_k}) \right)$, $n \in \mathbb{N}_0$ that $\mathcal{P}(\Phi) = \sum_{k=1}^{L} l_k (l_{k-1} + 1)$, $\mathcal{L}(\Phi) = L$, $\mathcal{I}(\Phi) = l_0$, $\mathcal{O}(\Phi) = l_L$, $\mathcal{H}(\Phi) = L - 1$, $\mathcal{D}(\Phi) = (l_0, l_1, \ldots, l_L)$, and*

$$\mathbb{D}_n(\Phi) = \begin{cases} l_n & : n \leq L \\ 0 & : n > L, \end{cases} \tag{2.52}$$

*and for every $L \in \mathbb{N}$, $l_0, l_1, \ldots, l_L \in \mathbb{N}$, $\Phi = ((W_1, B_1), (W_2, B_2), \ldots, (W_L, B_L)) \in \left( \bigtimes_{k=1}^{L} (\mathbb{R}^{l_k \times l_{k-1}} \times \mathbb{R}^{l_k}) \right)$ we denote by $\mathcal{W}_{(\cdot), \Phi} = (\mathcal{W}_{n, \Phi})_{n \in \{1, 2, \ldots, L\}} \colon \{1, 2, \ldots, L\} \to \left( \bigcup_{m, k \in \mathbb{N}} \mathbb{R}^{m \times k} \right)$ and $\mathcal{B}_{(\cdot), \Phi} = (\mathcal{B}_{n, \Phi})_{n \in \{1, 2, \ldots, L\}} \colon \{1, 2, \ldots, L\} \to \left( \bigcup_{m \in \mathbb{N}} \mathbb{R}^m \right)$ the functions which satisfy for all $n \in \{1, 2, \ldots, L\}$ that $\mathcal{W}_{n, \Phi} = W_n$ and $\mathcal{B}_{n, \Phi} = B_n$.*

**Definition 2.2.2.** *We say that $\Phi$ is a neural network if and only if it holds that $\Phi \in \mathbf{N}$.*

### 2.2.2 Realizations of ANNs

**Definition 2.2.3** (Realization associated to an ANN). *Let $a \in C(\mathbb{R}, \mathbb{R})$. Then we denote by $\mathcal{R}_a \colon \mathbf{N} \to \left( \bigcup_{k, l \in \mathbb{N}} C(\mathbb{R}^k, \mathbb{R}^l) \right)$ the function which satisfies for all $L \in \mathbb{N}$, $l_0, l_1, \ldots, l_L \in \mathbb{N}$, $\Phi = ((W_1, B_1), (W_2, B_2), \ldots, (W_L, B_L)) \in \left( \bigtimes_{k=1}^{L} (\mathbb{R}^{l_k \times l_{k-1}} \times \mathbb{R}^{l_k}) \right)$, $x_0 \in \mathbb{R}^{l_0}, x_1 \in \mathbb{R}^{l_1}, \ldots, x_L \in \mathbb{R}^{l_L}$ with $\forall k \in \{1, 2, \ldots, L\} \colon x_k = \mathfrak{M}_{a, l_k}(W_k x_{k-1} + B_k)$ that*

$$\mathcal{R}_a(\Phi) \in C(\mathbb{R}^{l_0}, \mathbb{R}^{l_L}) \qquad and \qquad (\mathcal{R}_a(\Phi))(x_0) = W_L x_{L-1} + B_L \tag{2.53}$$

*(cf. Definitions 2.1.4 and 2.2.1).*

**Lemma 2.2.4.** *Let $\Phi \in \mathbf{N}$ (cf. Definition 2.2.1). Then*

  (i) *it holds that $\mathcal{D}(\Phi) \in \mathbb{N}^{\mathcal{L}(\Phi)+1}$ and*

  (ii) *it holds for all $a \in C(\mathbb{R}, \mathbb{R})$ that $\mathcal{R}_a(\Phi) \in C(\mathbb{R}^{\mathcal{I}(\Phi)}, \mathbb{R}^{\mathcal{O}(\Phi)})$*

*(cf. Definition 2.2.3).*

*Proof of Lemma 2.2.4.* Note that the assumption that

$$\Phi \in \mathbf{N} = \bigcup_{L \in \mathbb{N}} \bigcup_{(l_0, l_1, \ldots, l_L) \in \mathbb{N}^{L+1}} \left( \bigtimes_{k=1}^{L} (\mathbb{R}^{l_k \times l_{k-1}} \times \mathbb{R}^{l_k}) \right) \tag{2.54}$$

ensures that there exist $L \in \mathbb{N}$, $l_0, l_1, \ldots, l_L \in \mathbb{N}$ such that

$$\Phi \in \left( \bigtimes_{k=1}^{L} (\mathbb{R}^{l_k \times l_{k-1}} \times \mathbb{R}^{l_k}) \right). \tag{2.55}$$

Observe that (2.55) assures that

$$\mathcal{L}(\Phi) = L, \qquad \mathcal{I}(\Phi) = l_0, \qquad \mathcal{O}(\Phi) = l_L, \tag{2.56}$$

$$\text{and} \qquad \mathcal{D}(\Phi) = (l_0, l_1, \ldots, l_L) \in \mathbb{N}^{L+1} = \mathbb{N}^{\mathcal{L}(\Phi)+1}. \tag{2.57}$$

This establishes item (i). Moreover, note that (2.56) and (2.53) show that $\mathcal{R}_a(\Phi) \in C(\mathbb{R}^{\mathcal{I}(\Phi)}, \mathbb{R}^{\mathcal{O}(\Phi)})$. This establishes item (ii). The proof of Lemma 2.2.4 is thus complete. $\square$

**Exercise 2.2.1.** *Prove or disprove the following statement: There exists $\Phi \in \mathbf{N}$ such that*

$$\mathcal{R}_{\tanh}(\Phi) = \mathfrak{l} \tag{2.58}$$

*(cf. Definitions 2.1.16, 2.1.20, 2.2.1, and 2.2.3).*

## 2.2.3 Compositions of ANNs

### 2.2.3.1 Standard compositions of ANNs

**Definition 2.2.5** (Composition of ANNs)**.** *We denote by $(\cdot) \bullet (\cdot) \colon \{(\Phi, \Psi) \in \mathbf{N} \times \mathbf{N} \colon \mathcal{I}(\Phi) = \mathcal{O}(\Psi)\} \to \mathbf{N}$ the function which satisfies for all $L, \mathfrak{L} \in \mathbb{N}$, $l_0, l_1, \ldots, l_L, \mathfrak{l}_0, \mathfrak{l}_1, \ldots, \mathfrak{l}_{\mathfrak{L}} \in \mathbb{N}$, $\Phi = ((W_1, B_1), (W_2, B_2), \ldots, (W_L, B_L)) \in \left( \bigtimes_{k=1}^{L} (\mathbb{R}^{l_k \times l_{k-1}} \times \mathbb{R}^{l_k}) \right)$, $\Psi = ((\mathscr{W}_1, \mathscr{B}_1), (\mathscr{W}_2, \mathscr{B}_2), \ldots, (\mathscr{W}_{\mathfrak{L}}, \mathscr{B}_{\mathfrak{L}})) \in \left( \bigtimes_{k=1}^{\mathfrak{L}} (\mathbb{R}^{\mathfrak{l}_k \times \mathfrak{l}_{k-1}} \times \mathbb{R}^{\mathfrak{l}_k}) \right)$ with $l_0 = \mathcal{I}(\Phi) = \mathcal{O}(\Psi) = \mathfrak{l}_{\mathfrak{L}}$ that*

$\Phi \bullet \Psi =$

$$\begin{cases} ((\mathscr{W}_1, \mathscr{B}_1), (\mathscr{W}_2, \mathscr{B}_2), \ldots, (\mathscr{W}_{\mathfrak{L}-1}, \mathscr{B}_{\mathfrak{L}-1}), (W_1 \mathscr{W}_{\mathfrak{L}}, W_1 \mathscr{B}_{\mathfrak{L}} + B_1), & \\ \qquad\qquad (W_2, B_2), (W_3, B_3), \ldots, (W_L, B_L)) & : (L > 1) \wedge (\mathfrak{L} > 1) \\ ((W_1 \mathscr{W}_1, W_1 \mathscr{B}_1 + B_1), (W_2, B_2), (W_3, B_3), \ldots, (W_L, B_L)) & : (L > 1) \wedge (\mathfrak{L} = 1) \\ ((\mathscr{W}_1, \mathscr{B}_1), (\mathscr{W}_2, \mathscr{B}_2), \ldots, (\mathscr{W}_{\mathfrak{L}-1}, \mathscr{B}_{\mathfrak{L}-1}), (W_1 \mathscr{W}_{\mathfrak{L}}, W_1 \mathscr{B}_{\mathfrak{L}} + B_1)) & : (L = 1) \wedge (\mathfrak{L} > 1) \\ ((W_1 \mathscr{W}_1, W_1 \mathscr{B}_1 + B_1)) & : (L = 1) \wedge (\mathfrak{L} = 1) \end{cases} \tag{2.59}$$

*(cf. Definition 2.2.1).*

### 2.2.3.2 Elementary properties of standard compositions of ANNs

**Lemma 2.2.6.** *Let $\Phi, \Psi \in \mathbf{N}$ satisfy $\mathcal{I}(\Phi) = \mathcal{O}(\Psi)$. Then*

*(i) it holds that $\mathcal{L}(\Phi \bullet \Psi) = \mathcal{L}(\Phi) + \mathcal{L}(\Psi) - 1$ and*

*(ii) it holds for all $i \in \{1, 2, \ldots, \mathcal{L}(\Phi \bullet \Psi)\}$ that*

$$
(\mathcal{W}_{i,(\Phi\bullet\Psi)}, \mathcal{B}_{i,(\Phi\bullet\Psi)}) = \begin{cases} (\mathcal{W}_{i,\Psi}, \mathcal{B}_{i,\Psi}) & : i < \mathcal{L}(\Psi) \\ (\mathcal{W}_{1,\Phi}\mathcal{W}_{\mathcal{L}(\Psi),\Psi}, \mathcal{W}_{1,\Phi}\mathcal{B}_{\mathcal{L}(\Psi),\Psi} + \mathcal{B}_{1,\Phi}) & : i = \mathcal{L}(\Psi) \\ (\mathcal{W}_{i-\mathcal{L}(\Psi)+1,\Phi}, \mathcal{B}_{i-\mathcal{L}(\Psi)+1,\Phi}) & : i > \mathcal{L}(\Psi). \end{cases} \quad (2.60)
$$

*Proof of Lemma 2.2.6.* Note that (2.59) clearly implies items (i) and (ii). The proof of Lemma 2.2.6 is thus complete. $\qquad\square$

**Proposition 2.2.7.** *Let $\Phi_1, \Phi_2 \in \mathbf{N}$ satisfy $\mathcal{I}(\Phi_1) = \mathcal{O}(\Phi_2)$ (cf. Definition 2.2.1). Then*

*(i) it holds that*

$$
\mathcal{D}(\Phi_1 \bullet \Phi_2) = (\mathbb{D}_0(\Phi_2), \mathbb{D}_1(\Phi_2), \ldots, \mathbb{D}_{\mathcal{H}(\Phi_2)}(\Phi_2), \mathbb{D}_1(\Phi_1), \mathbb{D}_2(\Phi_1), \ldots, \mathbb{D}_{\mathcal{L}(\Phi_1)}(\Phi_1)),
$$
$$(2.61)$$

*(ii) it holds that*

$$
[\mathcal{L}(\Phi_1 \bullet \Phi_2) - 1] = [\mathcal{L}(\Phi_1) - 1] + [\mathcal{L}(\Phi_2) - 1], \quad (2.62)
$$

*(iii) it holds that*

$$
\mathcal{H}(\Phi_1 \bullet \Phi_2) = \mathcal{H}(\Phi_1) + \mathcal{H}(\Phi_2), \quad (2.63)
$$

*(iv) it holds that*

$$
\begin{aligned}
\mathcal{P}(\Phi_1 \bullet \Phi_2) &= \mathcal{P}(\Phi_1) + \mathcal{P}(\Phi_2) + \mathbb{D}_1(\Phi_1)(\mathbb{D}_{\mathcal{L}(\Phi_2)-1}(\Phi_2) + 1) \\
&\quad - \mathbb{D}_1(\Phi_1)(\mathbb{D}_0(\Phi_1) + 1) - \mathbb{D}_{\mathcal{L}(\Phi_2)}(\Phi_2)(\mathbb{D}_{\mathcal{L}(\Phi_2)-1}(\Phi_2) + 1) \quad (2.64) \\
&\leq \mathcal{P}(\Phi_1) + \mathcal{P}(\Phi_2) + \mathbb{D}_1(\Phi_1)\mathbb{D}_{\mathcal{H}(\Phi_2)}(\Phi_2),
\end{aligned}
$$

*and*

*(v) it holds for all $a \in C(\mathbb{R}, \mathbb{R})$ that $\mathcal{R}_a(\Phi_1 \bullet \Phi_2) \in C(\mathbb{R}^{\mathcal{I}(\Phi_2)}, \mathbb{R}^{\mathcal{O}(\Phi_1)})$ and*

$$
\mathcal{R}_a(\Phi_1 \bullet \Phi_2) = [\mathcal{R}_a(\Phi_1)] \circ [\mathcal{R}_a(\Phi_2)] \quad (2.65)
$$

*(cf. Definitions 2.2.3 and 2.2.5).*

*Proof of Proposition 2.2.7.* Throughout this proof let $a \in C(\mathbb{R}, \mathbb{R})$, let $L_k \in \mathbb{N}$, $k \in \{1, 2\}$, satisfy for all $k \in \{1, 2\}$ that $L_k = \mathcal{L}(\Phi_k)$, let $l_{1,0}, l_{1,1}, \ldots, l_{1,\mathcal{L}(\Phi_1)}, l_{2,0}, l_{2,1}, \ldots, l_{2,\mathcal{L}(\Phi_2)} \in \mathbb{N}$, $\left((W_{k,1}, B_{k,1}), (W_{k,2}, B_{k,2}), \ldots, (W_{k,L_k}, B_{k,L_k})\right) \in (\bigtimes_{j=1}^{L_k}(\mathbb{R}^{l_{k,j} \times l_{k,j-1}} \times \mathbb{R}^{l_{k,j}}))$, $k \in \{1, 2\}$, satisfy for all $k \in \{1, 2\}$ that

$$
\Phi_k = \left((W_{k,1}, B_{k,1}), (W_{k,2}, B_{k,2}), \ldots, (W_{k,L_k}, B_{k,L_k})\right), \quad (2.66)
$$

let $L_3 \in \mathbb{N}$, $l_{3,0}, l_{3,1}, \ldots, l_{3,L_3} \in \mathbb{N}$, $\Phi_3 = \left((W_{3,1}, B_{3,1}), \ldots, (W_{3,L_3}, B_{3,L_3})\right) \in (\bigtimes_{j=1}^{L_3} (\mathbb{R}^{l_{3,j} \times l_{3,j-1}} \times \mathbb{R}^{l_{3,j}}))$ satisfy that $\Phi_3 = \Phi_1 \bullet \Phi_2$, let $x_0 \in \mathbb{R}^{l_{2,0}}, x_1 \in \mathbb{R}^{l_{2,1}}, \ldots, x_{L_2-1} \in \mathbb{R}^{l_{2,L_2-1}}$ satisfy

$$
\forall j \in \mathbb{N} \cap (0, L_2): \; x_j = \mathfrak{M}_{a,l_{2,j}}(W_{2,j}x_{j-1} + B_{2,j}) \quad (2.67)
$$

(cf. Definition 2.1.4), let $y_0 \in \mathbb{R}^{l_{1,0}}, y_1 \in \mathbb{R}^{l_{1,1}}, \ldots, y_{L_1-1} \in \mathbb{R}^{l_{1,L_1-1}}$ satisfy $y_0 = W_{2,L_2}x_{L_2-1} + B_{2,L_2}$ and

$$\forall\, j \in \mathbb{N} \cap (0, L_1): \; y_j = \mathfrak{M}_{a,l_{1,j}}(W_{1,j}y_{j-1} + B_{1,j}), \tag{2.68}$$

and let $z_0 \in \mathbb{R}^{l_{3,0}}, z_1 \in \mathbb{R}^{l_{3,1}}, \ldots, z_{L_3-1} \in \mathbb{R}^{l_{3,L_3-1}}$ satisfy $z_0 = x_0$ and

$$\forall\, j \in \mathbb{N} \cap (0, L_3): \; z_j = \mathfrak{M}_{a,l_{3,j}}(W_{3,j}z_{j-1} + B_{3,j}). \tag{2.69}$$

Note that (2.59) ensures that

$$\Phi_3 = \Phi_1 \bullet \Phi_2 =$$
$$\begin{cases} \begin{aligned} &((W_{2,1}, B_{2,1}),(W_{2,2}, B_{2,2}),\ldots,(W_{2,L_2-1}, B_{2,L_2-1}), \\ &\quad (W_{1,1}W_{2,L_2}, W_{1,1}B_{2,L_2} + B_{1,1}),(W_{1,2}, B_{1,2}), \\ &\quad\quad (W_{1,3}, B_{1,3}),\ldots,(W_{1,L_1}, B_{1,L_1})) \end{aligned} & : (L_1 > 1) \wedge (L_2 > 1) \\[2ex] \begin{aligned} &((W_{1,1}W_{2,1}, W_{1,1}B_{2,1} + B_{1,1}),(W_{1,2}, B_{1,2}), \\ &\quad (W_{1,3}, B_{1,3}),\ldots,(W_{1,L_1}, B_{1,L_1})) \end{aligned} & : (L_1 > 1) \wedge (L_2 = 1) \\[2ex] \begin{aligned} &((W_{2,1}, B_{2,1}),(W_{2,2}, B_{2,2}),\ldots,(W_{2,L_2-1}, B_{2,L_2-1}), \\ &\quad (W_{1,1}W_{2,L_2}, W_{1,1}B_{2,L_2} + B_{1,1})) \end{aligned} & : (L_1 = 1) \wedge (L_2 > 1) \\[2ex] (W_{1,1}W_{2,1}, W_{1,1}B_{2,1} + B_{1,1}) & : (L_1 = 1) \wedge (L_2 = 1). \end{cases} \tag{2.70}$$

Hence, we obtain that

$$\begin{aligned} [\mathcal{L}(\Phi_1 \bullet \Phi_2) - 1] &= [(L_2 - 1) + 1 + (L_1 - 1)] - 1 \\ &= [L_1 - 1] + [L_2 - 1] = [\mathcal{L}(\Phi_1) - 1] + [\mathcal{L}(\Phi_2) - 1] \end{aligned} \tag{2.71}$$

$$\text{and} \qquad \mathcal{D}(\Phi_1 \bullet \Phi_2) = (l_{2,0}, l_{2,1}, \ldots, l_{2,L_2-1}, l_{1,1}, l_{1,2}, \ldots, l_{1,L_1}). \tag{2.72}$$

This establishes items (i)–(iii). In addition, observe that (2.72) demonstrates that

$$\begin{aligned} \mathcal{P}(\Phi_1 \bullet \Phi_2) &= \sum_{j=1}^{L_3} l_{3,j}(l_{3,j-1} + 1) \\ &= \left[\sum_{j=1}^{L_2-1} l_{3,j}(l_{3,j-1} + 1)\right] + l_{3,L_2}(l_{3,L_2-1} + 1) + \left[\sum_{j=L_2+1}^{L_3} l_{3,j}(l_{3,j-1} + 1)\right] \\ &= \left[\sum_{j=1}^{L_2-1} l_{2,j}(l_{2,j-1} + 1)\right] + l_{1,1}(l_{2,L_2-1} + 1) + \left[\sum_{j=L_2+1}^{L_3} l_{1,j-L_2+1}(l_{1,j-L_2} + 1)\right] \\ &= \left[\sum_{j=1}^{L_2-1} l_{2,j}(l_{2,j-1} + 1)\right] + \left[\sum_{j=2}^{L_1} l_{1,j}(l_{1,j-1} + 1)\right] + l_{1,1}(l_{2,L_2-1} + 1) \\ &= \left[\sum_{j=1}^{L_2} l_{2,j}(l_{2,j-1} + 1)\right] + \left[\sum_{j=1}^{L_1} l_{1,j}(l_{1,j-1} + 1)\right] + l_{1,1}(l_{2,L_2-1} + 1) \\ &\quad - l_{2,L_2}(l_{2,L_2-1} + 1) - l_{1,1}(l_{1,0} + 1) \\ &= \mathcal{P}(\Phi_1) + \mathcal{P}(\Phi_2) + l_{1,1}(l_{2,L_2-1} + 1) - l_{2,L_2}(l_{2,L_2-1} + 1) \\ &\quad - l_{1,1}(l_{1,0} + 1) \\ &\leq \mathcal{P}(\Phi_1) + \mathcal{P}(\Phi_2) + l_{1,1}l_{2,L_2-1}. \end{aligned} \tag{2.73}$$

This establishes item (iv). Moreover, observe that (2.70) and the fact that $a \in C(\mathbb{R}, \mathbb{R})$ ensure that

$$\mathcal{R}_a(\Phi_1 \bullet \Phi_2) \in C(\mathbb{R}^{l_{2,0}}, \mathbb{R}^{l_{1,L_1}}) = C(\mathbb{R}^{\mathcal{I}(\Phi_2)}, \mathbb{R}^{\mathcal{O}(\Phi_1)}). \tag{2.74}$$

Next note that (2.71) implies that $L_3 = L_1 + L_2 - 1$. This, (2.70), and (2.72) ensure that

$$(l_{3,0}, l_{3,1}, \ldots, l_{3,L_1+L_2-1}) = (l_{2,0}, l_{2,1}, \ldots, l_{2,L_2-1}, l_{1,1}, l_{1,2}, \ldots, l_{1,L_1}), \tag{2.75}$$

$$\left[\forall\, j \in \mathbb{N} \cap (0, L_2)\colon\ (W_{3,j}, B_{3,j}) = (W_{2,j}, B_{2,j})\right], \tag{2.76}$$

$$(W_{3,L_2}, B_{3,L_2}) = (W_{1,1}W_{2,L_2}, W_{1,1}B_{2,L_2} + B_{1,1}), \tag{2.77}$$

and $\quad \left[\forall\, j \in \mathbb{N} \cap (L_2, L_1 + L_2)\colon\ (W_{3,j}, B_{3,j}) = (W_{1,j+1-L_2}, B_{1,j+1-L_2})\right]. \tag{2.78}$

This, (2.67), (2.69), and induction imply that for all $j \in \mathbb{N}_0 \cap [0, L_2)$ it holds that $z_j = x_j$. Combining this with (2.77) and the fact that $y_0 = W_{2,L_2}x_{L_2-1} + B_{2,L_2}$ ensures that

$$\begin{aligned}
W_{3,L_2}z_{L_2-1} + B_{3,L_2} &= W_{3,L_2}x_{L_2-1} + B_{3,L_2} \\
&= W_{1,1}W_{2,L_2}x_{L_2-1} + W_{1,1}B_{2,L_2} + B_{1,1} \\
&= W_{1,1}(W_{2,L_2}x_{L_2-1} + B_{2,L_2}) + B_{1,1} = W_{1,1}y_0 + B_{1,1}.
\end{aligned} \tag{2.79}$$

Next we claim that for all $j \in \mathbb{N} \cap [L_2, L_1 + L_2)$ it holds that

$$W_{3,j}z_{j-1} + B_{3,j} = W_{1,j+1-L_2}y_{j-L_2} + B_{1,j+1-L_2}. \tag{2.80}$$

We prove (2.80) by induction on $j \in \mathbb{N} \cap [L_2, L_1 + L_2)$. Note that (2.79) establishes (2.80) in the base case $j = L_2$. For the induction step note that the fact that $L_3 = L_1 + L_2 - 1$, (2.68), (2.69), (2.75), and (2.78) imply that for all $j \in \mathbb{N} \cap [L_2, \infty) \cap (0, L_1 + L_2 - 1)$ with

$$W_{3,j}z_{j-1} + B_{3,j} = W_{1,j+1-L_2}y_{j-L_2} + B_{1,j+1-L_2} \tag{2.81}$$

it holds that

$$\begin{aligned}
W_{3,j+1}z_j + B_{3,j+1} &= W_{3,j+1}\mathfrak{M}_{a,l_{3,j}}(W_{3,j}z_{j-1} + B_{3,j}) + B_{3,j+1} \\
&= W_{1,j+2-L_2}\mathfrak{M}_{a,l_{1,j+1-L_2}}(W_{1,j+1-L_2}y_{j-L_2} + B_{1,j+1-L_2}) + B_{1,j+2-L_2} \\
&= W_{1,j+2-L_2}y_{j+1-L_2} + B_{1,j+2-L_2}.
\end{aligned} \tag{2.82}$$

Induction hence proves (2.80). Next observe that (2.80) and the fact that $L_3 = L_1 + L_2 - 1$ assure that

$$W_{3,L_3}z_{L_3-1} + B_{3,L_3} = W_{3,L_1+L_2-1}z_{L_1+L_2-2} + B_{3,L_1+L_2-1} = W_{1,L_1}y_{L_1-1} + B_{1,L_1}. \tag{2.83}$$

The fact that $\Phi_3 = \Phi_1 \bullet \Phi_2$, (2.67), (2.68), and (2.69) therefore prove that

$$\begin{aligned}
[\mathcal{R}_a(\Phi_1 \bullet \Phi_2)](x_0) &= [\mathcal{R}_a(\Phi_3)](x_0) = [\mathcal{R}_a(\Phi_3)](z_0) = W_{3,L_3}z_{L_3-1} + B_{3,L_3} \\
&= W_{1,L_1}y_{L_1-1} + B_{1,L_1} = [\mathcal{R}_a(\Phi_1)](y_0) \\
&= [\mathcal{R}_a(\Phi_1)]\big(W_{2,L_2}x_{L_2-1} + B_{2,L_2}\big) \\
&= [\mathcal{R}_a(\Phi_1)]\big([\mathcal{R}_a(\Phi_2)](x_0)\big) = [(\mathcal{R}_a(\Phi_1)) \circ (\mathcal{R}_a(\Phi_2))](x_0).
\end{aligned} \tag{2.84}$$

Combining this with (2.74) establishes item (v). The proof of Proposition 2.2.7 is thus complete. $\qquad\square$

### 2.2.3.3 Associativity of standard compositions of ANNs

**Lemma 2.2.8.** *Let* $\Phi_1, \Phi_2, \Phi_3 \in \mathbf{N}$ *satisfy* $\mathcal{I}(\Phi_1) = \mathcal{O}(\Phi_2)$ *and* $\mathcal{I}(\Phi_2) = \mathcal{O}(\Phi_3)$ *(cf. Definition 2.2.1). Then it holds that*

$$(\Phi_1 \bullet \Phi_2) \bullet \Phi_3 = \Phi_1 \bullet (\Phi_2 \bullet \Phi_3) \tag{2.85}$$

*(cf. Definition 2.2.5).*

*Proof of Lemma 2.2.8.* Throughout this proof let $\Phi_4, \Phi_5, \Phi_6, \Phi_7 \in \mathbf{N}$ satisfy that $\Phi_4 = \Phi_1 \bullet \Phi_2$, $\Phi_5 = \Phi_2 \bullet \Phi_3$, $\Phi_6 = \Phi_4 \bullet \Phi_3$, and $\Phi_7 = \Phi_1 \bullet \Phi_5$, let $L_k \in \mathbb{N}$, $k \in \{1, 2, \ldots, 7\}$, satisfy for all $k \in \{1, 2, \ldots, 7\}$ that $L_k = \mathcal{L}(\Phi_k)$, let $l_{k,0}, l_{k,1}, \ldots, l_{k,L_k} \in \mathbb{N}$, $k \in \{1, 2, \ldots, 7\}$, and let $\big((W_{k,1}, B_{k,1}), (W_{k,2}, B_{k,2}), \ldots, (W_{k,L_k}, B_{k,L_k})\big) \in (\bigtimes_{j=1}^{L_k} (\mathbb{R}^{l_{k,j} \times l_{k,j-1}} \times \mathbb{R}^{l_{k,j}}))$, $k \in \{1, 2, \ldots, 7\}$, satisfy for all $k \in \{1, 2, \ldots, 7\}$ that

$$\Phi_k = \big((W_{k,1}, B_{k,1}), (W_{k,2}, B_{k,2}), \ldots, (W_{k,L_k}, B_{k,L_k})\big). \tag{2.86}$$

Observe that item (ii) in Proposition 2.2.7 and the fact that for all $k \in \{1, 2, 3\}$ it holds that $\mathcal{L}(\Phi_k) = L_k$ proves that

$$\begin{aligned}
\mathcal{L}(\Phi_6) &= \mathcal{L}((\Phi_1 \bullet \Phi_2) \bullet \Phi_3) = \mathcal{L}(\Phi_1 \bullet \Phi_2) + \mathcal{L}(\Phi_3) - 1 \\
&= \mathcal{L}(\Phi_1) + \mathcal{L}(\Phi_2) + \mathcal{L}(\Phi_3) - 2 = L_1 + L_2 + L_3 - 2 \\
&= \mathcal{L}(\Phi_1) + \mathcal{L}(\Phi_2 \bullet \Phi_3) - 1 = \mathcal{L}(\Phi_1 \bullet (\Phi_2 \bullet \Phi_3)) = \mathcal{L}(\Phi_7).
\end{aligned} \tag{2.87}$$

Next note that Lemma 2.2.6, (2.86), and the fact that $\Phi_4 = \Phi_1 \bullet \Phi_2$ imply that

$$\big[\forall j \in \mathbb{N} \cap (0, L_2) \colon (W_{4,j}, B_{4,j}) = (W_{2,j}, B_{2,j})\big], \tag{2.88}$$

$$(W_{4,L_2}, B_{4,L_2}) = (W_{1,1} W_{2,L_2}, W_{1,1} B_{2,L_2} + B_{1,1}), \tag{2.89}$$

$$\text{and} \qquad \big[\forall j \in \mathbb{N} \cap (L_2, L_1 + L_2) \colon (W_{4,j}, B_{4,j}) = (W_{1,j+1-L_2}, B_{1,j+1-L_2})\big]. \tag{2.90}$$

Hence, we obtain that

$$\big[\forall j \in \mathbb{N} \cap (L_3 - 1, L_2 + L_3 - 1) \colon (W_{4,j+1-L_3}, B_{4,j+1-L_3}) = (W_{2,j+1-L_3}, B_{2,j+1-L_3})\big], \tag{2.91}$$

$$(W_{4,L_2}, B_{4,L_2}) = (W_{1,1} W_{2,L_2}, W_{1,1} B_{2,L_2} + B_{1,1}), \tag{2.92}$$

and

$$\begin{aligned}
\big[\forall j \in \mathbb{N} \cap (L_2 + L_3 - 1, L_1 + L_2 + L_3 - 1) \colon& \\
(W_{4,j+1-L_3}, B_{4,j+1-L_3}) &= (W_{1,j+2-L_2-L_3}, B_{1,j+2-L_2-L_3})\big].
\end{aligned} \tag{2.93}$$

In addition, observe that Lemma 2.2.6, (2.86), and the fact that $\Phi_5 = \Phi_2 \bullet \Phi_3$ demonstrate that

$$\big[\forall j \in \mathbb{N} \cap (0, L_3) \colon (W_{5,j}, B_{5,j}) = (W_{3,j}, B_{3,j})\big], \tag{2.94}$$

$$(W_{5,L_3}, B_{5,L_3}) = (W_{2,1} W_{3,L_3}, W_{2,1} B_{3,L_3} + B_{2,1}), \tag{2.95}$$

$$\text{and} \qquad \big[\forall j \in \mathbb{N} \cap (L_3, L_2 + L_3) \colon (W_{5,j}, B_{5,j}) = (W_{2,j+1-L_3}, B_{2,j+1-L_3})\big]. \tag{2.96}$$

Moreover, note that Lemma 2.2.6, (2.86), and the fact that $\Phi_6 = \Phi_4 \bullet \Phi_3$ ensure that

$$\big[\forall j \in \mathbb{N} \cap (0, L_3) \colon (W_{6,j}, B_{6,j}) = (W_{3,j}, B_{3,j})\big], \tag{2.97}$$

$$(W_{6,L_3}, B_{6,L_3}) = (W_{4,1}W_{3,L_3}, W_{4,1}B_{3,L_3} + B_{4,1}), \tag{2.98}$$

$$\text{and} \qquad \left[\forall\, j \in \mathbb{N} \cap (L_3, L_4 + L_3)\colon (W_{6,j}, B_{6,j}) = (W_{4,j+1-L_3}, B_{4,j+1-L_3})\right]. \tag{2.99}$$

Furthermore, observe that Lemma 2.2.6, (2.86), and the fact that $\Phi_7 = \Phi_1 \bullet \Phi_5$ show that

$$\left[\forall\, j \in \mathbb{N} \cap (0, L_5)\colon (W_{7,j}, B_{7,j}) = (W_{5,j}, B_{5,j})\right], \tag{2.100}$$

$$(W_{7,L_5}, B_{7,L_5}) = (W_{1,1}W_{5,L_5}, W_{1,1}B_{5,L_5} + B_{1,1}), \tag{2.101}$$

$$\text{and} \qquad \left[\forall\, j \in \mathbb{N} \cap (L_5, L_1 + L_5)\colon (W_{7,j}, B_{7,j}) = (W_{1,j+1-L_5}, B_{1,j+1-L_5})\right]. \tag{2.102}$$

This, the fact that $L_3 \le L_2 + L_3 - 1 = L_5$, (2.94), and (2.97) imply that for all $j \in \mathbb{N} \cap (0, L_3)$ it holds that

$$(W_{6,j}, B_{6,j}) = (W_{3,j}, B_{3,j}) = (W_{5,j}, B_{5,j}) = (W_{7,j}, B_{7,j}). \tag{2.103}$$

In addition, observe that (2.88), (2.89), (2.94), (2.95), (2.98), (2.100), (2.101), and the fact that $L_5 = L_2 + L_3 - 1$ demonstrate that

$$
\begin{aligned}
(W_{6,L_3}, B_{6,L_3}) &= (W_{4,1}W_{3,L_3}, W_{4,1}B_{3,L_3} + B_{4,1}) \\
&= \begin{cases} (W_{2,1}W_{3,L_3}, W_{2,1}B_{3,L_3} + B_{2,1}) & : L_2 > 1 \\ (W_{1,1}W_{2,1}W_{3,L_3}, W_{1,1}W_{2,1}B_{3,L_3} + W_{1,1}B_{2,1} + B_{1,1}) & : L_2 = 1 \end{cases} \\
&= \begin{cases} (W_{2,1}W_{3,L_3}, W_{2,1}B_{3,L_3} + B_{2,1}) & : L_2 > 1 \\ (W_{1,1}(W_{2,1}W_{3,L_3}), W_{1,1}(W_{2,1}B_{3,L_3} + B_{2,1}) + B_{1,1}) & : L_2 = 1 \end{cases} \\
&= \begin{cases} (W_{5,L_3}, B_{5,L_3}) & : L_2 > 1 \\ (W_{1,1}W_{5,L_3}, W_{1,1}B_{5,L_3} + B_{1,1}) & : L_2 = 1 \end{cases} \\
&= (W_{7,L_3}, B_{7,L_3}).
\end{aligned}
\tag{2.104}
$$

Next note that the fact that $L_5 = L_2 + L_3 - 1 < L_1 + L_2 + L_3 - 1 = L_3 + L_4$, (2.99), (2.91), (2.96), and (2.100) ensure that for all $j \in \mathbb{N}$ with $L_3 < j < L_5$ it holds that

$$
\begin{aligned}
(W_{6,j}, B_{6,j}) &= (W_{4,j+1-L_3}, B_{4,j+1-L_3}) = (W_{2,j+1-L_3}, B_{2,j+1-L_3}) \\
&= (W_{5,j}, B_{5,j}) = (W_{7,j}, B_{7,j}).
\end{aligned}
\tag{2.105}
$$

Moreover, observe that the fact that $L_5 = L_2 + L_3 - 1 < L_1 + L_2 + L_3 - 1 = L_3 + L_4$, (2.99), (2.104), (2.89), (2.96), and (2.101) prove that

$$
\begin{aligned}
(W_{6,L_5}, B_{6,L_5}) &= \begin{cases} (W_{4,L_5+1-L_3}, B_{4,L_5+1-L_3}) & : L_2 > 1 \\ (W_{6,L_3}, B_{6,L_3}) & : L_2 = 1 \end{cases} \\
&= \begin{cases} (W_{4,L_2}, B_{4,L_2}) & : L_2 > 1 \\ (W_{7,L_3}, B_{7,L_3}) & : L_2 = 1 \end{cases} \\
&= \begin{cases} (W_{1,1}W_{2,L_2}, W_{1,1}B_{2,L_2} + B_{1,1}) & : L_2 > 1 \\ (W_{7,L_5}, B_{7,L_5}) & : L_2 = 1 \end{cases} \\
&= \begin{cases} (W_{1,1}W_{5,L_5}, W_{1,1}B_{5,L_5} + B_{1,1}) & : L_2 > 1 \\ (W_{7,L_5}, B_{7,L_5}) & : L_2 = 1 \end{cases} \\
&= (W_{7,L_5}, B_{7,L_5}).
\end{aligned}
\tag{2.106}
$$

Furthermore, note that (2.99), (2.93), (2.102), and the fact that $L_5 = L_2 + L_3 - 1 \geq L_3$ assure that for all $j \in \mathbb{N}$ with $L_5 < j \leq L_6$ it holds that

$$
\begin{aligned}
(W_{6,j}, B_{6,j}) &= (W_{4,j+1-L_3}, B_{4,j+1-L_3}) = (W_{1,j+2-L_2-L_3}, B_{1,j+2-L_2-L_3}) \\
&= (W_{1,j+1-L_5}, B_{1,j+1-L_5}) = (W_{7,j}, B_{7,j}).
\end{aligned}
\tag{2.107}
$$

Combining this with (2.87), (2.103), (2.104), (2.105), and (2.106) establishes that

$$
(\Phi_1 \bullet \Phi_2) \bullet \Phi_3 = \Phi_4 \bullet \Phi_3 = \Phi_6 = \Phi_7 = \Phi_1 \bullet \Phi_5 = \Phi_1 \bullet (\Phi_2 \bullet \Phi_3).
\tag{2.108}
$$

The proof of Lemma 2.2.8 is thus complete. $\qquad\qquad\square$

#### 2.2.3.4 Powers and extensions of ANNs

**Definition 2.2.9.** *Let $d \in \mathbb{N}$. Then we denote by $\mathrm{I}_d \in \mathbb{R}^{d \times d}$ the identity matrix in $\mathbb{R}^{d \times d}$.*

**Definition 2.2.10.** *We denote by $(\cdot)^{\bullet n} \colon \{\Phi \in \mathbf{N} \colon \mathcal{I}(\Phi) = \mathcal{O}(\Phi)\} \to \mathbf{N}$, $n \in \mathbb{N}_0$, the functions which satisfy for all $n \in \mathbb{N}_0$, $\Phi \in \mathbf{N}$ with $\mathcal{I}(\Phi) = \mathcal{O}(\Phi)$ that*

$$
\Phi^{\bullet n} = \begin{cases} \big(\mathrm{I}_{\mathcal{O}(\Phi)}, (0, 0, \dots, 0)\big) \in \mathbb{R}^{\mathcal{O}(\Phi) \times \mathcal{O}(\Phi)} \times \mathbb{R}^{\mathcal{O}(\Phi)} & : n = 0 \\ \Phi \bullet (\Phi^{\bullet(n-1)}) & : n \in \mathbb{N} \end{cases}
\tag{2.109}
$$

*(cf. Definitions 2.2.1, 2.2.5, and 2.2.9).*

### 2.2.4 Parallelizations of ANNs

#### 2.2.4.1 Parallelizations of ANNs with the same length

**Definition 2.2.11** (Parallelization of ANNs). *Let $n \in \mathbb{N}$. Then we denote by*

$$
\mathbf{P}_n \colon \big\{ (\Phi_1, \Phi_2, \dots, \Phi_n) \in \mathbf{N}^n \colon \mathcal{L}(\Phi_1) = \mathcal{L}(\Phi_2) = \dots = \mathcal{L}(\Phi_n) \big\} \to \mathbf{N}
\tag{2.110}
$$

*the function which satisfies for all $L \in \mathbb{N}$, $\Phi_1, \Phi_2, \dots, \Phi_n \in \mathbf{N}$ with $L = \mathcal{L}(\Phi_1) = \mathcal{L}(\Phi_2) = \dots = \mathcal{L}(\Phi_n)$ that*

$$
\mathbf{P}_n(\Phi_1, \Phi_2, \dots, \Phi_n) = \left( \left( \begin{pmatrix} \mathcal{W}_{1,\Phi_1} & 0 & 0 & \cdots & 0 \\ 0 & \mathcal{W}_{1,\Phi_2} & 0 & \cdots & 0 \\ 0 & 0 & \mathcal{W}_{1,\Phi_3} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \mathcal{W}_{1,\Phi_n} \end{pmatrix}, \begin{pmatrix} \mathcal{B}_{1,\Phi_1} \\ \mathcal{B}_{1,\Phi_2} \\ \mathcal{B}_{1,\Phi_3} \\ \vdots \\ \mathcal{B}_{1,\Phi_n} \end{pmatrix} \right), \right.
$$
$$
\left( \begin{pmatrix} \mathcal{W}_{2,\Phi_1} & 0 & 0 & \cdots & 0 \\ 0 & \mathcal{W}_{2,\Phi_2} & 0 & \cdots & 0 \\ 0 & 0 & \mathcal{W}_{2,\Phi_3} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \mathcal{W}_{2,\Phi_n} \end{pmatrix}, \begin{pmatrix} \mathcal{B}_{2,\Phi_1} \\ \mathcal{B}_{2,\Phi_2} \\ \mathcal{B}_{2,\Phi_3} \\ \vdots \\ \mathcal{B}_{2,\Phi_n} \end{pmatrix} \right), \dots,
\tag{2.111}
$$
$$
\left. \left( \begin{pmatrix} \mathcal{W}_{L,\Phi_1} & 0 & 0 & \cdots & 0 \\ 0 & \mathcal{W}_{L,\Phi_2} & 0 & \cdots & 0 \\ 0 & 0 & \mathcal{W}_{L,\Phi_3} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \mathcal{W}_{L,\Phi_n} \end{pmatrix}, \begin{pmatrix} \mathcal{B}_{L,\Phi_1} \\ \mathcal{B}_{L,\Phi_2} \\ \mathcal{B}_{L,\Phi_3} \\ \vdots \\ \mathcal{B}_{L,\Phi_n} \end{pmatrix} \right) \right)
$$

*(cf. Definition 2.2.1).*

**Lemma 2.2.12.** *Let* $n, L \in \mathbb{N}$, $\Phi_1, \Phi_2, \ldots, \Phi_n \in \mathbf{N}$ *satisfy* $L = \mathcal{L}(\Phi_1) = \mathcal{L}(\Phi_2) = \ldots = \mathcal{L}(\Phi_n)$ *(cf. Definition 2.2.1). Then it holds that*

$$\mathbf{P}_n(\Phi_1, \Phi_2, \ldots, \Phi_n) \in \left( \bigtimes_{k=1}^{L} \left( \mathbb{R}^{(\sum_{j=1}^n \mathbb{D}_k(\Phi_j)) \times (\sum_{j=1}^n \mathbb{D}_{k-1}(\Phi_j))} \times \mathbb{R}^{(\sum_{j=1}^n \mathbb{D}_k(\Phi_j))} \right) \right) \qquad (2.112)$$

*(cf. Definition 2.2.11).*

*Proof of Lemma 2.2.12.* Note that (2.111) proves (2.112). The proof of Lemma 2.2.12 is thus complete. $\qquad\square$

**Proposition 2.2.13.** *Let* $a \in C(\mathbb{R}, \mathbb{R})$, $n \in \mathbb{N}$, $\Phi = (\Phi_1, \Phi_2, \ldots, \Phi_n) \in \mathbf{N}^n$ *satisfy* $\mathcal{L}(\Phi_1) = \mathcal{L}(\Phi_2) = \ldots = \mathcal{L}(\Phi_n)$ *(cf. Definition 2.2.1). Then*

(i) *it holds that*

$$\mathcal{R}_a(\mathbf{P}_n(\Phi)) \in C\left( \mathbb{R}^{[\sum_{j=1}^n \mathcal{I}(\Phi_j)]}, \mathbb{R}^{[\sum_{j=1}^n \mathcal{O}(\Phi_j)]} \right) \qquad (2.113)$$

*and*

(ii) *it holds for all* $x_1 \in \mathbb{R}^{\mathcal{I}(\Phi_1)}, x_2 \in \mathbb{R}^{\mathcal{I}(\Phi_2)}, \ldots, x_n \in \mathbb{R}^{\mathcal{I}(\Phi_n)}$ *that*

$$\begin{aligned} &\left( \mathcal{R}_a(\mathbf{P}_n(\Phi)) \right)(x_1, x_2, \ldots, x_n) \\ &= \left( (\mathcal{R}_a(\Phi_1))(x_1), (\mathcal{R}_a(\Phi_2))(x_2), \ldots, (\mathcal{R}_a(\Phi_n))(x_n) \right) \in \mathbb{R}^{[\sum_{j=1}^n \mathcal{O}(\Phi_j)]} \end{aligned} \qquad (2.114)$$

*(cf. Definitions 2.2.3 and 2.2.11).*

*Proof of Proposition 2.2.13.* Throughout this proof let $L \in \mathbb{N}$ satisfy $L = \mathcal{L}(\Phi_1)$, let $l_{j,0}, l_{j,1}, \ldots, l_{j,L} \in \mathbb{N}$, $j \in \{1, 2, \ldots, n\}$, satisfy for all $j \in \{1, 2, \ldots, n\}$ that $\mathcal{D}(\Phi_j) = (l_{j,0}, l_{j,1}, \ldots, l_{j,L})$, let $((W_{j,1}, B_{j,1}), (W_{j,2}, B_{j,2}), \ldots, (W_{j,L}, B_{j,L})) \in (\bigtimes_{k=1}^{L} (\mathbb{R}^{l_{j,k} \times l_{j,k-1}} \times \mathbb{R}^{l_{j,k}}))$, $j \in \{1, 2, \ldots, n\}$, satisfy for all $j \in \{1, 2, \ldots, n\}$ that

$$\Phi_j = ((W_{j,1}, B_{j,1}), (W_{j,2}, B_{j,2}), \ldots, (W_{j,L}, B_{j,L})), \qquad (2.115)$$

let $\alpha_k \in \mathbb{N}$, $k \in \{0, 1, \ldots, L\}$, satisfy for all $k \in \{0, 1, \ldots, L\}$ that $\alpha_k = \sum_{j=1}^n l_{j,k}$, let $((A_1, b_1), (A_2, b_2), \ldots, (A_L, b_L)) \in (\bigtimes_{k=1}^{L} (\mathbb{R}^{\alpha_k \times \alpha_{k-1}} \times \mathbb{R}^{\alpha_k}))$ satisfy that

$$\mathbf{P}_n(\Phi) = ((A_1, b_1), (A_2, b_2), \ldots, (A_L, b_L)) \qquad (2.116)$$

(cf. Lemma 2.2.12), let $(x_{j,0}, x_{j,1}, \ldots, x_{j,L-1}) \in (\mathbb{R}^{l_{j,0}} \times \mathbb{R}^{l_{j,1}} \times \ldots \times \mathbb{R}^{l_{j,L-1}})$, $j \in \{1, 2, \ldots, n\}$, satisfy for all $j \in \{1, 2, \ldots, n\}$, $k \in \mathbb{N} \cap (0, L)$ that

$$x_{j,k} = \mathfrak{M}_{a,l_{j,k}}(W_{j,k} x_{j,k-1} + B_{j,k}) \qquad (2.117)$$

(cf. Definition 2.1.4), and let $\mathfrak{x}_0 \in \mathbb{R}^{\alpha_0}, \mathfrak{x}_1 \in \mathbb{R}^{\alpha_1}, \ldots, \mathfrak{x}_{L-1} \in \mathbb{R}^{\alpha_{L-1}}$ satisfy for all $k \in \{0, 1, \ldots, L-1\}$ that $\mathfrak{x}_k = (x_{1,k}, x_{2,k}, \ldots, x_{n,k})$. Observe that (2.116) demonstrates that $\mathcal{I}(\mathbf{P}_n(\Phi)) = \alpha_0$ and $\mathcal{O}(\mathbf{P}_n(\Phi)) = \alpha_L$. Combining this with item (ii) in Lemma 2.2.4, the fact that for all $k \in \{0, 1, \ldots, L\}$ it holds that $\alpha_k = \sum_{j=1}^n l_{j,k}$, the fact that for all $j \in \{1, 2, \ldots, n\}$ it holds that $\mathcal{I}(\Phi_j) = l_{j,0}$, and the fact that for all $j \in \{1, 2, \ldots, n\}$ it holds that $\mathcal{O}(\Phi_j) = l_{j,L}$ ensures that

$$\begin{aligned} \mathcal{R}_a(\mathbf{P}_n(\Phi)) &\in C(\mathbb{R}^{\alpha_0}, \mathbb{R}^{\alpha_L}) = C\left( \mathbb{R}^{[\sum_{j=1}^n l_{j,0}]}, \mathbb{R}^{[\sum_{j=1}^n l_{j,L}]} \right) \\ &= C\left( \mathbb{R}^{[\sum_{j=1}^n \mathcal{I}(\Phi_j)]}, \mathbb{R}^{[\sum_{j=1}^n \mathcal{O}(\Phi_j)]} \right). \end{aligned} \qquad (2.118)$$

This proves item (i). Moreover, observe that (2.111) and (2.116) demonstrate that for all $k \in \{1, 2, \ldots, L\}$ it holds that

$$A_k = \begin{pmatrix} W_{1,k} & 0 & 0 & \cdots & 0 \\ 0 & W_{2,k} & 0 & \cdots & 0 \\ 0 & 0 & W_{3,k} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & W_{n,k} \end{pmatrix} \quad \text{and} \quad b_k = \begin{pmatrix} B_{1,k} \\ B_{2,k} \\ B_{3,k} \\ \vdots \\ B_{n,k} \end{pmatrix}. \tag{2.119}$$

Combining this with (2.8), (2.117), and the fact that for all $k \in \mathbb{N} \cap [0, L)$ it holds that $\mathfrak{x}_k = (x_{1,k}, x_{2,k}, \ldots, x_{n,k})$ implies that for all $k \in \mathbb{N} \cap (0, L)$ it holds that

$$\mathfrak{M}_{a,\alpha_k}(A_k \mathfrak{x}_{k-1} + b_k) = \begin{pmatrix} \mathfrak{M}_{a,l_{1,k}}(W_{1,k} x_{1,k-1} + B_{1,k}) \\ \mathfrak{M}_{a,l_{2,k}}(W_{2,k} x_{2,k-1} + B_{2,k}) \\ \vdots \\ \mathfrak{M}_{a,l_{n,k}}(W_{n,k} x_{n,k-1} + B_{n,k}) \end{pmatrix} = \begin{pmatrix} x_{1,k} \\ x_{2,k} \\ \vdots \\ x_{n,k} \end{pmatrix} = \mathfrak{x}_k. \tag{2.120}$$

This, (2.53), (2.115), (2.116), (2.117), (2.119), the fact that $\mathfrak{x}_0 = (x_{1,0}, x_{2,0}, \ldots, x_{n,0})$, and the fact that $\mathfrak{x}_{L-1} = (x_{1,L-1}, x_{2,L-1}, \ldots, x_{n,L-1})$ ensure that

$$\big(\mathcal{R}_a\big(\mathbf{P}_n(\Phi)\big)\big)(x_{1,0}, x_{2,0}, \ldots, x_{n,0}) = \big(\mathcal{R}_a\big(\mathbf{P}_n(\Phi)\big)\big)(\mathfrak{x}_0)$$

$$= A_L \mathfrak{x}_{L-1} + b_L = \begin{pmatrix} W_{1,L} x_{1,L-1} + B_{1,L} \\ W_{2,L} x_{2,L-1} + B_{2,L} \\ \vdots \\ W_{n,L} x_{n,L-1} + B_{n,L} \end{pmatrix} = \begin{pmatrix} (\mathcal{R}_a(\Phi_1))(x_{1,0}) \\ (\mathcal{R}_a(\Phi_2))(x_{2,0}) \\ \vdots \\ (\mathcal{R}_a(\Phi_n))(x_{n,0}) \end{pmatrix}. \tag{2.121}$$

This establishes item (ii). The proof of Proposition 2.2.13 is thus complete. $\qquad \square$

**Proposition 2.2.14.** *Let* $n, L \in \mathbb{N}$, $\Phi_1, \Phi_2, \ldots, \Phi_n \in \mathbf{N}$ *satisfy* $L = \mathcal{L}(\Phi_1) = \mathcal{L}(\Phi_2) = \ldots = \mathcal{L}(\Phi_n)$ *(cf. Definition 2.2.1). Then*

*(i) it holds for all* $k \in \mathbb{N}_0$ *that*

$$\mathbb{D}_k(\mathbf{P}_n(\Phi_1, \Phi_2, \ldots, \Phi_n)) = \mathbb{D}_k(\Phi_1) + \mathbb{D}_k(\Phi_2) + \ldots + \mathbb{D}_k(\Phi_n), \tag{2.122}$$

*(ii) it holds that*

$$\mathcal{D}\big(\mathbf{P}_n(\Phi_1, \Phi_2, \ldots, \Phi_n)\big) = \mathcal{D}(\Phi_1) + \mathcal{D}(\Phi_2) + \cdots + \mathcal{D}(\Phi_n), \tag{2.123}$$

*and*

*(iii) it holds that*

$$\mathcal{P}\big(\mathbf{P}_n(\Phi_1, \Phi_2, \ldots, \Phi_n)\big) \leq \tfrac{1}{2}\big[\textstyle\sum_{j=1}^n \mathcal{P}(\Phi_j)\big]^2 \tag{2.124}$$

*(cf. Definition 2.2.11).*

*Proof of Proposition 2.2.14.* Throughout this proof let $l_{j,0}, l_{j,1}, \ldots, l_{j,L} \in \mathbb{N}$, $j \in \{1, 2, \ldots, n\}$, satisfy for all $j \in \{1, 2, \ldots, n\}$, $k \in \{0, 1, \ldots, L\}$ that $l_{j,k} = \mathbb{D}_k(\Phi_j)$. Note that

Lemma 2.2.12 establishes item (i). In addition, observe that item (i) implies item (ii). Moreover, note that item (i) demonstrates that

$$
\begin{aligned}
\mathcal{P}(\mathbf{P}_n(\Phi_1, \Phi_2, \ldots, \Phi_n)) &= \sum_{k=1}^{L} \left[ \textstyle\sum_{i=1}^{n} l_{i,k} \right] \left[ \left( \textstyle\sum_{i=1}^{n} l_{i,k-1} \right) + 1 \right] \\
&= \sum_{k=1}^{L} \left[ \textstyle\sum_{i=1}^{n} l_{i,k} \right] \left[ \left( \textstyle\sum_{j=1}^{n} l_{j,k-1} \right) + 1 \right] \\
&\leq \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{L} l_{i,k}(l_{j,k-1} + 1) \leq \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{k,\ell=1}^{L} l_{i,k}(l_{j,\ell-1} + 1) \\
&= \sum_{i=1}^{n} \sum_{j=1}^{n} \left[ \textstyle\sum_{k=1}^{L} l_{i,k} \right] \left[ \textstyle\sum_{\ell=1}^{L} (l_{j,\ell-1} + 1) \right] \\
&\leq \sum_{i=1}^{n} \sum_{j=1}^{n} \left[ \textstyle\sum_{k=1}^{L} \tfrac{1}{2} l_{i,k}(l_{i,k-1} + 1) \right] \left[ \textstyle\sum_{\ell=1}^{L} l_{j,\ell}(l_{j,\ell-1} + 1) \right] \\
&= \sum_{i=1}^{n} \sum_{j=1}^{n} \tfrac{1}{2} \mathcal{P}(\Phi_i) \mathcal{P}(\Phi_j) = \tfrac{1}{2} \left[ \textstyle\sum_{i=1}^{n} \mathcal{P}(\Phi_i) \right]^2.
\end{aligned}
\tag{2.125}
$$

The proof of Proposition 2.2.14 is thus complete. □

**Corollary 2.2.15.** *Let $n \in \mathbb{N}$, $\Phi = (\Phi_1, \Phi_2, \ldots, \Phi_n) \in \mathbf{N}^n$ satisfy that $\mathcal{D}(\Phi_1) = \mathcal{D}(\Phi_2) = \ldots = \mathcal{D}(\Phi_n)$ (cf. Definition 2.2.1). Then it holds that $\mathcal{P}(\mathbf{P}_n(\Phi)) \leq n^2 \mathcal{P}(\Phi_1)$ (cf. Definition 2.2.11).*

*Proof of Corollary 2.2.15.* Throughout this proof let $L \in \mathbb{N}$, $l_0, l_1, \ldots, l_L \in \mathbb{N}$ satisfy that $\mathcal{D}(\Phi_1) = (l_0, l_1, \ldots, l_L)$. Note that item (ii) in Proposition 2.2.14 and the fact that $\forall\, j \in \{1, 2, \ldots, n\} \colon \mathcal{D}(\Phi_j) = (l_0, l_1, \ldots, l_L)$ demonstrate that

$$
\begin{aligned}
\mathcal{P}(\mathbf{P}_n(\Phi_1, \Phi_2, \ldots, \Phi_n)) &= \sum_{j=1}^{L} (n l_j)\big((n l_{j-1}) + 1\big) \leq \sum_{j=1}^{L} (n l_j)\big((n l_{j-1}) + n\big) \\
&= n^2 \left[ \sum_{j=1}^{L} l_j (l_{j-1} + 1) \right] = n^2 \mathcal{P}(\Phi_1).
\end{aligned}
\tag{2.126}
$$

The proof of Corollary 2.2.15 is thus complete. □

### 2.2.4.2 Parallelizations of ANNs with different lengths

**Definition 2.2.16** (Parallelization of ANNs with different length). *Let $n \in \mathbb{N}$, $\Psi = (\Psi_1, \Psi_2, \ldots, \Psi_n) \in \mathbf{N}^n$ satisfy for all $j \in \{1, 2, \ldots, n\}$ that $\mathcal{H}(\Psi_j) = 1$ and $\mathcal{I}(\Psi_j) = \mathcal{O}(\Psi_j)$. Then we denote by*

$$
\mathrm{P}_{n,\Psi} \colon \{(\Phi_1, \Phi_2, \ldots, \Phi_n) \in \mathbf{N}^n \colon (\forall\, j \in \{1, 2, \ldots, n\} \colon \mathcal{O}(\Phi_j) = \mathcal{I}(\Psi_j))\} \to \mathbf{N} \tag{2.127}
$$

*the function which satisfies for all $\Phi = (\Phi_1, \Phi_2, \ldots, \Phi_n) \in \mathbf{N}^n$ with $\forall\, j \in \{1, 2, \ldots, n\} \colon \mathcal{O}(\Phi_j) = \mathcal{I}(\Psi_j)$ that*

$$
\mathrm{P}_{n,\Psi}(\Phi) = \mathbf{P}_n\big(\mathcal{E}_{\max_{k \in \{1,2,\ldots,n\}} \mathcal{L}(\Phi_k), \Psi_1}(\Phi_1), \ldots, \mathcal{E}_{\max_{k \in \{1,2,\ldots,n\}} \mathcal{L}(\Phi_k), \Psi_n}(\Phi_n)\big) \tag{2.128}
$$

*(cf. Definitions 2.2.1, 2.2.11, and 16.2.1 and Lemma 16.2.2).*

**Corollary 2.2.17.** *Let* $a \in C(\mathbb{R}, \mathbb{R})$, $n \in \mathbb{N}$, $\mathbb{I} = (\mathbb{I}_1, \mathbb{I}_2, \ldots, \mathbb{I}_n)$, $\Phi = (\Phi_1, \Phi_2, \ldots, \Phi_n) \in$ $\mathbf{N}^n$ *satisfy for all* $j \in \{1, 2, \ldots, n\}$, $x \in \mathbb{R}^{\mathcal{O}(\Phi_j)}$ *that* $\mathcal{H}(\mathbb{I}_j) = 1$, $\mathcal{I}(\mathbb{I}_j) = \mathcal{O}(\mathbb{I}_j) = \mathcal{O}(\Phi_j)$, *and* $(\mathcal{R}_a(\mathbb{I}_j))(x) = x$ *(cf. Definitions 2.2.1 and 2.2.3). Then*

*(i) it holds that*

$$\mathcal{R}_a\big(\mathrm{P}_{n,\mathbb{I}}(\Phi)\big) \in C\big(\mathbb{R}^{[\sum_{j=1}^n \mathcal{I}(\Phi_j)]}, \mathbb{R}^{[\sum_{j=1}^n \mathcal{O}(\Phi_j)]}\big) \tag{2.129}$$

*and*

*(ii) it holds for all* $x_1 \in \mathbb{R}^{\mathcal{I}(\Phi_1)}, x_2 \in \mathbb{R}^{\mathcal{I}(\Phi_2)}, \ldots, x_n \in \mathbb{R}^{\mathcal{I}(\Phi_n)}$ *that*

$$\begin{aligned}
&\big(\mathcal{R}_a(\mathrm{P}_{n,\mathbb{I}}(\Phi))\big)(x_1, x_2, \ldots, x_n) \\
&= \big((\mathcal{R}_a(\Phi_1))(x_1), (\mathcal{R}_a(\Phi_2))(x_2), \ldots, (\mathcal{R}_a(\Phi_n))(x_n)\big) \in \mathbb{R}^{[\sum_{j=1}^n \mathcal{O}(\Phi_j)]}
\end{aligned} \tag{2.130}$$

*(cf. Definition 2.2.16).*

*Proof of Corollary 2.2.17.* Throughout this proof let $L \in \mathbb{N}$ satisfy $L = \max_{j \in \{1,2,\ldots,n\}}$ $\mathcal{L}(\Phi_j)$. Note that item (ii) in Lemma 16.2.2, the assumption that for all $j \in \{1, 2, \ldots, n\}$ it holds that $\mathcal{H}(\mathbb{I}_j) = 1$, (16.5), (2.62), and item (ii) in Lemma 16.2.3 demonstrate

(I) that for all $j \in \{1, 2, \ldots, n\}$ it holds that $\mathcal{L}(\mathcal{E}_{L,\mathbb{I}_j}(\Phi_j)) = L$ and $\mathcal{R}_a(\mathcal{E}_{L,\mathbb{I}_j}(\Phi_j)) \in$ $C(\mathbb{R}^{\mathcal{I}(\Phi_j)}, \mathbb{R}^{\mathcal{O}(\Phi_j)})$ and

(II) that for all $j \in \{1, 2, \ldots, n\}$, $x \in \mathbb{R}^{\mathcal{I}(\Phi_j)}$ it holds that

$$\big(\mathcal{R}_a(\mathcal{E}_{L,\mathbb{I}_j}(\Phi_j))\big)(x) = (\mathcal{R}_a(\Phi_j))(x) \tag{2.131}$$

(cf. Definition 16.2.1). Items (i)–(ii) in Proposition 2.2.13 therefore imply

(A) that

$$\mathcal{R}_a\big(\mathbf{P}_n\big(\mathcal{E}_{L,\mathbb{I}_1}(\Phi_1), \mathcal{E}_{L,\mathbb{I}_2}(\Phi_2), \ldots, \mathcal{E}_{L,\mathbb{I}_n}(\Phi_n)\big)\big) \in C\big(\mathbb{R}^{[\sum_{j=1}^n \mathcal{I}(\Phi_j)]}, \mathbb{R}^{[\sum_{j=1}^n \mathcal{O}(\Phi_j)]}\big) \tag{2.132}$$

and

(B) that for all $x_1 \in \mathbb{R}^{\mathcal{I}(\Phi_1)}, x_2 \in \mathbb{R}^{\mathcal{I}(\Phi_2)}, \ldots, x_n \in \mathbb{R}^{\mathcal{I}(\Phi_n)}$ it holds that

$$\begin{aligned}
&\big(\mathcal{R}_a\big(\mathbf{P}_n\big(\mathcal{E}_{L,\mathbb{I}_1}(\Phi_1), \mathcal{E}_{L,\mathbb{I}_2}(\Phi_2), \ldots, \mathcal{E}_{L,\mathbb{I}_n}(\Phi_n)\big)\big)\big)(x_1, x_2, \ldots, x_n) \\
&= \Big((\mathcal{R}_a(\mathcal{E}_{L,\mathbb{I}_1}(\Phi_1)))(x_1), (\mathcal{R}_a(\mathcal{E}_{L,\mathbb{I}_2}(\Phi_2)))(x_2), \ldots, (\mathcal{R}_a(\mathcal{E}_{L,\mathbb{I}_n}(\Phi_n)))(x_n)\Big) \\
&= \Big((\mathcal{R}_a(\Phi_1))(x_1), (\mathcal{R}_a(\Phi_2))(x_2), \ldots, (\mathcal{R}_a(\Phi_n))(x_n)\Big)
\end{aligned} \tag{2.133}$$

(cf. Definition 2.2.11). Combining this with (2.128) and the fact that $L = \max_{j \in \{1,2,\ldots,n\}}$ $\mathcal{L}(\Phi_j)$ ensures

(C) that

$$\mathcal{R}_a\big(\mathrm{P}_{n,\mathbb{I}}(\Phi)\big) \in C\big(\mathbb{R}^{[\sum_{j=1}^n \mathcal{I}(\Phi_j)]}, \mathbb{R}^{[\sum_{j=1}^n \mathcal{O}(\Phi_j)]}\big) \tag{2.134}$$

and

(D) that for all $x_1 \in \mathbb{R}^{\mathcal{I}(\Phi_1)}, x_2 \in \mathbb{R}^{\mathcal{I}(\Phi_2)}, \ldots, x_n \in \mathbb{R}^{\mathcal{I}(\Phi_n)}$ it holds that

$$\begin{aligned}
&\big(\mathcal{R}_a\big(\mathrm{P}_{n,\mathbb{I}}(\Phi)\big)\big)(x_1, x_2, \ldots, x_n) \\
&= \big(\mathcal{R}_a\big(\mathbf{P}_n\big(\mathcal{E}_{L,\mathbb{I}_1}(\Phi_1), \mathcal{E}_{L,\mathbb{I}_2}(\Phi_2), \ldots, \mathcal{E}_{L,\mathbb{I}_n}(\Phi_n)\big)\big)\big)(x_1, x_2, \ldots, x_n) \\
&= \Big((\mathcal{R}_a(\Phi_1))(x_1), (\mathcal{R}_a(\Phi_2))(x_2), \ldots, (\mathcal{R}_a(\Phi_n))(x_n)\Big).
\end{aligned} \tag{2.135}$$

This establishes items (i)–(ii). The proof of Corollary 2.2.17 is thus complete. $\square$

## 2.2.5 Representations of the identities with rectifier functions

**Definition 2.2.18.** *We denote by $\mathfrak{I}_d \in \mathbf{N}$, $n \in \mathbb{N}$, the neural networks which satisfy for all $d \in \mathbb{N}$ that*

$$\mathfrak{I}_1 = \left( \left( \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \end{pmatrix} \right), \left( \begin{pmatrix} 1 & -1 \end{pmatrix}, 0 \right) \right) \in \left( (\mathbb{R}^{2\times 1} \times \mathbb{R}^2) \times (\mathbb{R}^{1\times 2} \times \mathbb{R}^1) \right) \qquad (2.136)$$

*and*

$$\mathfrak{I}_d = \mathbf{P}_d(\mathfrak{I}_1, \mathfrak{I}_1, \dots, \mathfrak{I}_1) \qquad (2.137)$$

*(cf. Definitions 2.2.1 and 2.2.11).*

**Lemma 2.2.19.** *Let $d \in \mathbb{N}$. Then*

*(i) it holds that $\mathcal{D}(\mathfrak{I}_d) = (d, 2d, d) \in \mathbb{N}^3$,*

*(ii) it holds that $\mathcal{R}_{\mathfrak{r}}(\mathfrak{I}_d) \in C(\mathbb{R}^d, \mathbb{R}^d)$, and*

*(iii) it holds for all $x \in \mathbb{R}^d$ that*

$$(\mathcal{R}_{\mathfrak{r}}(\mathfrak{I}_d))(x) = x \qquad (2.138)$$

*(cf. Definitions 2.2.1, 2.2.3, and 2.2.18).*

*Proof of Lemma 2.2.19.* Throughout this proof let $L = 2$, $l_0 = 1$, $l_1 = 2$, $l_2 = 1$. Note that (2.136) ensures that

$$\mathcal{D}(\mathfrak{I}_1) = (1, 2, 1) = (l_0, l_1, l_2). \qquad (2.139)$$

This and Lemma 2.2.12 prove that

$$\begin{aligned} &\mathbf{P}_d(\mathfrak{I}_1, \mathfrak{I}_1, \dots, \mathfrak{I}_1) \\ &\in \left( \bigtimes_{k=1}^{L} \left( \mathbb{R}^{(dl_k)\times(dl_{k-1})} \times \mathbb{R}^{(dl_k)} \right) \right) = \left( \left( \mathbb{R}^{(2d)\times d} \times \mathbb{R}^{2d} \right) \times \left( \mathbb{R}^{d\times(2d)} \times \mathbb{R}^d \right) \right) \end{aligned} \qquad (2.140)$$

(cf. Definition 2.2.11). Hence, we obtain that $\mathcal{D}(\mathfrak{I}_d) = (d, 2d, d) \in \mathbb{N}^3$. This establishes item (i). Next note that (2.136) assures that for all $x \in \mathbb{R}$ it holds that

$$(\mathcal{R}_{\mathfrak{r}}(\mathfrak{I}_1))(x) = \mathfrak{r}(x) - \mathfrak{r}(-x) = \max\{x, 0\} - \max\{-x, 0\} = x. \qquad (2.141)$$

Combining this and Proposition 2.2.13 demonstrates that for all $x = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$ it holds that $\mathcal{R}_{\mathfrak{r}}(\mathfrak{I}_d) \in C(\mathbb{R}^d, \mathbb{R}^d)$ and

$$\begin{aligned} (\mathcal{R}_{\mathfrak{r}}(\mathfrak{I}_d))(x) &= \left( \mathcal{R}_{\mathfrak{r}}\left( \mathbf{P}_d(\mathfrak{I}_1, \mathfrak{I}_1, \dots, \mathfrak{I}_1) \right) \right)(x_1, x_2, \dots, x_d) \\ &= \left( (\mathcal{R}_{\mathfrak{r}}(\mathfrak{I}_1))(x_1), (\mathcal{R}_{\mathfrak{r}}(\mathfrak{I}_1))(x_2), \dots, (\mathcal{R}_{\mathfrak{r}}(\mathfrak{I}_1))(x_d) \right) \\ &= (x_1, x_2, \dots, x_d) = x. \end{aligned} \qquad (2.142)$$

This establishes items (ii)–(iii). The proof of Lemma 2.2.19 is thus complete. $\qquad\square$

## 2.2.6 Scalar multiplications of ANNs

### 2.2.6.1 Affine transformations as ANNs

**Definition 2.2.20** (Affine linear transformation ANN)**.** *Let* $m, n \in \mathbb{N}$, $W \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^m$*. Then we denote by* $\mathbf{A}_{W,B} \in (\mathbb{R}^{m \times n} \times \mathbb{R}^m) \subseteq \mathbf{N}$ *the neural network given by* $\mathbf{A}_{W,B} = (W, B)$ *(cf. Definitions 2.2.1 and 2.2.2).*

**Lemma 2.2.21.** *Let* $m, n \in \mathbb{N}$, $W \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^m$*. Then*

   *(i) it holds that* $\mathcal{D}(\mathbf{A}_{W,B}) = (n, m) \in \mathbb{N}^2$,

   *(ii) it holds for all* $a \in C(\mathbb{R}, \mathbb{R})$ *that* $\mathcal{R}_a(\mathbf{A}_{W,B}) \in C(\mathbb{R}^n, \mathbb{R}^m)$*, and*

   *(iii) it holds for all* $a \in C(\mathbb{R}, \mathbb{R})$, $x \in \mathbb{R}^n$ *that* $(\mathcal{R}_a(\mathbf{A}_{W,B}))(x) = Wx + B$

*(cf. Definitions 2.2.1, 2.2.3, and 2.2.20).*

*Proof of Lemma 2.2.21.* Note the fact that $\mathbf{A}_{W,B} \in (\mathbb{R}^{m \times n} \times \mathbb{R}^m) \subseteq \mathbf{N}$ ensures that $\mathcal{D}(\mathbf{A}_{W,B}) = (n, m) \in \mathbb{N}^2$. This establishes item (i). Next observe that the fact that $\mathbf{A}_{W,B} = (W, B) \in (\mathbb{R}^{m \times n} \times \mathbb{R}^m)$ and (2.53) prove that for all $a \in C(\mathbb{R}, \mathbb{R})$, $x \in \mathbb{R}^n$ it holds that $\mathcal{R}_a(\mathbf{A}_{W,B}) \in C(\mathbb{R}^n, \mathbb{R}^m)$ and

$$(\mathcal{R}_a(\mathbf{A}_{W,B}))(x) = Wx + B. \tag{2.143}$$

This establishes items (ii) and (iii). The proof of Lemma 2.2.21 is thus complete. □

**Lemma 2.2.22.** *Let* $\Phi \in \mathbf{N}$ *(cf. Definition 2.2.1). Then*

   *(i) it holds for all* $m \in \mathbb{N}$, $W \in \mathbb{R}^{m \times \mathcal{O}(\Phi)}$, $B \in \mathbb{R}^m$ *that*

$$\mathcal{D}(\mathbf{A}_{W,B} \bullet \Phi) = (\mathbb{D}_0(\Phi), \mathbb{D}_1(\Phi), \ldots, \mathbb{D}_{\mathcal{H}(\Phi)}(\Phi), m), \tag{2.144}$$

   *(ii) it holds for all* $a \in C(\mathbb{R}, \mathbb{R})$, $m \in \mathbb{N}$, $W \in \mathbb{R}^{m \times \mathcal{O}(\Phi)}$, $B \in \mathbb{R}^m$ *that* $\mathcal{R}_a(\mathbf{A}_{W,B} \bullet \Phi) \in C(\mathbb{R}^{\mathcal{I}(\Phi)}, \mathbb{R}^m)$,

   *(iii) it holds for all* $a \in C(\mathbb{R}, \mathbb{R})$, $m \in \mathbb{N}$, $W \in \mathbb{R}^{m \times \mathcal{O}(\Phi)}$, $B \in \mathbb{R}^m$, $x \in \mathbb{R}^{\mathcal{I}(\Phi)}$ *that*

$$(\mathcal{R}_a(\mathbf{A}_{W,B} \bullet \Phi))(x) = W((\mathcal{R}_a(\Phi))(x)) + B, \tag{2.145}$$

   *(iv) it holds for all* $n \in \mathbb{N}$, $W \in \mathbb{R}^{\mathcal{I}(\Phi) \times n}$, $B \in \mathbb{R}^{\mathcal{I}(\Phi)}$ *that*

$$\mathcal{D}(\Phi \bullet \mathbf{A}_{W,B}) = (n, \mathbb{D}_1(\Phi), \mathbb{D}_2(\Phi), \ldots, \mathbb{D}_{\mathcal{L}(\Phi)}(\Phi)), \tag{2.146}$$

   *(v) it holds for all* $a \in C(\mathbb{R}, \mathbb{R})$, $n \in \mathbb{N}$, $W \in \mathbb{R}^{\mathcal{I}(\Phi) \times n}$, $B \in \mathbb{R}^{\mathcal{I}(\Phi)}$ *that* $\mathcal{R}_a(\Phi \bullet \mathbf{A}_{W,B}) \in C(\mathbb{R}^n, \mathbb{R}^{\mathcal{O}(\Phi)})$*, and*

   *(vi) it holds for all* $a \in C(\mathbb{R}, \mathbb{R})$, $n \in \mathbb{N}$, $W \in \mathbb{R}^{\mathcal{I}(\Phi) \times n}$, $B \in \mathbb{R}^{\mathcal{I}(\Phi)}$, $x \in \mathbb{R}^n$ *that*

$$(\mathcal{R}_a(\Phi \bullet \mathbf{A}_{W,B}))(x) = (\mathcal{R}_a(\Phi))(Wx + B), \tag{2.147}$$

*(cf. Definitions 2.2.3, 2.2.5, and 2.2.20).*

*Proof of Lemma 2.2.22.* Note that Lemma 2.2.21 demonstrates that for all $m, n \in \mathbb{N}$, $W \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^m$, $a \in C(\mathbb{R}, \mathbb{R})$, $x \in \mathbb{R}^n$ it holds that $\mathcal{R}_a(\mathbf{A}_{W,B}) \in C(\mathbb{R}^n, \mathbb{R}^m)$ and

$$(\mathcal{R}_a(\mathbf{A}_{W,B}))(x) = Wx + B. \tag{2.148}$$

Combining this and Proposition 2.2.7 establishes items (i), (ii), (iii), (iv), (v), and (vi). The proof of Lemma 2.2.22 is thus complete. □

**2.2.6.2  Scalar multiplications of ANNs**

**Definition 2.2.23** (Scalar multiplications of ANNs)**.** *We denote by* $(\cdot)\circledast(\cdot)\colon \mathbb{R}\times\mathbf{N}\to\mathbf{N}$ *the function which satisfies for all* $\lambda\in\mathbb{R}$, $\Phi\in\mathbf{N}$ *that*

$$\lambda\circledast\Phi = \mathbf{A}_{\lambda\,\mathrm{I}_{\mathcal{O}(\Phi)},0}\bullet\Phi \tag{2.149}$$

*(cf. Definitions 2.2.1, 2.2.5, 2.2.9, and 2.2.20).*

**Lemma 2.2.24.** *Let* $\lambda\in\mathbb{R}$, $\Phi\in\mathbf{N}$ *(cf. Definition 2.2.1). Then*

(i) *it holds that* $\mathcal{D}(\lambda\circledast\Phi) = \mathcal{D}(\Phi)$,

(ii) *it holds for all* $a\in C(\mathbb{R},\mathbb{R})$ *that* $\mathcal{R}_a(\lambda\circledast\Phi)\in C(\mathbb{R}^{\mathcal{I}(\Phi)},\mathbb{R}^{\mathcal{O}(\Phi)})$, *and*

(iii) *it holds for all* $a\in C(\mathbb{R},\mathbb{R})$, $x\in\mathbb{R}^{\mathcal{I}(\Phi)}$ *that*

$$(\mathcal{R}_a(\lambda\circledast\Phi))(x) = \lambda\big((\mathcal{R}_a(\Phi))(x)\big) \tag{2.150}$$

*(cf. Definitions 2.2.3 and 2.2.23).*

*Proof of Lemma 2.2.24.* Throughout this proof let $L\in\mathbb{N}$, $l_0,l_1,\ldots,l_L\in\mathbb{N}$ satisfy $L=\mathcal{L}(\Phi)$ and $(l_0,l_1,\ldots,l_L)=\mathcal{D}(\Phi)$. Note that item (i) in Lemma 2.2.21 proves that

$$\mathcal{D}(\mathbf{A}_{\lambda\,\mathrm{I}_{\mathcal{O}(\Phi)},0}) = (\mathcal{O}(\Phi),\mathcal{O}(\Phi)) \tag{2.151}$$

(cf. Definitions 2.2.9 and 2.2.20). Combining this and item (i) in Lemma 2.2.22 assures that

$$\mathcal{D}(\lambda\circledast\Phi) = \mathcal{D}(\mathbf{A}_{\lambda\,\mathrm{I}_{\mathcal{O}(\Phi)},0}\bullet\Phi) = (l_0,l_1,\ldots,l_{L-1},\mathcal{O}(\Phi)) = \mathcal{D}(\Phi). \tag{2.152}$$

This establishes item (i). Moreover, observe that items (ii)–(iii) in Lemma 2.2.22 demonstrate that for all $a\in C(\mathbb{R},\mathbb{R})$, $x\in\mathbb{R}^{\mathcal{I}(\Phi)}$ it holds that $\mathcal{R}_a(\lambda\circledast\Phi)\in C(\mathbb{R}^{\mathcal{I}(\Phi)},\mathbb{R}^{\mathcal{O}(\Phi)})$ and

$$\begin{aligned}(\mathcal{R}_a(\lambda\circledast\Phi))(x) &= (\mathcal{R}_a(\mathbf{A}_{\lambda\,\mathrm{I}_{\mathcal{O}(\Phi)},0}\bullet\Phi))(x)\\ &= \lambda\,\mathrm{I}_{\mathcal{O}(\Phi)}\big((\mathcal{R}_a(\Phi))(x)\big) = \lambda\big((\mathcal{R}_a(\Phi))(x)\big).\end{aligned} \tag{2.153}$$

This establishes items (ii)–(iii). The proof of Lemma 2.2.24 is thus complete. □

## 2.2.7  Sums of ANNs with the same length

### 2.2.7.1  Sums of vectors as neural networks

**Definition 2.2.25.** *Let* $m,n\in\mathbb{N}$. *Then we denote by* $\mathbb{S}_{m,n}\in(\mathbb{R}^{m\times(mn)}\times\mathbb{R}^m)$ *the neural network given by*

$$\mathbb{S}_{m,n} = \mathbf{A}_{(\mathrm{I}_m\ \mathrm{I}_m\ \ldots\ \mathrm{I}_m),0} \tag{2.154}$$

*(cf. Definitions 2.2.9 and 2.2.20).*

**Lemma 2.2.26.** *Let* $m,n\in\mathbb{N}$. *Then*

(i) *it holds that* $\mathbb{S}_{m,n}\in\mathbf{N}$,

(ii) *it holds that* $\mathcal{D}(\mathbb{S}_{m,n}) = (mn,m)\in\mathbb{N}^2$,

*(iii)  it holds for all $a \in C(\mathbb{R},\mathbb{R})$ that $\mathcal{R}_a(\mathbb{S}_{m,n}) \in C(\mathbb{R}^{mn}, \mathbb{R}^m)$, and*

*(iv)  it holds for all $a \in C(\mathbb{R},\mathbb{R})$, $x_1, x_2, \ldots, x_n \in \mathbb{R}^m$ that*

$$(\mathcal{R}_a(\mathbb{S}_{m,n}))(x_1, x_2, \ldots, x_n) = \sum_{k=1}^{n} x_k \tag{2.155}$$

*(cf. Definitions 2.2.1, 2.2.3, and 2.2.25).*

*Proof of Lemma 2.2.26.* Note that the fact that $\mathbb{S}_{m,n} \in (\mathbb{R}^{m \times (mn)} \times \mathbb{R}^m)$ ensures that $\mathbb{S}_{m,n} \in \mathbf{N}$ and $\mathcal{D}(\mathbb{S}_{m,n}) = (mn, m) \in \mathbb{N}^2$. This establishes items (i) and (ii). Next observe that items (ii) and (iii) in Lemma 2.2.21 prove that for all $a \in C(\mathbb{R},\mathbb{R})$, $x_1, x_2, \ldots, x_n \in \mathbb{R}^m$ it holds that $\mathcal{R}_a(\mathbb{S}_{m,n}) \in C(\mathbb{R}^{mn}, \mathbb{R}^m)$ and

$$\begin{aligned}(\mathcal{R}_a(\mathbb{S}_{m,n}))(x_1, x_2, \ldots, x_n) &= \big(\mathcal{R}_a\big(\mathbf{A}_{(\mathrm{I}_m \ \mathrm{I}_m \ \ldots \ \mathrm{I}_m),0}\big)\big)(x_1, x_2, \ldots, x_n) \\ &= (\mathrm{I}_m \ \ \mathrm{I}_m \ \ \ldots \ \ \mathrm{I}_m)(x_1, x_2, \ldots, x_n) = \sum_{k=1}^{n} x_k\end{aligned} \tag{2.156}$$

*(cf. Definition 2.2.9 and Definition 2.2.20).* This establishes items (iii) and (iv). The proof of Lemma 2.2.26 is thus complete. $\qquad \square$

**Lemma 2.2.27.** *Let $m, n \in \mathbb{N}$, $a \in C(\mathbb{R},\mathbb{R})$, $\Phi \in \mathbf{N}$ satisfy $\mathcal{O}(\Phi) = nm$ (cf. Definition 2.2.1). Then*

*(i)  it holds that $\mathcal{R}_a(\mathbb{S}_{m,n} \bullet \Phi) \in C(\mathbb{R}^{\mathcal{I}(\Phi)}, \mathbb{R}^m)$ and*

*(ii)  it holds for all $x \in \mathbb{R}^{\mathcal{I}(\Phi)}$, $y_1, y_2, \ldots, y_n \in \mathbb{R}^m$ with $(\mathcal{R}_a(\Phi))(x) = (y_1, y_2, \ldots, y_n)$ that*

$$\big(\mathcal{R}_a(\mathbb{S}_{m,n} \bullet \Phi)\big)(x) = \sum_{k=1}^{n} y_k \tag{2.157}$$

*(cf. Definitions 2.2.3, 2.2.5, and 2.2.25).*

*Proof of Lemma 2.2.27.* Note that Lemma 2.2.26 ensures that for all $x_1, x_2, \ldots, x_n \in \mathbb{R}^m$ it holds that $\mathcal{R}_a(\mathbb{S}_{m,n}) \in C(\mathbb{R}^{nm}, \mathbb{R}^m)$ and

$$(\mathcal{R}_a(\mathbb{S}_{m,n}))(x_1, x_2, \ldots, x_n) = \sum_{k=1}^{n} x_k. \tag{2.158}$$

Combining this and item (v) in Proposition 2.2.7 establishes items (i)–(ii). The proof of Lemma 2.2.27 is thus complete. $\qquad \square$

**Lemma 2.2.28.** *Let $n \in \mathbb{N}$, $a \in C(\mathbb{R},\mathbb{R})$, $\Phi \in \mathbf{N}$ (cf. Definition 2.2.1). Then*

*(i)  it holds that $\mathcal{R}_a(\Phi \bullet \mathbb{S}_{\mathcal{I}(\Phi),n}) \in C(\mathbb{R}^{n\mathcal{I}(\Phi)}, \mathbb{R}^{\mathcal{O}(\Phi)})$ and*

*(ii)  it holds for all $x_1, x_2, \ldots, x_n \in \mathbb{R}^{\mathcal{I}(\Phi)}$ that*

$$\big(\mathcal{R}_a(\Phi \bullet \mathbb{S}_{\mathcal{I}(\Phi),n})\big)(x_1, x_2, \ldots, x_n) = (\mathcal{R}_a(\Phi))\left(\sum_{k=1}^{n} x_k\right) \tag{2.159}$$

*(cf. Definitions 2.2.3, 2.2.5, and 2.2.25).*

*Proof of Lemma 2.2.28.* Note that Lemma 2.2.26 demonstrates that for all $m \in \mathbb{N}$, $x_1, x_2, \ldots, x_n \in \mathbb{R}^m$ it holds that $\mathcal{R}_a(\mathbb{S}_{m,n}) \in C(\mathbb{R}^{mn}, \mathbb{R}^m)$ and

$$(\mathcal{R}_a(\mathbb{S}_{m,n}))(x_1, x_2, \ldots, x_n) = \sum_{k=1}^{n} x_k. \tag{2.160}$$

Combining this and item (v) in Proposition 2.2.7 establishes items (i) and (ii). The proof of Lemma 2.2.28 is thus complete. $\qquad \square$

### 2.2.7.2   Concatenation of vectors as neural networks

**Definition 2.2.29.** *Let $m, n \in \mathbb{N}$, $A \in \mathbb{R}^{m \times n}$. Then we denote by $A^* \in \mathbb{R}^{n \times m}$ the transpose of $A$.*

**Definition 2.2.30.** *Let $m, n \in \mathbb{N}$. Then we denote by $\mathbb{T}_{m,n} \in (\mathbb{R}^{(mn) \times m} \times \mathbb{R}^{mn})$ the neural network given by*

$$\mathbb{T}_{m,n} = \mathbf{A}_{(\mathrm{I}_m \ \mathrm{I}_m \ \dots \ \mathrm{I}_m)^*,0} \tag{2.161}$$

*(cf. Definitions 2.2.9, 2.2.20, and 2.2.29).*

**Lemma 2.2.31.** *Let $m, n \in \mathbb{N}$. Then*

(i) *it holds that $\mathbb{T}_{m,n} \in \mathbf{N}$,*

(ii) *it holds that $\mathcal{D}(\mathbb{T}_{m,n}) = (m, mn) \in \mathbb{N}^2$,*

(iii) *it holds for all $a \in C(\mathbb{R}, \mathbb{R})$ that $\mathcal{R}_a(\mathbb{T}_{m,n}) \in C(\mathbb{R}^m, \mathbb{R}^{mn})$, and*

(iv) *it holds for all $a \in C(\mathbb{R}, \mathbb{R})$, $x \in \mathbb{R}^m$ that*

$$(\mathcal{R}_a(\mathbb{T}_{m,n}))(x) = (x, x, \dots, x) \tag{2.162}$$

*(cf. Definitions 2.2.1, 2.2.3, and 2.2.30).*

*Proof of Lemma 2.2.31.* Note that the fact that $\mathbb{T}_{m,n} \in (\mathbb{R}^{(mn) \times m} \times \mathbb{R}^{mn})$ ensures that $\mathbb{T}_{m,n} \in \mathbf{N}$ and $\mathcal{D}(\mathbb{T}_{m,n}) = (m, mn) \in \mathbb{N}^2$. This establishes items (i)–(ii). Next observe that items (v)–(vi) in Lemma 2.2.22 prove that for all $a \in C(\mathbb{R}, \mathbb{R})$, $x \in \mathbb{R}^m$ it holds that $\mathcal{R}_a(\mathbb{T}_{m,n}) \in C(\mathbb{R}^m, \mathbb{R}^{mn})$ and

$$\begin{aligned}
(\mathcal{R}_a(\mathbb{T}_{m,n}))(x) &= \left(\mathcal{R}_a\left(\mathbf{A}_{(\mathrm{I}_m \ \mathrm{I}_m \ \dots \ \mathrm{I}_m)^*,0}\right)\right)(x) \\
&= (\mathrm{I}_m \ \mathrm{I}_m \ \dots \ \mathrm{I}_m)^* x = (x, x, \dots, x)
\end{aligned} \tag{2.163}$$

(cf. Definitions 2.2.9 and 2.2.20). This establishes items (iii) and (iv). The proof of Lemma 2.2.31 is thus complete. $\square$

**Lemma 2.2.32.** *Let $n \in \mathbb{N}$, $a \in C(\mathbb{R}, \mathbb{R})$, $\Phi \in \mathbf{N}$ (cf. Definition 2.2.1). Then*

(i) *it holds that $\mathcal{R}_a(\mathbb{T}_{\mathcal{O}(\Phi),n} \bullet \Phi) \in C(\mathbb{R}^{\mathcal{I}(\Phi)}, \mathbb{R}^{n\mathcal{O}(\Phi)})$ and*

(ii) *it holds for all $x \in \mathbb{R}^{\mathcal{I}(\Phi)}$ that*

$$\left(\mathcal{R}_a(\mathbb{T}_{\mathcal{O}(\Phi),n} \bullet \Phi)\right)(x) = \left((\mathcal{R}_a(\Phi))(x), (\mathcal{R}_a(\Phi))(x), \dots, (\mathcal{R}_a(\Phi))(x)\right) \tag{2.164}$$

*(cf. Definitions 2.2.3, 2.2.5, and 2.2.30).*

*Proof of Lemma 2.2.32.* Note that Lemma 2.2.31 ensures that for all $m \in \mathbb{N}$, $x \in \mathbb{R}^m$ it holds that $\mathcal{R}_a(\mathbb{T}_{m,n}) \in C(\mathbb{R}^m, \mathbb{R}^{mn})$ and

$$(\mathcal{R}_a(\mathbb{T}_{m,n}))(x) = (x, x, \dots, x). \tag{2.165}$$

Combining this and item (v) in Proposition 2.2.7 establishes items (i) and (ii). The proof of Lemma 2.2.32 is thus complete. $\square$

**Lemma 2.2.33.** *Let $m, n \in \mathbb{N}$, $a \in C(\mathbb{R}, \mathbb{R})$, $\Phi \in \mathbf{N}$ satisfy $\mathcal{I}(\Phi) = mn$ (cf. Definition 2.2.1). Then*

(i) *it holds that $\mathcal{R}_a(\Phi \bullet \mathbb{T}_{m,n}) \in C(\mathbb{R}^m, \mathbb{R}^{\mathcal{O}(\Phi)})$ and*

(ii) *it holds for all $x \in \mathbb{R}^m$ that*

$$\big(\mathcal{R}_a(\Phi \bullet \mathbb{T}_{m,n})\big)(x) = (\mathcal{R}_a(\Phi))(x, x, \ldots, x) \tag{2.166}$$

*(cf. Definitions 2.2.3, 2.2.5, and 2.2.30).*

*Proof of Lemma 2.2.33.* Observe that Lemma 2.2.31 demonstrates that for all $x \in \mathbb{R}^m$ it holds that $\mathcal{R}_a(\mathbb{T}_{m,n}) \in C(\mathbb{R}^m, \mathbb{R}^{mn})$ and

$$(\mathcal{R}_a(\mathbb{T}_{m,n}))(x) = (x, x, \ldots, x). \tag{2.167}$$

Combining this and item (v) in Proposition 2.2.7 establishes items (i) and (ii). The proof of Lemma 2.2.33 is thus complete. □

### 2.2.7.3   Sums of ANNs

**Definition 2.2.34** (Sums of ANNs with the same length)**.** *Let $n \in \mathbb{Z}$, $m \in \{n, n+1, \ldots\}$, $\Phi_n, \Phi_{n+1}, \ldots, \Phi_m \in \mathbf{N}$ satisfy for all $k \in \{n, n+1, \ldots, m\}$ that $\mathcal{L}(\Phi_k) = \mathcal{L}(\Phi_n)$, $\mathcal{I}(\Phi_k) = \mathcal{I}(\Phi_n)$, and $\mathcal{O}(\Phi_k) = \mathcal{O}(\Phi_n)$. Then we denote by $\bigoplus_{k=n}^m \Phi_k$ (we denote by $\Phi_n \oplus \Phi_{n+1} \oplus \ldots \oplus \Phi_m$) the neural network given by*

$$\bigoplus_{k=n}^m \Phi_k = \big(\mathbb{S}_{\mathcal{O}(\Phi_n),m-n+1} \bullet \big[\mathbf{P}_{m-n+1}(\Phi_n, \Phi_{n+1}, \ldots, \Phi_m)\big] \bullet \mathbb{T}_{\mathcal{I}(\Phi_n),m-n+1}\big) \in \mathbf{N} \tag{2.168}$$

*(cf. Definitions 2.2.1, 2.2.2, 2.2.5, 2.2.11, 2.2.25, and 2.2.30).*

**Lemma 2.2.35.** *Let $n \in \mathbb{Z}$, $m \in \{n, n+1, \ldots\}$, $\Phi_n, \Phi_{n+1}, \ldots, \Phi_m \in \mathbf{N}$ satisfy for all $k \in \{n, n+1, \ldots, m\}$ that $\mathcal{L}(\Phi_k) = \mathcal{L}(\Phi_n)$, $\mathcal{I}(\Phi_k) = \mathcal{I}(\Phi_n)$, and $\mathcal{O}(\Phi_k) = \mathcal{O}(\Phi_n)$ (cf. Definition 2.2.1). Then*

(i) *it holds that $\mathcal{L}(\bigoplus_{k=n}^m \Phi_k) = \mathcal{L}(\Phi_n)$,*

(ii) *it holds that*

$$\mathcal{D}\bigg(\bigoplus_{k=n}^m \Phi_k\bigg) = \bigg(\mathcal{I}(\Phi_n), \sum_{k=n}^m \mathbb{D}_1(\Phi_k), \sum_{k=n}^m \mathbb{D}_2(\Phi_k), \ldots, \sum_{k=n}^m \mathbb{D}_{\mathcal{H}(\Phi_n)}(\Phi_k), \mathcal{O}(\Phi_n)\bigg), \tag{2.169}$$

(iii) *it holds for all $a \in C(\mathbb{R}, \mathbb{R})$ that $\mathcal{R}_a(\bigoplus_{k=n}^m \Phi_k) \in C(\mathbb{R}^{\mathcal{I}(\Phi_n)}, \mathbb{R}^{\mathcal{O}(\Phi_n)})$, and*

(iv) *it holds for all $a \in C(\mathbb{R}, \mathbb{R})$, $x \in \mathbb{R}^{\mathcal{I}(\Phi_n)}$ that*

$$\bigg(\mathcal{R}_a\bigg(\bigoplus_{k=n}^m \Phi_k\bigg)\bigg)(x) = \sum_{k=n}^m (\mathcal{R}_a(\Phi_k))(x) \tag{2.170}$$

*(cf. Definitions 2.2.3 and 2.2.34).*

*Proof of Lemma 2.2.35.* First, note that Lemma 2.2.12 proves that

$$
\begin{aligned}
&\mathcal{D}\big(\mathbf{P}_{m-n+1}(\Phi_n, \Phi_{n+1}, \ldots, \Phi_m)\big)\\
&= \left(\sum_{k=n}^{m} \mathbb{D}_0(\Phi_k), \sum_{k=n}^{m} \mathbb{D}_1(\Phi_k), \ldots, \sum_{k=n}^{m} \mathbb{D}_{\mathcal{L}(\Phi_n)-1}(\Phi_k), \sum_{k=n}^{m} \mathbb{D}_{\mathcal{L}(\Phi_n)}(\Phi_k)\right)\\
&= \left((m-n+1)\mathcal{I}(\Phi_n), \sum_{k=n}^{m} \mathbb{D}_1(\Phi_k), \sum_{k=n}^{m} \mathbb{D}_2(\Phi_k), \ldots, \right.\\
&\qquad \left. \sum_{k=n}^{m} \mathbb{D}_{\mathcal{L}(\Phi_n)-1}(\Phi_k), (m-n+1)\mathcal{O}(\Phi_n)\right)
\end{aligned}
\tag{2.171}
$$

(cf. Definition 2.2.11). Moreover, observe that item (ii) in Lemma 2.2.26 ensures that

$$
\mathcal{D}\big(\mathbb{S}_{\mathcal{O}(\Phi_n),m-n+1}\big) = ((m-n+1)\mathcal{O}(\Phi_n), \mathcal{O}(\Phi_n))
\tag{2.172}
$$

(cf. Definition 2.2.25). This, (2.171), and item (i) in Proposition 2.2.7 demonstrate that

$$
\begin{aligned}
&\mathcal{D}\big(\mathbb{S}_{\mathcal{O}(\Phi_n),m-n+1} \bullet \big[\mathbf{P}_{m-n+1}(\Phi_n, \Phi_{n+1}, \ldots, \Phi_m)\big]\big)\\
&= \left((m-n+1)\mathcal{I}(\Phi_n), \sum_{k=n}^{m} \mathbb{D}_1(\Phi_k), \sum_{k=n}^{m} \mathbb{D}_2(\Phi_k), \ldots, \sum_{k=n}^{m} \mathbb{D}_{\mathcal{L}(\Phi_n)-1}(\Phi_k), \mathcal{O}(\Phi_n)\right).
\end{aligned}
\tag{2.173}
$$

Next note that item (ii) in Lemma 2.2.31 assures that

$$
\mathcal{D}\big(\mathbb{T}_{\mathcal{I}(\Phi_n),m-n+1}\big) = (\mathcal{I}(\Phi_n), (m-n+1)\mathcal{I}(\Phi_n))
\tag{2.174}
$$

(cf. Definition 2.2.30). Combining this, (2.173), and, item (i) in Proposition 2.2.7 proves that

$$
\begin{aligned}
&\mathcal{D}\left(\bigoplus_{k=n}^{m} \Phi_k\right)\\
&= \mathcal{D}\big(\mathbb{S}_{\mathcal{O}(\Phi_n),(m-n+1)} \bullet \big[\mathbf{P}_{m-n+1}(\Phi_n, \Phi_{n+1}, \ldots, \Phi_m)\big] \bullet \mathbb{T}_{\mathcal{I}(\Phi_n),(m-n+1)}\big)\\
&= \left(\mathcal{I}(\Phi_n), \sum_{k=n}^{m} \mathbb{D}_1(\Phi_k), \sum_{k=n}^{m} \mathbb{D}_2(\Phi_k), \ldots, \sum_{k=n}^{m} \mathbb{D}_{\mathcal{L}(\Phi_n)-1}(\Phi_k), \mathcal{O}(\Phi_n)\right).
\end{aligned}
\tag{2.175}
$$

This establishes items (i) and (ii). Next observe that Lemma 2.2.33 and (2.171) ensure that for all $a \in C(\mathbb{R}, \mathbb{R})$, $x \in \mathbb{R}^{\mathcal{I}(\Phi_n)}$ it holds that

$$
\mathcal{R}_a([\mathbf{P}_{m-n+1}(\Phi_n, \Phi_{n+1}, \ldots, \Phi_m)] \bullet \mathbb{T}_{\mathcal{I}(\Phi_n),m-n+1}) \in C(\mathbb{R}^{\mathcal{I}(\Phi_n)}, \mathbb{R}^{(m-n+1)\mathcal{O}(\Phi_n)})
\tag{2.176}
$$

and

$$
\begin{aligned}
&\big(\mathcal{R}_a\big([\mathbf{P}_{m-n+1}(\Phi_n, \Phi_{n+1}, \ldots, \Phi_m)] \bullet \mathbb{T}_{\mathcal{I}(\Phi_n),m-n+1}\big)\big)(x)\\
&= \big(\mathcal{R}_a\big(\mathbf{P}_{m-n+1}(\Phi_n, \Phi_{n+1}, \ldots, \Phi_m)\big)\big)(x, x, \ldots, x).
\end{aligned}
\tag{2.177}
$$

Combining this with item (ii) in Proposition 2.2.13 proves that for all $a \in C(\mathbb{R}, \mathbb{R})$, $x \in \mathbb{R}^{\mathcal{I}(\Phi_n)}$ it holds that

$$
\begin{aligned}
&\big(\mathcal{R}_a\big([\mathbf{P}_{m-n+1}(\Phi_n, \Phi_{n+1}, \ldots, \Phi_m)] \bullet \mathbb{T}_{\mathcal{I}(\Phi_n),m-n+1}\big)\big)(x)\\
&= \big((\mathcal{R}_a(\Phi_n))(x), (\mathcal{R}_a(\Phi_{n+1}))(x), \ldots, (\mathcal{R}_a(\Phi_m))(x)\big) \in \mathbb{R}^{(m-n+1)\mathcal{O}(\Phi_n)}.
\end{aligned}
\tag{2.178}
$$

Lemma 2.2.27, (2.172), and Lemma 2.2.8 therefore demonstrate that for all $a \in C(\mathbb{R}, \mathbb{R})$, $x \in \mathbb{R}^{\mathcal{I}(\Phi_n)}$ it holds that $\mathcal{R}_a(\bigoplus_{k=n}^{m} \Phi_k) \in C(\mathbb{R}^{\mathcal{I}(\Phi_n)}, \mathbb{R}^{\mathcal{O}(\Phi_n)})$ and

$$
\begin{aligned}
&\left(\mathcal{R}_a\left(\bigoplus_{k=n}^{m} \Phi_k\right)\right)(x) \\
&= \left(\mathcal{R}_a\left(\mathbb{S}_{\mathcal{O}(\Phi_n),m-n+1} \bullet [\mathbf{P}_{m-n+1}(\Phi_n, \Phi_{n+1}, \ldots, \Phi_m)] \bullet \mathbb{T}_{\mathcal{I}(\Phi_n),m-n+1}\right)\right)(x) \\
&= \sum_{k=n}^{m} (\mathcal{R}_a(\Phi_k))(x).
\end{aligned}
\tag{2.179}
$$

This establishes items (iii)–(iv). The proof of Lemma 2.2.35 is thus complete. □

## 2.2.8 On the connection to the vectorized description of ANNs

**Definition 2.2.36.** *We denote by $\mathcal{T} \colon \mathbf{N} \to \left(\bigcup_{d\in\mathbb{N}} \mathbb{R}^d\right)$ the function which satisfies for all $L, d \in \mathbb{N}$, $l_0, l_1, \ldots, l_L \in \mathbb{N}$, $\Phi = ((W_1, B_1), (W_2, B_2), \ldots, (W_L, B_L)) \in \left(\bigtimes_{m=1}^{L}(\mathbb{R}^{l_m \times l_{m-1}} \times \mathbb{R}^{l_m})\right)$, $\theta = (\theta_1, \theta_2, \ldots, \theta_d) \in \mathbb{R}^d$, $k \in \{1, 2, \ldots, L\}$ with $\mathcal{T}(\Phi) = \theta$ that*

$$
d = \mathcal{P}(\Phi), \qquad B_k = \begin{pmatrix} \theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+l_k l_{k-1}+1} \\ \theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+l_k l_{k-1}+2} \\ \theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+l_k l_{k-1}+3} \\ \vdots \\ \theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+l_k l_{k-1}+l_k} \end{pmatrix},
\tag{2.180}
$$

*and*

$$
W_k = 
\begin{pmatrix}
\theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+1} & \theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+2} & \cdots & \theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+l_{k-1}} \\
\theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+l_{k-1}+1} & \theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+l_{k-1}+2} & \cdots & \theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+2l_{k-1}} \\
\theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+2l_{k-1}+1} & \theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+2l_{k-1}+2} & \cdots & \theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+3l_{k-1}} \\
\vdots & \vdots & \ddots & \vdots \\
\theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+(l_k-1)l_{k-1}+1} & \theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+(l_k-1)l_{k-1}+2} & \cdots & \theta_{(\sum_{i=1}^{k-1} l_i(l_{i-1}+1))+l_k l_{k-1}}
\end{pmatrix},
\tag{2.181}
$$

*(cf. Definition 2.2.1).*

**Lemma 2.2.37.** *Let $a, b \in \mathbb{N}$, $W = (W_{i,j})_{(i,j)\in\{1,2,\ldots,a\}\times\{1,2,\ldots,b\}} \in \mathbb{R}^{a\times b}$, $B = (B_i)_{i\in\{1,2,\ldots,a\}} \in \mathbb{R}^a$. Then*

$$
\begin{aligned}
&\mathcal{T}(\mathbf{A}_{W,B}) = \\
&\left(W_{1,1}, W_{1,2}, \ldots, W_{1,b}, W_{2,1}, W_{2,2}, \ldots, W_{2,b}, \ldots, W_{a,1}, W_{a,2}, \ldots, W_{a,b}, B_1, B_2, \ldots, B_a\right)
\end{aligned}
\tag{2.182}
$$

*(cf. Definitions 2.2.20 and 2.2.36).*

*Proof of Lemma 2.2.37.* Observe that (2.180) clearly establishes (2.182). The proof of Lemma 2.2.37 is thus complete. □

**Lemma 2.2.38.** *Let $L \in \mathbb{N}$, $l_0, l_1, \ldots, l_L \in \mathbb{N}$, let $W_k = (W_{k,i,j})_{(i,j)\in\{1,2,\ldots,l_k\}\times\{1,2,\ldots,l_{k-1}\}} \in \mathbb{R}^{l_k \times l_{k-1}}$, $k \in \{1, 2, \ldots, L\}$, and let $B_k = (B_{k,i})_{i\in\{1,2,\ldots,l_k\}} \in \mathbb{R}^{l_k}$, $k \in \{1, 2, \ldots, L\}$. Then*

(i)  *it holds for all $k \in \{1, 2, \ldots, L\}$ that*

$$\mathcal{T}\big(((W_k, B_k))\big) = \big(W_{k,1,1}, W_{k,1,2}, \ldots, W_{k,1,l_{k-1}}, W_{k,2,1}, W_{k,2,2}, \ldots, W_{k,2,l_{k-1}}, \ldots,$$
$$W_{k,l_k,1}, W_{k,l_k,2}, \ldots, W_{k,l_k,l_{k-1}}, B_{k,1}, B_{k,2}, \ldots, B_{k,l_k}\big) \quad (2.183)$$

*and*

(ii)  *it holds that*

$$\mathcal{T}\Big(((W_1, B_1), (W_2, B_2), \ldots, (W_L, B_L))\Big)$$
$$= \Big(W_{1,1,1}, W_{1,1,2}, \ldots, W_{1,1,l_0}, \ldots, W_{1,l_1,1}, W_{1,l_1,2}, \ldots, W_{1,l_1,l_0}, B_{1,1}, B_{1,2}, \ldots, B_{1,l_1},$$
$$W_{2,1,1}, W_{2,1,2}, \ldots, W_{2,1,l_1}, \ldots, W_{2,l_2,1}, W_{2,l_2,2}, \ldots, W_{2,l_2,l_1}, B_{2,1}, B_{2,2}, \ldots, B_{2,l_2},$$
$$\ldots,$$
$$W_{L,1,1}, W_{L,1,2}, \ldots, W_{L,1,l_{L-1}}, \ldots W_{L,l_L,1}, W_{L,l_L,2}, \ldots, W_{L,l_L,l_{L-1}},$$
$$B_{L,1}, B_{L,2}, \ldots, B_{L,l_L}\Big)$$
$$(2.184)$$

*(cf. Definition 2.2.36).*

*Proof of Lemma 2.2.38.* Note that Lemma 2.2.37 proves item (i). Moreover, observe that (2.180) establishes item (ii). The proof of Lemma 2.2.38 is thus complete. $\qquad\square$

**Exercise 2.2.2.** *Prove or disprove the following statement: The function $\mathcal{T}$ is injective (cf. Definition 2.2.36).*

**Exercise 2.2.3.** *Prove or disprove the following statement: The function $\mathcal{T}$ is surjective (cf. Definition 2.2.36).*

**Exercise 2.2.4.** *Prove or disprove the following statement: The function $\mathcal{T}$ is bijective (cf. Definition 2.2.36).*

**Lemma 2.2.39.** *Let $a \in C(\mathbb{R}, \mathbb{R})$, $\Phi \in \mathbf{N}$, $L \in \mathbb{N}$, $l_0, l_1, \ldots, l_L \in \mathbb{N}$ satisfy $\mathcal{D}(\Phi) = (l_0, l_1, \ldots, l_L)$ (cf. Definition 2.2.1). Then it holds for all $x \in \mathbb{R}^{l_0}$ that*

$$(\mathcal{R}_a(\Phi))(x) = \begin{cases} \big(\mathcal{N}_{\mathrm{id}_{\mathbb{R}^{l_L}}}^{\mathcal{T}(\Phi), l_0}\big)(x) & : L = 1 \\ \big(\mathcal{N}_{\mathfrak{M}_{a,l_1}, \mathfrak{M}_{a,l_2}, \ldots, \mathfrak{M}_{a,l_{L-1}}, \mathrm{id}_{\mathbb{R}^{l_L}}}^{\mathcal{T}(\Phi), l_0}\big)(x) & : L > 1 \end{cases} \quad (2.185)$$

*(cf. Definitions 2.1.2, 2.1.4, 2.2.3, and 2.2.36).*

*Proof of Lemma 2.2.39.* Throughout this proof let $((W_1, B_1), (W_2, B_2), \ldots, (W_L, B_L)) \in \big(\bigtimes_{k=1}^L (\mathbb{R}^{l_k \times l_{k-1}} \times \mathbb{R}^{l_k})\big)$ satisfy $\Phi = ((W_1, B_1), (W_2, B_2), \ldots, (W_L, B_L))$. Note that (2.180) shows that for all $k \in \{1, 2, \ldots, L\}$, $x \in \mathbb{R}^{l_{k-1}}$ it holds that

$$W_k x + B_k = \big(\mathcal{A}_{l_k, l_{k-1}}^{\mathcal{T}(\Phi), \sum_{i=1}^{k-1} l_i(l_{i-1}+1)}\big)(x) \quad (2.186)$$

(cf. Definitions 2.1.1 and 2.2.36). This demonstrates that for all $x_0 \in \mathbb{R}^{l_0}$, $x_1 \in \mathbb{R}^{l_1}, \ldots, x_L \in \mathbb{R}^{l_L}$ with $\forall\, k \in \{1, 2, \ldots, L\}$: $x_k = \mathfrak{M}_{a,l_k}(W_k x_{k-1} + B_k)$ it holds that

$$x_{L-1} = \tag{2.187}$$

$$
\begin{cases}
x_0 & : L = 1 \\
\Big(\mathfrak{M}_{a,l_{L-1}} \circ \mathcal{A}_{l_{L-1},l_{L-2}}^{\mathcal{T}(\Phi),\sum_{i=1}^{L-2} l_i(l_{i-1}+1)} \circ \mathfrak{M}_{a,l_{L-2}} \circ \mathcal{A}_{l_{L-2},l_{L-3}}^{\mathcal{T}(\Phi),\sum_{i=1}^{L-3} l_i(l_{i-1}+1)} \circ \\
\qquad \ldots \circ \mathfrak{M}_{a,l_1} \circ \mathcal{A}_{l_1,l_0}^{\mathcal{T}(\Phi),0}\Big)(x_0) & : L > 1
\end{cases}
$$

(cf. Definition 2.1.4). Combining this and (2.186) with (2.3) and (2.53) proves that for all $x_0 \in \mathbb{R}^{l_0}$, $x_1 \in \mathbb{R}^{l_1}, \ldots, x_L \in \mathbb{R}^{l_L}$ with $\forall\, k \in \{1, 2, \ldots, L\}$: $x_k = \mathfrak{M}_{a,l_k}(W_k x_{k-1} + B_k)$ it holds that

$$
\begin{aligned}
\big(\mathcal{R}_a(\Phi)\big)(x_0) = W_L x_{L-1} + B_L &= \Big(\mathcal{A}_{l_L,l_{L-1}}^{\mathcal{T}(\Phi),\sum_{i=1}^{L-1} l_i(l_{i-1}+1)}\Big)(x_{L-1}) \\
&= \begin{cases}
\big(\mathcal{N}_{\mathrm{id}_{\mathbb{R}^{l_L}}}^{\mathcal{T}(\Phi),l_0}\big)(x_0) & : L = 1 \\
\big(\mathcal{N}_{\mathfrak{M}_{a,l_1},\mathfrak{M}_{a,l_2},\ldots,\mathfrak{M}_{a,l_{L-1}},\mathrm{id}_{\mathbb{R}^{l_L}}}^{\mathcal{T}(\Phi),l_0}\big)(x_0) & : L > 1
\end{cases}
\end{aligned}
\tag{2.188}
$$

(cf. Definitions 2.1.2 and 2.2.3). The proof of Lemma 2.2.39 is thus complete. $\qquad\square$

**Corollary 2.2.40.** *Let $\Phi \in \mathbf{N}$ (cf. Definition 2.2.1). Then it holds for all $x \in \mathbb{R}^{\mathcal{I}(\Phi)}$ that*

$$\big(\mathcal{N}_{-\infty,\infty}^{\mathcal{T}(\Phi),\mathcal{D}(\Phi)}\big)(x) = (\mathcal{R}_{\mathfrak{r}}(\Phi))(x) \tag{2.189}$$

*(cf. Definitions 2.1.6, 2.1.27, 2.2.3, and 2.2.36).*

*Proof of Corollary 2.2.40.* Note that Lemma 2.2.39, (2.50), (2.11), and the fact that for all $d \in \mathbb{N}$ it holds that $\mathfrak{C}_{-\infty,\infty,d} = \mathrm{id}_{\mathbb{R}^d}$ establish (2.189) (cf. Definition 2.1.12). The proof of Corollary 2.2.40 is thus complete. $\qquad\square$

# Chapter 3

# Low-dimensional neural network approximation results

## 3.1 One-dimensional neural network approximation results

### 3.1.1 Linear interpolation of one-dimensional functions

#### 3.1.1.1 On the modulus of continuity

**Definition 3.1.1.** *Let $A \subseteq \mathbb{R}$ be a set and let $f\colon A \to \mathbb{R}$ be a function. Then we denote by $w_f\colon [0,\infty] \to [0,\infty]$ the function which satisfies for all $h \in [0,\infty]$ that*

$$w_f(h) = \sup\bigl(\bigl\{|f(x) - f(y)| \in [0,\infty)\colon (x, y \in A \text{ with } |x - y| \le h)\bigr\} \cup \{0\}\bigr) \qquad (3.1)$$

*and we call $w_f$ the modulus of continuity of $f$.*

**Lemma 3.1.2.** *Let $a \in [-\infty, \infty]$, $b \in [a, \infty]$ and let $f\colon ([a,b] \cap \mathbb{R}) \to \mathbb{R}$ be a function. Then*

  *(i) it holds that $w_f$ is non-decreasing,*

  *(ii) it holds that $f$ is uniformly continuous if and only if $\lim_{h \searrow 0} w_f(h) = 0$,*

  *(iii) it holds that $f$ is globally bounded if and only if $w_f(\infty) < \infty$,*

  *(iv) it holds for all $x, y \in [a,b] \cap \mathbb{R}$ that $|f(x) - f(y)| \le w_f(|x - y|)$, and*

  *(v) it holds for all $h, \mathfrak{h} \in [0,\infty]$ that $w_f(h + \mathfrak{h}) \le w_f(h) + w_f(\mathfrak{h})$*

*(cf. Definition 3.1.1).*

*Proof of Lemma 3.1.2.* First, observe that (3.1) implies items (i), (ii), (iii), and (iv). Moreover, note that (3.1) and the triangle inequality assure that for all $h, \mathfrak{h} \in [0,\infty]$ it

holds that

$$
\begin{aligned}
&w_f(h + \mathfrak{h})\\
&= \sup\big(\big\{|f(x) - f(y)| \in [0, \infty)\colon\\
&\qquad (x, y \in [a, b] \cap \mathbb{R} \text{ with } x \leq y \text{ and } |x - y| \leq (h + \mathfrak{h}))\big\} \cup \{0\}\big)\\
&= \sup\big(\big\{|f(x) - f(y)| \in [0, \infty)\colon\\
&\qquad \big(x, y, z \in [a, b] \cap \mathbb{R} \text{ with } x \leq z \leq y, |x - z| \leq h, \text{ and } |y - z| \leq \mathfrak{h}\big)\big\} \cup \{0\}\big)\\
&\leq \sup\big(\big\{|f(x) - f(z)| + |f(z) - f(y)| \in [0, \infty)\colon\\
&\qquad (x, y, z \in [a, b] \cap \mathbb{R} \text{ with } x \leq z \leq y, |x - z| \leq h, \text{ and } |z - y| \leq \mathfrak{h})\big\} \cup \{0\}\big)\\
&\leq \sup\big(\big\{|f(x) - f(z)| \in [0, \infty)\colon (x, z \in [a, b] \cap \mathbb{R} \text{ with } x \leq z \text{ and } |x - z| \leq h)\big\} \cup \{0\}\big)\\
&\quad + \sup\big(\big\{|f(z) - f(y)| \in [0, \infty)\colon (z, y \in [a, b] \cap \mathbb{R} \text{ with } z \leq y \text{ and } |z - y| \leq \mathfrak{h})\big\} \cup \{0\}\big)\\
&= w_f(h) + w_f(\mathfrak{h})
\end{aligned}
\tag{3.2}
$$

(cf. Definition 3.1.1). This establishes item (v). The proof of Lemma 3.1.2 is thus complete. $\qquad\square$

**Lemma 3.1.3.** *Let $A \subseteq \mathbb{R}$, $L \in [0, \infty)$ and let $f\colon A \to \mathbb{R}$ satisfy for all $x, y \in A$ that $|f(x) - f(y)| \leq L|x - y|$. Then it holds for all $h \in [0, \infty)$ that $w_f(h) \leq Lh$.*

*Proof of Lemma 3.1.3.* Observe that the assumption that for all $x, y \in A$ it holds that $|f(x) - f(y)| \leq L|x - y|$ and (3.1) imply that for all $h \in [0, \infty)$ it holds that

$$
\begin{aligned}
w_f(h) &= \sup\big(\big\{|f(x) - f(y)| \in [0, \infty)\colon (x, y \in A \text{ with } |x - y| \leq h)\big\} \cup \{0\}\big)\\
&\leq \sup\big(\big\{L|x - y| \in [0, \infty)\colon (x, y \in A \text{ with } |x - y| \leq h)\big\} \cup \{0\}\big)\\
&\leq \sup(\{Lh, 0\}) = Lh.
\end{aligned}
\tag{3.3}
$$

The proof of Lemma 3.1.3 is thus complete. $\qquad\square$

#### 3.1.1.2 Linear interpolation of one-dimensional functions

**Definition 3.1.4** (Linear interpolation operator). *Let $K \in \mathbb{N}$, $\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K, f_0, f_1, \ldots, f_K \in \mathbb{R}$ satisfy $\mathfrak{x}_0 < \mathfrak{x}_1 < \ldots < \mathfrak{x}_K$. Then we denote by $\mathscr{L}^{f_0, f_1, \ldots, f_K}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K}\colon \mathbb{R} \to \mathbb{R}$ the function which satisfies for all $k \in \{1, 2, \ldots, K\}$, $x \in (-\infty, \mathfrak{x}_0)$, $y \in [\mathfrak{x}_{k-1}, \mathfrak{x}_k)$, $z \in [\mathfrak{x}_K, \infty)$ that $(\mathscr{L}^{f_0, f_1, \ldots, f_K}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K})(x) = f_0$, $(\mathscr{L}^{f_0, f_1, \ldots, f_K}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K})(z) = f_K$, and*

$$
(\mathscr{L}^{f_0, f_1, \ldots, f_K}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K})(y) = f_{k-1} + \big(\tfrac{y - \mathfrak{x}_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}}\big)(f_k - f_{k-1}).
\tag{3.4}
$$

**Lemma 3.1.5.** *Let $K \in \mathbb{N}$, $\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K, f_0, f_1, \ldots, f_K \in \mathbb{R}$ satisfy $\mathfrak{x}_0 < \mathfrak{x}_1 < \ldots < \mathfrak{x}_K$. Then*

*(i) it holds for all $k \in \{0, 1, \ldots, K\}$ that*

$$
(\mathscr{L}^{f_0, f_1, \ldots, f_K}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K})(\mathfrak{x}_k) = f_k,
\tag{3.5}
$$

*(ii) it holds for all $k \in \{1, 2, \ldots, K\}$, $x \in [\mathfrak{x}_{k-1}, \mathfrak{x}_k]$ that*

$$
(\mathscr{L}^{f_0, f_1, \ldots, f_K}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K})(x) = f_{k-1} + \big(\tfrac{x - \mathfrak{x}_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}}\big)(f_k - f_{k-1}),
\tag{3.6}
$$

*and*

*(iii) it holds for all $k \in \{1, 2, \ldots, K\}$, $x \in [\mathfrak{x}_{k-1}, \mathfrak{x}_k]$ that*

$$(\mathscr{L}^{f_0, f_1, \ldots, f_K}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K})(x) = \left(\tfrac{\mathfrak{x}_k - x}{\mathfrak{x}_k - \mathfrak{x}_{k-1}}\right) f_{k-1} + \left(\tfrac{x - \mathfrak{x}_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}}\right) f_k. \tag{3.7}$$

*(cf. Definition 3.1.4).*

*Proof of Lemma 3.1.5.* Observe that (3.4) implies items (i) and (ii). Moreover, note that item (ii) implies that for all $k \in \{1, 2, \ldots, K\}$, $x \in [\mathfrak{x}_{k-1}, \mathfrak{x}_k]$ it holds that

$$\begin{aligned}
(\mathscr{L}^{f_0, f_1, \ldots, f_K}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K})(x) &= \left[\left(\tfrac{\mathfrak{x}_k - \mathfrak{x}_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}}\right) - \left(\tfrac{x - \mathfrak{x}_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}}\right)\right] f_{k-1} + \left(\tfrac{x - \mathfrak{x}_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}}\right) f_k \\
&= \left(\tfrac{\mathfrak{x}_k - x}{\mathfrak{x}_k - \mathfrak{x}_{k-1}}\right) f_{k-1} + \left(\tfrac{x - \mathfrak{x}_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}}\right) f_k.
\end{aligned} \tag{3.8}$$

This proves item (iii). The proof of Lemma 3.1.5 is thus complete. $\qquad\square$

**Lemma 3.1.6.** *Let $K \in \mathbb{N}$, $\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K \in \mathbb{R}$ satisfy $\mathfrak{x}_0 < \mathfrak{x}_1 < \ldots < \mathfrak{x}_K$ and let $f\colon [\mathfrak{x}_0, \mathfrak{x}_K] \to \mathbb{R}$ be a function. Then*

*(i) it holds for all $x, y \in \mathbb{R}$ with $x \neq y$ that*

$$\begin{aligned}
&\left|(\mathscr{L}^{f(\mathfrak{x}_0), f(\mathfrak{x}_1), \ldots, f(\mathfrak{x}_K)}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K})(x) - (\mathscr{L}^{f(\mathfrak{x}_0), f(\mathfrak{x}_1), \ldots, f(\mathfrak{x}_K)}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K})(y)\right| \\
&\leq \left(\max_{k \in \{1, 2, \ldots, K\}} \left(\frac{w_f(|\mathfrak{x}_k - \mathfrak{x}_{k-1}|)}{|\mathfrak{x}_k - \mathfrak{x}_{k-1}|}\right)\right) |x - y|
\end{aligned} \tag{3.9}$$

*and*

*(ii) it holds that $\sup_{x \in [\mathfrak{x}_0, \mathfrak{x}_K]} \left|(\mathscr{L}^{f(\mathfrak{x}_0), f(\mathfrak{x}_1), \ldots, f(\mathfrak{x}_K)}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K})(x) - f(x)\right| \leq w_f(\max_{k \in \{1, 2, \ldots, K\}} |\mathfrak{x}_k - \mathfrak{x}_{k-1}|)$*

*(cf. Definitions 3.1.1 and 3.1.4).*

*Proof of Lemma 3.1.6.* Throughout this proof let $\mathfrak{l}\colon \mathbb{R} \to \mathbb{R}$ satisfy for all $x \in \mathbb{R}$ that $\mathfrak{l}(x) = (\mathscr{L}^{f(\mathfrak{x}_0), f(\mathfrak{x}_1), \ldots, f(\mathfrak{x}_K)}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K})(x)$ and let $L \in [0, \infty]$ satisfy

$$L = \max_{k \in \{1, 2, \ldots, K\}} \left(\frac{w_f(|\mathfrak{x}_k - \mathfrak{x}_{k-1}|)}{|\mathfrak{x}_k - \mathfrak{x}_{k-1}|}\right) \tag{3.10}$$

(cf. Definitions 3.1.1 and 3.1.4). Observe that item (ii) in Lemma 3.1.5, item (iv) in Lemma 3.1.2, and (3.10) assure that for all $k \in \{1, 2, \ldots, K\}$, $x, y \in [\mathfrak{x}_{k-1}, \mathfrak{x}_k]$ with $x \neq y$ it holds that

$$\begin{aligned}
|\mathfrak{l}(x) - \mathfrak{l}(y)| &= \left|\left(\tfrac{x - \mathfrak{x}_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}}\right)(f(\mathfrak{x}_k) - f(\mathfrak{x}_{k-1})) - \left(\tfrac{y - \mathfrak{x}_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}}\right)(f(\mathfrak{x}_k) - f(\mathfrak{x}_{k-1}))\right| \\
&= \left|\left(\frac{f(\mathfrak{x}_k) - f(\mathfrak{x}_{k-1})}{\mathfrak{x}_k - \mathfrak{x}_{k-1}}\right)(x - y)\right| \leq \left(\frac{w_f(|\mathfrak{x}_k - \mathfrak{x}_{k-1}|)}{|\mathfrak{x}_k - \mathfrak{x}_{k-1}|}\right) |x - y| \leq L|x - y|.
\end{aligned} \tag{3.11}$$

This, item (iv) in Lemma 3.1.2, Lemma 3.1.5, and (3.10) ensure that for all $k, l \in$

$\{1, 2, \ldots, K\}$, $x \in [\mathfrak{x}_{k-1}, \mathfrak{x}_k]$, $y \in [\mathfrak{x}_{l-1}, \mathfrak{x}_l]$ with $k < l$ and $x \neq y$ it holds that

$$
\begin{aligned}
|\mathfrak{l}(x) &- \mathfrak{l}(y)| \\
&\leq |\mathfrak{l}(x) - \mathfrak{l}(\mathfrak{x}_k)| + |\mathfrak{l}(\mathfrak{x}_k) - \mathfrak{l}(\mathfrak{x}_{l-1})| + |\mathfrak{l}(\mathfrak{x}_{l-1}) - \mathfrak{l}(y)| \\
&= |\mathfrak{l}(x) - \mathfrak{l}(\mathfrak{x}_k)| + |f(\mathfrak{x}_k) - f(\mathfrak{x}_{l-1})| + |\mathfrak{l}(\mathfrak{x}_{l-1}) - \mathfrak{l}(y)| \\
&\leq |\mathfrak{l}(x) - \mathfrak{l}(\mathfrak{x}_k)| + \left( \sum_{j=k+1}^{l-1} |f(\mathfrak{x}_j) - f(\mathfrak{x}_{j-1})| \right) + |\mathfrak{l}(\mathfrak{x}_{l-1}) - \mathfrak{l}(y)| \\
&\leq |\mathfrak{l}(x) - \mathfrak{l}(\mathfrak{x}_k)| + \left( \sum_{j=k+1}^{l-1} w_f(|\mathfrak{x}_j - \mathfrak{x}_{j-1}|) \right) + |\mathfrak{l}(\mathfrak{x}_{l-1}) - \mathfrak{l}(y)| \\
&\leq L \left( (\mathfrak{x}_k - x) + \left( \sum_{j=k+1}^{l-1} (\mathfrak{x}_j - \mathfrak{x}_{j-1}) \right) + (y - \mathfrak{x}_{l-1}) \right) = L|x - y|.
\end{aligned}
\tag{3.12}
$$

Combining this and (3.11) shows that for all $x, y \in [\mathfrak{x}_0, \mathfrak{x}_K]$ with $x \neq y$ it holds that $|\mathfrak{l}(x) - \mathfrak{l}(y)| \leq L|x - y|$. This, the fact that for all $x, y \in (-\infty, \mathfrak{x}_0]$ with $x \neq y$ it holds that $|\mathfrak{l}(x) - \mathfrak{l}(y)| = 0 \leq L|x - y|$, the fact that for all $x, y \in [\mathfrak{x}_K, \infty)$ with $x \neq y$ it holds that $|\mathfrak{l}(x) - \mathfrak{l}(y)| = 0 \leq L|x - y|$, and the triangle inequality hence demonstrate that for all $x, y \in \mathbb{R}$ with $x \neq y$ it holds that $|\mathfrak{l}(x) - \mathfrak{l}(y)| \leq L|x - y|$. This proves item (i). Moreover, note that (3.1), Lemma 3.1.2, and item (iii) in Lemma 3.1.5 assure that for all $k \in \{1, 2, \ldots, K\}$, $x \in [\mathfrak{x}_{k-1}, \mathfrak{x}_k]$ it holds that

$$
\begin{aligned}
|\mathfrak{l}(x) - f(x)| &= \left| \left( \frac{\mathfrak{x}_k - x}{\mathfrak{x}_k - \mathfrak{x}_{k-1}} \right) f(\mathfrak{x}_k) + \left( \frac{x - \mathfrak{x}_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}} \right) f(\mathfrak{x}_{k-1}) - f(x) \right| \\
&= \left| \left( \frac{\mathfrak{x}_k - x}{\mathfrak{x}_k - \mathfrak{x}_{k-1}} \right) (f(\mathfrak{x}_k) - f(x)) + \left( \frac{x - \mathfrak{x}_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}} \right) (f(\mathfrak{x}_{k-1}) - f(x)) \right| \\
&\leq \left( \frac{\mathfrak{x}_k - x}{\mathfrak{x}_k - \mathfrak{x}_{k-1}} \right) |f(\mathfrak{x}_k) - f(x)| + \left( \frac{x - \mathfrak{x}_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}} \right) |f(\mathfrak{x}_{k-1}) - f(x)| \\
&\leq w_f(|\mathfrak{x}_k - \mathfrak{x}_{k-1}|) \left( \frac{\mathfrak{x}_k - x}{\mathfrak{x}_k - \mathfrak{x}_{k-1}} + \frac{x - \mathfrak{x}_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}} \right) \\
&= w_f(|\mathfrak{x}_k - \mathfrak{x}_{k-1}|) \leq w_f(\max_{j \in \{1,2,\ldots,K\}} |\mathfrak{x}_j - \mathfrak{x}_{j-1}|).
\end{aligned}
\tag{3.13}
$$

This establishes item (ii). The proof of Lemma 3.1.6 is thus complete. $\qquad \square$

**Lemma 3.1.7.** *Let $K \in \mathbb{N}$, $L, \mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K \in \mathbb{R}$ satisfy $\mathfrak{x}_0 < \mathfrak{x}_1 < \ldots < \mathfrak{x}_K$ and let $f \colon [\mathfrak{x}_0, \mathfrak{x}_K] \to \mathbb{R}$ satisfy for all $x, y \in [\mathfrak{x}_0, \mathfrak{x}_K]$ that $|f(x) - f(y)| \leq L|x - y|$. Then*

*(i) it holds for all $x, y \in \mathbb{R}$ that*

$$
\left| (\mathscr{L}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K}^{f(\mathfrak{x}_0), f(\mathfrak{x}_1), \ldots, f(\mathfrak{x}_K)})(x) - (\mathscr{L}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K}^{f(\mathfrak{x}_0), f(\mathfrak{x}_1), \ldots, f(\mathfrak{x}_K)})(y) \right| \leq L|x - y|
\tag{3.14}
$$

*and*

*(ii) it holds that $\sup_{x \in [\mathfrak{x}_0, \mathfrak{x}_K]} \left| (\mathscr{L}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K}^{f(\mathfrak{x}_0), f(\mathfrak{x}_1), \ldots, f(\mathfrak{x}_K)})(x) - f(x) \right| \leq L(\max_{k \in \{1,2,\ldots,K\}} |\mathfrak{x}_k - \mathfrak{x}_{k-1}|)$*

*(cf. Definition 3.1.4).*

*Proof of Lemma 3.1.7.* Note that the assumption that for all $x, y \in [\mathfrak{r}_0, \mathfrak{r}_K]$ it holds that $|f(x) - f(y)| \le L|x - y|$, Lemma 3.1.3, and item (i) in Lemma 3.1.6 demonstrate that for all $x, y \in \mathbb{R}$ it holds that

$$
\begin{aligned}
&\left|(\mathscr{L}_{\mathfrak{r}_0, \mathfrak{r}_1, \ldots, \mathfrak{r}_K}^{f(\mathfrak{r}_0), f(\mathfrak{r}_1), \ldots, f(\mathfrak{r}_K)})(x) - (\mathscr{L}_{\mathfrak{r}_0, \mathfrak{r}_1, \ldots, \mathfrak{r}_K}^{f(\mathfrak{r}_0), f(\mathfrak{r}_1), \ldots, f(\mathfrak{r}_K)})(y)\right| \\
&\le \left(\max_{k \in \{1, 2, \ldots, K\}} \left(\frac{L|\mathfrak{r}_k - \mathfrak{r}_{k-1}|}{|\mathfrak{r}_k - \mathfrak{r}_{k-1}|}\right)\right)|x - y| = L|x - y|.
\end{aligned}
\tag{3.15}
$$

This proves item (i). Moreover, observe that the assumption that for all $x, y \in [\mathfrak{r}_0, \mathfrak{r}_K]$ it holds that $|f(x) - f(y)| \le L|x - y|$, Lemma 3.1.3, and item (ii) in Lemma 3.1.6 assure that

$$
\sup_{x \in [\mathfrak{r}_0, \mathfrak{r}_K]} \left|(\mathscr{L}_{\mathfrak{r}_0, \mathfrak{r}_1, \ldots, \mathfrak{r}_K}^{f(\mathfrak{r}_0), f(\mathfrak{r}_1), \ldots, f(\mathfrak{r}_K)})(x) - f(x)\right| \le L\left(\max_{k \in \{1, 2, \ldots, K\}} |\mathfrak{r}_k - \mathfrak{r}_{k-1}|\right).
\tag{3.16}
$$

This establishes item (ii). The proof of Lemma 3.1.7 is thus complete. $\qquad\square$

## 3.1.2 Activation functions as neural networks

**Definition 3.1.8** (Activation functions as neural networks). *Let $n \in \mathbb{N}$. Then we denote by $\mathfrak{i}_n \in ((\mathbb{R}^{n \times n} \times \mathbb{R}^n) \times (\mathbb{R}^{n \times n} \times \mathbb{R}^n)) \subseteq \mathbf{N}$ the neural network given by $\mathfrak{i}_n = ((I_n, 0), (I_n, 0))$ (cf. Definitions 2.2.1 and 2.2.9).*

**Lemma 3.1.9.** *Let $n \in \mathbb{N}$. Then*

*(i) it holds that $\mathcal{D}(\mathfrak{i}_n) = (n, n, n) \in \mathbb{N}^3$,*

*(ii) it holds for all $a \in C(\mathbb{R}, \mathbb{R})$ that $\mathcal{R}_a(\mathfrak{i}_n) \in C(\mathbb{R}^n, \mathbb{R}^n)$, and*

*(iii) it holds for all $a \in C(\mathbb{R}, \mathbb{R})$ that $\mathcal{R}_a(\mathfrak{i}_n) = \mathfrak{M}_{a,n}$*

*(cf. Definitions 2.1.4, 2.2.1, 2.2.3, and 3.1.8).*

*Proof of Lemma 3.1.9.* Note the fact that $\mathfrak{i}_n \in ((\mathbb{R}^{n \times n} \times \mathbb{R}^n) \times (\mathbb{R}^{n \times n} \times \mathbb{R}^n)) \subseteq \mathbf{N}$ ensures that $\mathcal{D}(\mathfrak{i}_n) = (n, n, n) \in \mathbb{N}^3$. This establishes item (i). Next observe the fact that $\mathfrak{i}_n = ((I_n, 0), (I_n, 0)) \in ((\mathbb{R}^{n \times n} \times \mathbb{R}^n) \times (\mathbb{R}^{n \times n} \times \mathbb{R}^n))$ and (2.53) prove that for all $a \in C(\mathbb{R}, \mathbb{R})$, $x \in \mathbb{R}^n$ it holds that $\mathcal{R}_a(\mathfrak{i}_n) \in C(\mathbb{R}^n, \mathbb{R}^n)$ and

$$
(\mathcal{R}_a(\mathfrak{i}_n))(x) = I_n(\mathfrak{M}_{a,n}(I_n x + 0)) + 0 = \mathfrak{M}_{a,n}(x).
\tag{3.17}
$$

This establishes items (ii) and (iii). The proof of Lemma 3.1.9 is thus complete. $\qquad\square$

**Lemma 3.1.10.** *Let $\Phi \in \mathbf{N}$ (cf. Definition 2.2.1). Then*

*(i) it holds that*

$$
\begin{aligned}
&\mathcal{D}(\mathfrak{i}_{\mathcal{O}(\Phi)} \bullet \Phi) \\
&= (\mathbb{D}_0(\Phi), \mathbb{D}_1(\Phi), \mathbb{D}_2(\Phi), \ldots, \mathbb{D}_{\mathcal{L}(\Phi)-1}(\Phi), \mathbb{D}_{\mathcal{L}(\Phi)}(\Phi), \mathbb{D}_{\mathcal{L}(\Phi)}(\Phi)) \in \mathbb{N}^{\mathcal{L}(\Phi)+2},
\end{aligned}
\tag{3.18}
$$

*(ii) it holds for all $a \in C(\mathbb{R}, \mathbb{R})$ that $\mathcal{R}_a(\mathfrak{i}_{\mathcal{O}(\Phi)} \bullet \Phi) \in C(\mathbb{R}^{\mathcal{I}(\Phi)}, \mathbb{R}^{\mathcal{O}(\Phi)})$,*

*(iii) it holds for all $a \in C(\mathbb{R}, \mathbb{R})$ that $\mathcal{R}_a(\mathfrak{i}_{\mathcal{O}(\Phi)} \bullet \Phi) = \mathfrak{M}_{a,\mathcal{O}(\Phi)} \circ (\mathcal{R}_a(\Phi))$,*

(iv) it holds that

$$
\begin{aligned}
&\mathcal{D}(\Phi \bullet \mathfrak{i}_{\mathcal{I}(\Phi)}) \\
&= (\mathbb{D}_0(\Phi), \mathbb{D}_0(\Phi), \mathbb{D}_1(\Phi), \mathbb{D}_2(\Phi), \ldots, \mathbb{D}_{\mathcal{L}(\Phi)-1}(\Phi), \mathbb{D}_{\mathcal{L}(\Phi)}(\Phi)) \in \mathbb{N}^{\mathcal{L}(\Phi)+2},
\end{aligned}
\tag{3.19}
$$

(v) it holds for all $a \in C(\mathbb{R}, \mathbb{R})$ that $\mathcal{R}_a(\Phi \bullet \mathfrak{i}_{\mathcal{I}(\Phi)}) \in C(\mathbb{R}^{\mathcal{I}(\Phi)}, \mathbb{R}^{\mathcal{O}(\Phi)})$, and

(vi) it holds for all $a \in C(\mathbb{R}, \mathbb{R})$ that $\mathcal{R}_a(\Phi \bullet \mathfrak{i}_{\mathcal{I}(\Phi)}) = (\mathcal{R}_a(\Phi)) \circ \mathfrak{M}_{a,\mathcal{I}(\Phi)}$

*(cf. Definitions 2.1.4, 2.2.3, 2.2.5, and 3.1.8).*

*Proof of Lemma 3.1.10.* Note that Lemma 3.1.9 demonstrates that for all $n \in \mathbb{N}$, $a \in C(\mathbb{R}, \mathbb{R})$, $x \in \mathbb{R}^n$ it holds that $\mathcal{R}_a(\mathfrak{i}_n) \in C(\mathbb{R}^n, \mathbb{R}^n)$ and

$$
(\mathcal{R}_a(\mathfrak{i}_n))(x) = \mathfrak{M}_{a,n}(x)
\tag{3.20}
$$

(cf. Definitions 2.1.4, 2.2.3, and 3.1.8). Combining this and Proposition 2.2.7 establishes items (i), (ii), (iii), (iv), (v), and (vi). The proof of Lemma 3.1.10 is thus complete. $\square$

### 3.1.3 Linear interpolation with neural networks

**Lemma 3.1.11.** *Let $\alpha, \beta, h \in \mathbb{R}$, $\mathbf{H} \in \mathbf{N}$ satisfy $\mathbf{H} = h \circledast (\mathfrak{i}_1 \bullet \mathbf{A}_{\alpha,\beta})$ (cf. Definitions 2.2.1, 2.2.5, 2.2.20, 2.2.23, and 3.1.8). Then*

(i) *it holds that $\mathbf{H} = ((\alpha, \beta), (h, 0))$,*

(ii) *it holds that $\mathcal{D}(\mathbf{H}) = (1, 1, 1) \in \mathbb{N}^3$,*

(iii) *it holds that $\mathcal{R}_{\mathfrak{r}}(\mathbf{H}) \in C(\mathbb{R}, \mathbb{R})$, and*

(iv) *it holds for all $x \in \mathbb{R}$ that $(\mathcal{R}_{\mathfrak{r}}(\mathbf{H}))(x) = h \max\{\alpha x + \beta, 0\}$*

*(cf. Definitions 2.1.6 and 2.2.3).*

*Proof of Lemma 3.1.11.* Note that Lemma 2.2.21 ensures that $\mathbf{A}_{\alpha,\beta} = (\alpha, \beta)$, $\mathcal{D}(\mathbf{A}_{\alpha,\beta}) = (1, 1) \in \mathbb{N}^2$, $\mathcal{R}_{\mathfrak{r}}(\mathbf{A}_{\alpha,\beta}) \in C(\mathbb{R}, \mathbb{R})$, and $\forall x \in \mathbb{R} \colon (\mathcal{R}_{\mathfrak{r}}(\mathbf{A}_{\alpha,\beta}))(x) = \alpha x + \beta$ (cf. Definitions 2.1.6 and 2.2.3). Lemmas 3.1.9 and 3.1.10, (2.10), (2.53), and (2.59) therefore imply that $\mathfrak{i}_1 \bullet \mathbf{A}_{\alpha,\beta} = ((\alpha, \beta), (1, 0))$, $\mathcal{D}(\mathfrak{i}_1 \bullet \mathbf{A}_{\alpha,\beta}) = (1, 1, 1) \in \mathbb{N}^3$, $\mathcal{R}_{\mathfrak{r}}(\mathfrak{i}_1 \bullet \mathbf{A}_{\alpha,\beta}) \in C(\mathbb{R}, \mathbb{R})$, and

$$
\forall x \in \mathbb{R} \colon (\mathcal{R}_{\mathfrak{r}}(\mathfrak{i}_1 \bullet \mathbf{A}_{\alpha,\beta}))(x) = \mathfrak{r}(\mathcal{R}_{\mathfrak{r}}(\mathbf{A}_{\alpha,\beta})(x)) = \max\{\alpha x + \beta, 0\}.
\tag{3.21}
$$

This, Lemma 2.2.24, and (2.149) ensure that $h \circledast (\mathfrak{i}_1 \bullet \mathbf{A}_{\alpha,\beta}) = ((\alpha, \beta), (h, 0))$, $\mathcal{R}_{\mathfrak{r}}(\mathbf{H}) \in C(\mathbb{R}, \mathbb{R})$, $\mathcal{D}(\mathbf{H}) = (1, 1, 1)$, and

$$
(\mathcal{R}_{\mathfrak{r}}(\mathbf{H}))(x) = h((\mathcal{R}_{\mathfrak{r}}(\mathfrak{i}_1 \bullet \mathbf{A}_{\alpha,\beta}))(x)) = h \max\{\alpha x + \beta, 0\}.
\tag{3.22}
$$

This establishes items (i)–(iv). The proof of Lemma 3.1.11 is thus complete. $\square$

**Lemma 3.1.12.** *Let $K \in \mathbb{N}$, $f_0, f_1, \ldots, f_K, \mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K \in \mathbb{R}$ satisfy $\mathfrak{x}_0 < \mathfrak{x}_1 < \ldots < \mathfrak{x}_K$ and let $\mathbf{F} \in \mathbf{N}$ satisfy*

$$
\mathbf{F} = \mathbf{A}_{1,f_0} \bullet \left( \bigoplus_{k=0}^{K} \left( \left( \frac{(f_{\min\{k+1,K\}} - f_k)}{(\mathfrak{x}_{\min\{k+1,K\}} - \mathfrak{x}_{\min\{k,K-1\}})} - \frac{(f_k - f_{\max\{k-1,0\}})}{(\mathfrak{x}_{\max\{k,1\}} - \mathfrak{x}_{\max\{k-1,0\}})} \right) \circledast (\mathfrak{i}_1 \bullet \mathbf{A}_{1,-\mathfrak{x}_k}) \right) \right)
\tag{3.23}
$$

*(cf. Definitions 2.2.1, 2.2.5, 2.2.20, 2.2.23, 2.2.34, and 3.1.8). Then*

(i) it holds that $\mathcal{D}(\mathbf{F}) = (1, K+1, 1) \in \mathbb{N}^3$,

(ii) it holds that $\mathcal{R}_{\mathfrak{r}}(\mathbf{F}) \in C(\mathbb{R}, \mathbb{R})$,

(iii) it holds that $\mathcal{R}_{\mathfrak{r}}(\mathbf{F}) = \mathcal{L}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K}^{f_0, f_1, \ldots, f_K}$, and

(iv) it holds that $\mathcal{P}(\mathbf{F}) = 3K + 4$

*(cf. Definitions 2.1.6, 2.2.3, and 3.1.4).*

*Proof of Lemma 3.1.12.* Throughout this proof let $c_0, c_1, \ldots, c_K \in \mathbb{R}$ satisfy for all $k \in \{0, 1, \ldots, K\}$ that

$$c_k = \frac{(f_{\min\{k+1,K\}} - f_k)}{(\mathfrak{x}_{\min\{k+1,K\}} - \mathfrak{x}_{\min\{k,K-1\}})} - \frac{(f_k - f_{\max\{k-1,0\}})}{(\mathfrak{x}_{\max\{k,1\}} - \mathfrak{x}_{\max\{k-1,0\}})} \tag{3.24}$$

and let $\Phi_0, \Phi_1, \ldots, \Phi_K \in ((\mathbb{R}^{1\times1} \times \mathbb{R}^1) \times (\mathbb{R}^{1\times1} \times \mathbb{R}^1)) \subseteq \mathbf{N}$ satisfy for all $k \in \{0, 1, \ldots, K\}$ that $\Phi_k = c_k \circledast (\mathfrak{i}_1 \bullet \mathbf{A}_{1, -\mathfrak{x}_k})$. Observe that Lemma 3.1.11 assures that for all $k \in \{0, 1, \ldots, K\}$ it holds that $\mathcal{R}_{\mathfrak{r}}(\Phi_k) \in C(\mathbb{R}, \mathbb{R})$, $\mathcal{D}(\Phi_k) = (1, 1, 1) \in \mathbb{N}^3$, and $\forall x \in \mathbb{R} \colon (\mathcal{R}_{\mathfrak{r}}(\Phi_k))(x) = c_k \max\{x - \mathfrak{x}_k, 0\}$ (cf. Definitions 2.1.6 and 2.2.3). This, Lemma 2.2.22, Lemma 2.2.35, and (3.23) assure that $\mathcal{D}(\mathbf{F}) = (1, K+1, 1) \in \mathbb{N}^3$ and $\mathcal{R}_{\mathfrak{r}}(\mathbf{F}) \in C(\mathbb{R}, \mathbb{R})$. This establishes items (i) and (ii). Moreover, note that item (i) and (2.52) imply that

$$\mathcal{P}(\mathbf{F}) = 2(K+1) + (K+2) = 3K + 4. \tag{3.25}$$

This proves item (iv). Next observe that (3.24), Lemma 2.2.22, and Lemma 2.2.35 ensure that for all $x \in \mathbb{R}$ it holds that

$$(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) = f_0 + \sum_{k=0}^{K} (\mathcal{R}_{\mathfrak{r}}(\Phi_k))(x) = f_0 + \sum_{k=0}^{K} c_k \max\{x - \mathfrak{x}_k, 0\}. \tag{3.26}$$

This and the fact that $\forall k \in \{0, 1, \ldots, K\} \colon \mathfrak{x}_0 \leq \mathfrak{x}_k$ assure that for all $x \in (-\infty, \mathfrak{x}_0]$ it holds that

$$(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) = f_0 + 0 = f_0. \tag{3.27}$$

Next we claim that for all $k \in \{1, 2, \ldots, K\}$ it holds that

$$\sum_{n=0}^{k-1} c_n = \frac{f_k - f_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}}. \tag{3.28}$$

We now prove (3.28) by induction on $k \in \{1, 2, \ldots, K\}$. For the base case $k = 1$ observe that (3.24) assures that $\sum_{n=0}^{0} c_n = c_0 = \frac{f_1 - f_0}{\mathfrak{x}_1 - \mathfrak{x}_0}$. This proves (3.28) in the base case $k = 1$. For the induction step note that (3.24) ensures that for all $k \in \{2, 3, \ldots, K\}$ with $\sum_{n=0}^{k-2} c_n = \frac{f_{k-1} - f_{k-2}}{\mathfrak{x}_{k-1} - \mathfrak{x}_{k-2}}$ it holds that

$$\sum_{n=0}^{k-1} c_n = c_{k-1} + \sum_{n=0}^{k-2} c_n = \frac{f_k - f_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}} - \frac{f_{k-1} - f_{k-2}}{\mathfrak{x}_{k-1} - \mathfrak{x}_{k-2}} + \frac{f_{k-1} - f_{k-2}}{\mathfrak{x}_{k-1} - \mathfrak{x}_{k-2}} = \frac{f_k - f_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}}. \tag{3.29}$$

Induction thus proves (3.28). In addition, observe that (3.26), (3.28), and the fact that $\forall k \in \{1, 2, \ldots, K\} \colon \mathfrak{x}_{k-1} < \mathfrak{x}_k$ show that for all $k \in \{1, 2, \ldots, K\}$, $x \in [\mathfrak{x}_{k-1}, \mathfrak{x}_k]$ it holds

that

$$(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - (\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(\mathfrak{x}_{k-1}) = \sum_{n=0}^{K} c_n(\max\{x - \mathfrak{x}_n, 0\} - \max\{\mathfrak{x}_{k-1} - \mathfrak{x}_n, 0\})$$

$$= \sum_{n=0}^{k-1} c_n[(x - \mathfrak{x}_n) - (\mathfrak{x}_{k-1} - \mathfrak{x}_n)] = \sum_{n=0}^{k-1} c_n(x - \mathfrak{x}_{k-1}) \tag{3.30}$$

$$= (\tfrac{f_k - f_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}})(x - \mathfrak{x}_{k-1}).$$

Next we claim that for all $k \in \{1, 2, \ldots, K\}$, $x \in [\mathfrak{x}_{k-1}, \mathfrak{x}_k]$ it holds that

$$(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) = f_{k-1} + (\tfrac{f_k - f_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}})(x - \mathfrak{x}_{k-1}). \tag{3.31}$$

We now prove (3.31) by induction on $k \in \{1, 2, \ldots, K\}$. For the base case $k = 1$ observe that (3.27) and (3.30) demonstrate that for all $x \in [\mathfrak{x}_0, \mathfrak{x}_1]$ it holds that

$$(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) = (\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(\mathfrak{x}_0) + (\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - (\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(\mathfrak{x}_0) = f_0 + (\tfrac{f_1 - f_0}{\mathfrak{x}_1 - \mathfrak{x}_0})(x - \mathfrak{x}_0). \tag{3.32}$$

This proves (3.31) in the base case $k = 1$. For the induction step note that (3.30) implies that for all $k \in \{2, 3, \ldots, K\}$, $x \in [\mathfrak{x}_{k-1}, \mathfrak{x}_k]$ with $\forall y \in [\mathfrak{x}_{k-2}, \mathfrak{x}_{k-1}]: (\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(y) = f_{k-2} + (\tfrac{f_{k-1} - f_{k-2}}{\mathfrak{x}_{k-1} - \mathfrak{x}_{k-2}})(y - \mathfrak{x}_{k-2})$ it holds that

$$(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) = (\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(\mathfrak{x}_{k-1}) + (\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - (\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(\mathfrak{x}_{k-1})$$

$$= f_{k-2} + (\tfrac{f_{k-1} - f_{k-2}}{\mathfrak{x}_{k-1} - \mathfrak{x}_{k-2}})(\mathfrak{x}_{k-1} - \mathfrak{x}_{k-2}) + (\tfrac{f_k - f_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}})(x - \mathfrak{x}_{k-1}) \tag{3.33}$$

$$= f_{k-1} + (\tfrac{f_k - f_{k-1}}{\mathfrak{x}_k - \mathfrak{x}_{k-1}})(x - \mathfrak{x}_{k-1}).$$

Induction thus proves (3.31). Furthermore, observe that (3.24) and (3.28) ensure that

$$\sum_{n=0}^{K} c_n = c_K + \sum_{n=0}^{K-1} c_n = -\tfrac{f_K - f_{K-1}}{\mathfrak{x}_K - \mathfrak{x}_{K-1}} + \tfrac{f_K - f_{K-1}}{\mathfrak{x}_K - \mathfrak{x}_{K-1}} = 0. \tag{3.34}$$

The fact that $\forall\, k \in \{0, 1, \ldots, K\}: \mathfrak{x}_k \leq \mathfrak{x}_K$ and (3.26) hence imply that for all $x \in [\mathfrak{x}_K, \infty)$ it holds that

$$(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - (\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(\mathfrak{x}_K) = \left[\sum_{n=0}^{K} c_n(\max\{x - \mathfrak{x}_n, 0\} - \max\{\mathfrak{x}_K - \mathfrak{x}_n, 0\})\right]$$

$$= \sum_{n=0}^{K} c_n[(x - \mathfrak{x}_n) - (\mathfrak{x}_K - \mathfrak{x}_n)] = \sum_{n=0}^{K} c_n(x - \mathfrak{x}_K) = 0. \tag{3.35}$$

This and (3.31) show that for all $x \in [\mathfrak{x}_K, \infty)$ it holds that

$$(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) = (\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(\mathfrak{x}_K) = f_{K-1} + (\tfrac{f_K - f_{K-1}}{\mathfrak{x}_K - \mathfrak{x}_{K-1}})(\mathfrak{x}_K - \mathfrak{x}_{K-1}) = f_K. \tag{3.36}$$

Combining this, (3.27), (3.31), and (3.4) establishes item (iii). The proof of Lemma 3.1.12 is thus complete. $\qquad\square$

**Exercise 3.1.1.** *Prove or disprove the following statement: There exists $\Phi \in \mathbf{N}$ such that $\mathcal{P}(\Phi) \leq 16$ and*

$$\sup_{x \in [-2\pi, 2\pi]} \left|\cos(x) - (\mathcal{R}_{\mathfrak{r}}(\Phi))(x)\right| \leq \tfrac{1}{2} \tag{3.37}$$

*(cf. Definitions 2.1.6, 2.2.1, and 2.2.3).*

**Exercise 3.1.2.** *Prove or disprove the following statement: There exists $\Phi \in \mathbf{N}$ such that $\mathcal{I}(\Phi) = 4$, $\mathcal{O}(\Phi) = 1$, $\mathcal{P}(\Phi) \leq 60$, and $\forall x, y, u, v \in \mathbb{R}: (\mathcal{R}_\mathfrak{r}(\Phi))(x, y, u, v) = \max\{x, y, u, v\}$ (cf. Definitions 2.1.6, 2.2.1, and 2.2.3).*

**Exercise 3.1.3.** *Prove or disprove the following statement: For every $m \in \mathbb{N}$ there exists $\Phi \in \mathbf{N}$ such that $\mathcal{I}(\Phi) = 2^m$, $\mathcal{O}(\Phi) = 1$, $\mathcal{P}(\Phi) \leq 3(2^m(2^m + 1))$, and $\forall x = (x_1, x_2, \ldots, x_{2^m}) \in \mathbb{R}: (\mathcal{R}_\mathfrak{r}(\Phi))(x) = \max\{x_1, x_2, \ldots, x_{2^m}\}$ (cf. Definitions 2.1.6, 2.2.1, and 2.2.3).*

## 3.1.4 Neural network approximations for one-dimensional functions

**Lemma 3.1.13.** *Let $K \in \mathbb{N}$, $L, a, \mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K \in \mathbb{R}$, $b \in (a, \infty)$ satisfy for all $k \in \{0, 1, \ldots, K\}$ that $\mathfrak{x}_k = a + \frac{k(b-a)}{K}$, let $f \colon [a, b] \to \mathbb{R}$ satisfy for all $x, y \in [a, b]$ that $|f(x) - f(y)| \leq L|x - y|$, and let $\mathbf{F} \in \mathbf{N}$ satisfy*

$$\mathbf{F} = \mathbf{A}_{1, f(\mathfrak{x}_0)} \bullet \left( \bigoplus_{k=0}^{K} \left( \left( \frac{K(f(\mathfrak{x}_{\min\{k+1, K\}}) - 2f(\mathfrak{x}_k) + f(\mathfrak{x}_{\max\{k-1, 0\}}))}{(b-a)} \right) \circledast (\mathfrak{i}_1 \bullet \mathbf{A}_{1, -\mathfrak{x}_k}) \right) \right) \quad (3.38)$$

*(cf. Definitions 2.2.1, 2.2.5, 2.2.20, 2.2.23, 2.2.34, and 3.1.8). Then*

*(i) it holds that $\mathcal{D}(\mathbf{F}) = (1, K + 1, 1)$,*

*(ii) it holds that $\mathcal{R}_\mathfrak{r}(\mathbf{F}) \in C(\mathbb{R}, \mathbb{R})$,*

*(iii) it holds that $\mathcal{R}_\mathfrak{r}(\mathbf{F}) = \mathcal{L}_{\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K}^{f(\mathfrak{x}_0), f(\mathfrak{x}_1), \ldots, f(\mathfrak{x}_K)}$,*

*(iv) it holds for all $x, y \in \mathbb{R}$ that $|(\mathcal{R}_\mathfrak{r}(\mathbf{F}))(x) - (\mathcal{R}_\mathfrak{r}(\mathbf{F}))(y)| \leq L|x - y|$,*

*(v) it holds that $\sup_{x \in [a,b]}|(\mathcal{R}_\mathfrak{r}(\mathbf{F}))(x) - f(x)| \leq L(b-a)K^{-1}$, and*

*(vi) it holds that $\mathcal{P}(\mathbf{F}) = 3K + 4$*

*(cf. Definitions 2.1.6, 2.2.3, and 3.1.4).*

*Proof of Lemma 3.1.13.* Note that the fact that $\forall k \in \{0, 1, \ldots, K\}: \mathfrak{x}_{\min\{k+1, K\}} - \mathfrak{x}_{\min\{k, K-1\}} = \mathfrak{x}_{\max\{k, 1\}} - \mathfrak{x}_{\max\{k-1, 0\}} = (b - a)K^{-1}$ assures that for all $k \in \{0, 1, \ldots, K\}$ it holds that

$$\frac{(f(\mathfrak{x}_{\min\{k+1, K\}}) - f(\mathfrak{x}_k))}{(\mathfrak{x}_{\min\{k+1, K\}} - \mathfrak{x}_{\min\{k, K-1\}})} - \frac{(f(\mathfrak{x}_k) - f(\mathfrak{x}_{\max\{k-1, 0\}}))}{(\mathfrak{x}_{\max\{k, 1\}} - \mathfrak{x}_{\max\{k-1, 0\}})} = \frac{K(f(\mathfrak{x}_{\min\{k+1, K\}}) - 2f(\mathfrak{x}_k) + f(\mathfrak{x}_{\max\{k-1, 0\}}))}{(b-a)}. \quad (3.39)$$

This and items (i), (ii), (iii), and (iv) in Lemma 3.1.12 prove items (i), (ii), (iii), and (vi). Combining item (iii) with the assumption that $\forall x, y \in [a, b]: |f(x) - f(y)| \leq L|x - y|$ and item (i) in Lemma 3.1.7 establishes item (iv). Moreover, note that item (iii), the assumption that $\forall x, y \in [a, b]: |f(x) - f(y)| \leq L|x - y|$, item (ii) in Lemma 3.1.7, and the fact that $\forall k \in \{1, 2, \ldots, K\}: \mathfrak{x}_k - \mathfrak{x}_{k-1} = (b - a)K^{-1}$ demonstrate that for all $x \in [a, b]$ it holds that

$$|(\mathcal{R}_\mathfrak{r}(\mathbf{F}))(x) - f(x)| \leq L\left( \max_{k \in \{1, 2, \ldots, K\}} |\mathfrak{x}_k - \mathfrak{x}_{k-1}| \right) = L(b - a)K^{-1}. \quad (3.40)$$

This establishes item (v). The proof of Lemma 3.1.13 is thus complete. $\square$

**Lemma 3.1.14.** *Let $L, a \in \mathbb{R}$, $b \in [a, \infty)$, $\xi \in [a, b]$, let $f\colon [a,b] \to \mathbb{R}$ satisfy for all $x, y \in [a,b]$ that $|f(x) - f(y)| \leq L|x-y|$, and let $\mathbf{F} \in \mathbf{N}$ satisfy $\mathbf{F} = \mathbf{A}_{1,f(\xi)} \bullet (0 \circledast (\mathfrak{i}_1 \bullet \mathbf{A}_{1,-\xi}))$ (cf. Definitions 2.2.1, 2.2.5, 2.2.20, 2.2.23, and 3.1.8). Then*

*(i) it holds that $\mathcal{D}(\mathbf{F}) = (1, 1, 1)$,*

*(ii) it holds that $\mathcal{R}_{\mathfrak{r}}(\mathbf{F}) \in C(\mathbb{R}, \mathbb{R})$,*

*(iii) it holds for all $x \in \mathbb{R}$ that $(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) = f(\xi)$,*

*(iv) it holds that $\sup_{x \in [a,b]} |(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - f(x)| \leq L \max\{\xi - a, b - \xi\}$, and*

*(v) it holds that $\mathcal{P}(\mathbf{F}) = 4$*

*(cf. Definitions 2.1.6 and 2.2.3).*

*Proof of Lemma 3.1.14.* Note that items (i) and (ii) in Lemma 2.2.22, and items (ii) and (iii) in Lemma 3.1.11 establish items (i) and (ii). In addition, observe that item (iii) in Lemma 2.2.22 and item (iii) in Lemma 2.2.24 assure that for all $x \in \mathbb{R}$ it holds that

$$\begin{aligned} (\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) &= (\mathcal{R}_{\mathfrak{r}}(0 \circledast (\mathfrak{i}_1 \bullet \mathbf{A}_{1,-\xi})))(x) + f(\xi) \\ &= 0\big((\mathcal{R}_{\mathfrak{r}}(\mathfrak{i}_1 \bullet \mathbf{A}_{1,-\xi}))(x)\big) + f(\xi) = f(\xi) \end{aligned} \tag{3.41}$$

(cf. Definitions 2.1.6 and 2.2.3). This establishes item (iii). Next note that (3.41), the fact that $\xi \in [a,b]$, and the fact that for all $x, y \in [a,b]$ it holds that $|f(x) - f(y)| \leq L|x-y|$ assure that for all $x \in [a,b]$ it holds that

$$|(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - f(x)| = |f(\xi) - f(x)| \leq L|x - \xi| \leq L \max\{\xi - a, b - \xi\}. \tag{3.42}$$

This establishes item (iv). Moreover, note that (2.52) and item (i) assure that

$$\mathcal{P}(\mathbf{F}) = 1(1 + 1) + 1(1 + 1) = 4. \tag{3.43}$$

This establishes item (v). The proof of Lemma 3.1.14 it thus completed. $\square$

**Corollary 3.1.15.** *Let $\varepsilon \in (0, \infty)$, $L, a \in \mathbb{R}$, $b \in (a, \infty)$, $K \in \mathbb{N}_0 \cap [\frac{L(b-a)}{\varepsilon}, \frac{L(b-a)}{\varepsilon} + 1)$, $\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K \in \mathbb{R}$ satisfy for all $k \in \{0, 1, \ldots, K\}$ that $\mathfrak{x}_k = a + \frac{k(b-a)}{\max\{K,1\}}$, let $f\colon [a,b] \to \mathbb{R}$ satisfy for all $x, y \in [a,b]$ that $|f(x) - f(y)| \leq L|x-y|$, and let $\mathbf{F} \in \mathbf{N}$ satisfy*

$$\mathbf{F} = \mathbf{A}_{1,f(\mathfrak{x}_0)} \bullet \left( \bigoplus_{k=0}^{K} \left( \left( \tfrac{K(f(\mathfrak{x}_{\min\{k+1,K\}}) - 2f(\mathfrak{x}_k) + f(\mathfrak{x}_{\max\{k-1,0\}}))}{(b-a)} \right) \circledast (\mathfrak{i}_1 \bullet \mathbf{A}_{1,-\mathfrak{x}_k}) \right) \right) \tag{3.44}$$

*(cf. Definitions 2.2.1, 2.2.5, 2.2.20, 2.2.23, 2.2.34, and 3.1.8). Then*

*(i) it holds that $\mathcal{D}(\mathbf{F}) = (1, K + 1, 1)$,*

*(ii) it holds that $\mathcal{R}_{\mathfrak{r}}(\mathbf{F}) \in C(\mathbb{R}, \mathbb{R})$,*

*(iii) it holds for all $x, y \in \mathbb{R}$ that $|(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - (\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(y)| \leq L|x-y|$,*

*(iv) it holds that $\sup_{x \in [a,b]} |(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - f(x)| \leq \frac{L(b-a)}{\max\{K,1\}} \leq \varepsilon$, and*

*(v) it holds that $\mathcal{P}(\mathbf{F}) = 3K + 4 \leq 3L(b-a)\varepsilon^{-1} + 7$*

*(cf. Definitions 2.1.6, 2.2.1, and 2.2.3).*

*Proof of Corollary 3.1.15.* Note that the fact that $K \in \mathbb{N}_0 \cap [\frac{L(b-a)}{\varepsilon}, \frac{L(b-a)}{\varepsilon} + 1)$ implies that $\frac{L(b-a)}{\max\{K,1\}} \leq \varepsilon$. This, items (i), (ii), (iv), and (v) in Lemma 3.1.13, and items (i), (ii), (iii), and (iv) in Lemma 3.1.14 establish items (i), (ii), (iii), and (iv). Moreover, note that the fact that $K \leq 1 + \frac{L(b-a)}{\varepsilon}$, item (vi) in Lemma 3.1.13, and item (v) in Lemma 3.1.14 assure that

$$\mathcal{P}(\mathbf{F}) = 3K + 4 \leq \frac{3L(b-a)}{\varepsilon} + 7. \tag{3.45}$$

This establishes item (v). The proof of Corollary 3.1.15 is thus complete. $\qquad\square$

**Definition 3.1.16** (*p*-norm). *We denote by $\|\cdot\|_p \colon (\bigcup_{d=1}^{\infty} \mathbb{R}^d) \to \mathbb{R}$, $p \in [1, \infty]$, the functions which satisfy for all $p \in [1, \infty)$, $d \in \mathbb{N}$, $\theta = (\theta_1, \theta_2, \ldots, \theta_d) \in \mathbb{R}^d$ that $\|\theta\|_p = [\sum_{i=1}^{d} |\theta_i|^p]^{1/p}$ and $\|\theta\|_{\infty} = \max_{i \in \{1,2,\ldots,d\}} |\theta_i|$.*

**Corollary 3.1.17.** *Let $\varepsilon \in (0, \infty)$, $L \in [0, \infty)$, $a \in \mathbb{R}$, $b \in [a, \infty)$ and let $f \colon [a,b] \to \mathbb{R}$ satisfy for all $x, y \in [a,b]$ that $|f(x) - f(y)| \leq L|x-y|$. Then there exists $\mathbf{F} \in \mathbf{N}$ such that*

(i) *it holds that $\mathcal{R}_{\mathfrak{r}}(\mathbf{F}) \in C(\mathbb{R}, \mathbb{R})$,*

(ii) *it holds that $\mathcal{H}(\mathbf{F}) = 1$,*

(iii) *it holds that $\mathbb{D}_1(\mathbf{F}) \leq L(b-a)\varepsilon^{-1} + 2$,*

(iv) *it holds for all $x, y \in \mathbb{R}$ that $|(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - (\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(y)| \leq L|x-y|$,*

(v) *it holds that $\sup_{x \in [a,b]} |(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - f(x)| \leq \varepsilon$,*

(vi) *it holds that $\mathcal{P}(\mathbf{F}) = 3(\mathbb{D}_1(\mathbf{F})) + 1 \leq 3L(b-a)\varepsilon^{-1} + 7$, and*

(vii) *it holds that $\|\mathcal{T}(\mathbf{F})\|_{\infty} \leq \max\{1, |a|, |b|, 2L, |f(a)|\}$*

*(cf. Definitions 2.1.6, 2.2.1, 2.2.3, 2.2.36, and 3.1.16).*

*Proof of Corollary 3.1.17.* Throughout this proof assume w.l.o.g. that $a < b$, let $K \in \mathbb{N}_0 \cap [\frac{L(b-a)}{\varepsilon}, \frac{L(b-a)}{\varepsilon} + 1)$, $\mathfrak{x}_0, \mathfrak{x}_1, \ldots, \mathfrak{x}_K \in \mathbb{R}$, $c_0, c_1, \ldots, c_K \in \mathbb{R}$ satisfy for all $k \in \{0, 1, \ldots, K\}$ that $\mathfrak{x}_k = a + \frac{k(b-a)}{\max\{K,1\}}$ and

$$c_k = \frac{K(f(\mathfrak{x}_{\min\{k+1,K\}}) - 2f(\mathfrak{x}_k) + f(\mathfrak{x}_{\max\{k-1,0\}}))}{(b-a)}, \tag{3.46}$$

and let $\mathbf{F} \in \mathbf{N}$ satisfy

$$\mathbf{F} = \mathbf{A}_{1,f(\mathfrak{x}_0)} \bullet \left( \bigoplus_{k=0}^{K} (c_k \circledast (\mathfrak{i}_1 \bullet \mathbf{A}_{1,-\mathfrak{x}_k})) \right) \tag{3.47}$$

*(cf. Definitions 2.2.1, 2.2.5, 2.2.20, 2.2.23, 2.2.34, and 3.1.8).* Note that Corollary 3.1.15 implies that

(I) *it holds that $\mathcal{D}(\mathbf{F}) = (1, K+1, 1)$,*

(II) *it holds that $\mathcal{R}_{\mathfrak{r}}(\mathbf{F}) \in C(\mathbb{R}, \mathbb{R})$,*

(III) it holds for all $x, y \in \mathbb{R}$ that $|(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - (\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(y)| \leq L|x - y|$,

(IV) it holds that $\sup_{x \in [a,b]}|(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - f(x)| \leq \varepsilon$, and

(V) it holds that $\mathcal{P}(\mathbf{F}) = 3K + 4$

(cf. Definitions 2.1.6 and 2.2.3). This establishes items (i), (iv), and (v). Next note that item (I) and the fact that $K \leq 1 + \frac{L(b-a)}{\varepsilon}$ prove items (ii) and (iii). Next observe that items (I) and (V) imply that

$$\mathcal{P}(\mathbf{F}) = 3K + 4 = 3(K + 1) + 1 = 3(\mathbb{D}_1(\mathbf{F})) + 1 \leq \tfrac{3L(b-a)}{\varepsilon} + 7. \tag{3.48}$$

This establishes item (vi). In the next step we observe that Lemma 3.1.11 shows that for all $k \in \{0, 1, \ldots, K\}$ it holds that

$$c_k \circledast (\mathfrak{i}_1 \bullet \mathbf{A}_{1,-\mathfrak{x}_k}) = ((1, -\mathfrak{x}_k), (c_k, 0)). \tag{3.49}$$

Combining this with (2.168), (2.161), (2.154), and Lemma 2.2.6 demonstrates that

$$\mathbf{F} = \left( \left( \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}, \begin{pmatrix} -\mathfrak{x}_0 \\ -\mathfrak{x}_1 \\ \vdots \\ -\mathfrak{x}_K \end{pmatrix} \right), \left( \begin{pmatrix} c_0 & c_1 & \cdots & c_K \end{pmatrix}, f(\mathfrak{x}_0) \right) \right) \tag{3.50}$$
$$\in (\mathbb{R}^{(K+1)\times 1} \times \mathbb{R}^{K+1}) \times (\mathbb{R}^{1 \times (K+1)} \times \mathbb{R}).$$

Lemma 2.2.38 therefore ensures that

$$\|\mathcal{T}(\mathbf{F})\|_\infty = \max\{|\mathfrak{x}_0|, |\mathfrak{x}_1|, \ldots, |\mathfrak{x}_K|, |c_0|, |c_1|, \ldots, |c_K|, |f(\mathfrak{x}_0)|, 1\} \tag{3.51}$$

(cf. Definitions 2.2.36 and 3.1.16). In addition, note that the assumption that for all $x, y \in [a,b]$ it holds that $|f(x) - f(y)| \leq L|x - y|$ and the fact that $\forall\, k \in \mathbb{N} \cap (0, K+1)\colon \mathfrak{x}_k - \mathfrak{x}_{k-1} = (b-a)[\max\{K, 1\}]^{-1}$ imply that for all $k \in \{0, 1, \ldots, K\}$ it holds that

$$\begin{aligned}
|c_k| &\leq \frac{K(|f(\mathfrak{x}_{\min\{k+1,K\}}) - f(\mathfrak{x}_k)| + |f(\mathfrak{x}_{\max\{k-1,0\}})) - f(\mathfrak{x}_k)|}{(b-a)} \\
&\leq \frac{KL(|\mathfrak{x}_{\min\{k+1,K\}} - \mathfrak{x}_k| + |\mathfrak{x}_{\max\{k-1,0\}} - \mathfrak{x}_k|)}{(b-a)} \\
&\leq \frac{2KL(b-a)[\max\{K, 1\}]^{-1}}{(b-a)} \leq 2L.
\end{aligned} \tag{3.52}$$

This and (3.51) establish item (vii). The proof of Corollary 3.1.17 is thus complete. $\qquad\square$

**Corollary 3.1.18.** *Let $L, a \in \mathbb{R}$, $b \in [a, \infty)$ and let $f\colon [a,b] \to \mathbb{R}$ satisfy for all $x, y \in [a,b]$ that $|f(x) - f(y)| \leq L|x - y|$. Then there exist $C \in \mathbb{R}$ and $\mathbf{F} = (\mathbf{F}_\varepsilon)_{\varepsilon \in (0,1]}\colon (0,1] \to \mathbf{N}$ such that for all $\varepsilon \in (0,1]$ it holds that $\mathcal{R}_{\mathfrak{r}}(\mathbf{F}_\varepsilon) \in C(\mathbb{R}, \mathbb{R})$, $\sup_{x \in [a,b]}|(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}_\varepsilon))(x) - f(x)| \leq \varepsilon$, $\mathcal{H}(\mathbf{F}_\varepsilon) = 1$, $\|\mathcal{T}(\mathbf{F}_\varepsilon)\|_\infty \leq \max\{1, |a|, |b|, 2L, |f(a)|\}$, and $\mathcal{P}(\mathbf{F}_\varepsilon) \leq C\varepsilon^{-1}$ (cf. Definitions 2.1.6, 2.2.1, 2.2.3, 2.2.36, and 3.1.16).*

*Proof of Corollary 3.1.18.* Throughout this proof assume w.l.o.g. that $L \geq 0$ and let $C = 3L(b-a) + 7$. Observe that for all $\varepsilon \in (0,1]$ it holds that

$$3L(b-a)\varepsilon^{-1} + 7 \leq 3L(b-a)\varepsilon^{-1} + 7\varepsilon^{-1} = C\varepsilon^{-1}. \tag{3.53}$$

This and Corollary 3.1.17 establish that there exists $\mathbf{F} = (\mathbf{F}_\varepsilon)_{\varepsilon \in (0,1]} \colon (0,1] \to \mathbf{N}$ such that for all $\varepsilon \in (0,1]$ it holds that $\mathcal{R}_\mathfrak{r}(\mathbf{F}_\varepsilon) \in C(\mathbb{R}, \mathbb{R})$, $\sup_{x \in [a,b]} |(\mathcal{R}_\mathfrak{r}(\mathbf{F}_\varepsilon))(x) - f(x)| \leq \varepsilon$, $\mathcal{H}(\mathbf{F}_\varepsilon) = 1$, $\|\mathcal{T}(\mathbf{F}_\varepsilon)\|_\infty \leq \max\{|a|, |b|, 2L, |f(a)|\}$, and

$$\mathcal{P}(\mathbf{F}_\varepsilon) \leq 3L(b-a)\varepsilon^{-1} + 7 \leq C\varepsilon^{-1} \tag{3.54}$$

(cf. Definitions 2.1.6, 2.2.1, 2.2.3, 2.2.36, and 3.1.16). The proof of Corollary 3.1.18 is thus complete. $\qquad\square$

**Exercise 3.1.4.** *Prove or disprove the following statement: There exists $\Phi \in \mathbf{N}$ such that $\mathcal{P}(\Phi) \leq 10$ and*

$$\sup_{x \in [0,10]} \left| \sqrt{x} - (\mathcal{R}_\mathfrak{r}(\Phi))(x) \right| \leq \tfrac{1}{4} \tag{3.55}$$

*(cf. Definitions 2.1.6, 2.2.1, and 2.2.3).*

# 3.2 Multi-dimensional ANN approximation results

## 3.2.1 Approximations for Lipschitz continuous functions

**Lemma 3.2.1.** *Let $(E, \delta)$ be a metric space, let $L \in [0, \infty)$, $D \subseteq E$, $\mathcal{M} \subseteq E$ satisfy $\emptyset \neq \mathcal{M} \subseteq D$, let $f \colon D \to \mathbb{R}$ satisfy for all $x \in D$, $y \in \mathcal{M}$ that $|f(x) - f(y)| \leq L\delta(x,y)$, and let $F \colon E \to \mathbb{R} \cup \{\infty\}$ satisfy for all $x \in E$ that*

$$F(x) = \sup_{y \in \mathcal{M}} [f(y) - L\delta(x,y)]. \tag{3.56}$$

*Then*

*(i) it holds for all $x \in \mathcal{M}$ that $F(x) = f(x)$,*

*(ii) it holds for all $x \in D$ that $F(x) \leq f(x)$,*

*(iii) it holds for all $x \in E$ that $F(x) < \infty$,*

*(iv) it holds for all $x, y \in E$ that $|F(x) - F(y)| \leq L\delta(x,y)$, and*

*(v) it holds for all $x \in D$ that*

$$|F(x) - f(x)| \leq 2L \left[ \inf_{y \in \mathcal{M}} \delta(x,y) \right]. \tag{3.57}$$

*Proof of Lemma 3.2.1.* First, observe that the assumption that $\forall\, x \in D$, $y \in \mathcal{M} : |f(x) - f(y)| \leq L\delta(x,y)$ ensures that for all $x \in D$, $y \in \mathcal{M}$ it holds that

$$f(y) + L\delta(x,y) \geq f(x) \geq f(y) - L\delta(x,y). \tag{3.58}$$

Hence, we obtain that for all $x \in D$ it holds that

$$f(x) \geq \sup_{y \in \mathcal{M}} [f(y) - L\delta(x,y)] = F(x). \tag{3.59}$$

This establishes item (ii). Moreover, note that (3.56) implies that for all $x \in \mathcal{M}$ it holds that

$$F(x) \geq f(x) - L\delta(x,x) = f(x). \tag{3.60}$$

This and (3.59) establish item (i). Next note that (3.58) (applied for all $y, z \in \mathcal{M}$ with $x \curvearrowright y$, $y \curvearrowright z$) and the triangle inequality ensure that for all $x \in E$, $y, z \in \mathcal{M}$ it holds that

$$f(y) - L\delta(x,y) \leq f(z) + L\delta(y,z) - L\delta(x,y) \leq f(z) + L\delta(x,z). \tag{3.61}$$

Hence, we obtain that for all $x \in E$, $z \in \mathcal{M}$ it holds that

$$F(x) = \sup_{y \in \mathcal{M}}[f(y) - L\delta(x,y)] \leq f(z) + L\delta(x,z) < \infty. \tag{3.62}$$

This proves item (iii). Combining item (iii) with (3.56) and the triangle inequality shows that for all $x, y \in E$ it holds that

$$
\begin{aligned}
F(x) - F(y) &= \left[\sup_{v \in \mathcal{M}}(f(v) - L\delta(x,v))\right] - \left[\sup_{w \in \mathcal{M}}(f(w) - L\delta(y,w))\right] \\
&= \sup_{v \in \mathcal{M}}\left[f(v) - L\delta(x,v) - \sup_{w \in \mathcal{M}}(f(w) - L\delta(y,w))\right] \\
&\leq \sup_{v \in \mathcal{M}}\left[f(v) - L\delta(x,v) - (f(v) - L\delta(y,v))\right] \\
&= \sup_{v \in \mathcal{M}}(L\delta(y,v) - L\delta(x,v)) \\
&\leq \sup_{v \in \mathcal{M}}(L\delta(y,x) + L\delta(x,v) - L\delta(x,v)) = L\delta(x,y).
\end{aligned}
\tag{3.63}
$$

This and the fact that for all $x, y \in E$ it holds that $\delta(x,y) = \delta(y,x)$ establish item (iv). Moreover, observe that items (i) and (iv), the triangle inequality, and the assumption that $\forall\, x \in D, y \in \mathcal{M} \colon |f(x) - f(y)| \leq L\delta(x,y)$ ensure that for all $x \in D$ it holds that

$$
\begin{aligned}
|F(x) - f(x)| &= \inf_{y \in \mathcal{M}}|F(x) - F(y) + f(y) - f(x)| \\
&\leq \inf_{y \in \mathcal{M}}(|F(x) - F(y)| + |f(y) - f(x)|) \\
&\leq \inf_{y \in \mathcal{M}}(2L\delta(x,y)) = 2L\left[\inf_{y \in \mathcal{M}}\delta(x,y)\right].
\end{aligned}
\tag{3.64}
$$

This establishes item (v). The proof of Lemma 3.2.1 is thus complete. □

**Corollary 3.2.2.** *Let $(E, \delta)$ be a metric space, let $L \in [0, \infty)$, $\mathcal{M} \subseteq E$ satisfy $\mathcal{M} \neq \emptyset$, let $f\colon E \to \mathbb{R}$ satisfy for all $x \in E$, $y \in \mathcal{M}$ that $|f(x) - f(y)| \leq L\delta(x,y)$, and let $F\colon E \to \mathbb{R} \cup \{\infty\}$ satisfy for all $x \in E$ that*

$$F(x) = \sup_{y \in \mathcal{M}}[f(y) - L\delta(x,y)]. \tag{3.65}$$

*Then*

*(i) it holds for all $x \in \mathcal{M}$ that $F(x) = f(x)$,*

*(ii) it holds for all $x \in E$ that $F(x) \leq f(x)$,*

*(iii) it holds for all $x, y \in E$ that $|F(x) - F(y)| \leq L\delta(x,y)$, and*

*(iv) it holds for all $x \in E$ that*

$$|F(x) - f(x)| \leq 2L\left[\inf_{y \in \mathcal{M}} \delta(x,y)\right]. \tag{3.66}$$

*Proof of Corollary 3.2.2.* Note that Lemma 3.2.1 establishes items (i), (ii), (iii), and (iv). The proof of Corollary 3.2.2 is thus complete. $\square$

**Exercise 3.2.1.** *Prove or disprove the following statement: There exists $\Phi \in \mathbf{N}$ such that $\mathcal{I}(\Phi) = 2$, $\mathcal{O}(\Phi) = 1$, $\mathcal{P}(\Phi) \leq 60\,000\,000$, and*

$$\sup_{x,y \in [0,2\pi]} |\sin(x)\sin(y) - (\mathcal{R}_{\mathfrak{r}}(\Phi))(x,y)| \leq \tfrac{1}{5}. \tag{3.67}$$

## 3.2.2  ANN representations

### 3.2.2.1  ANN representations for the 1-norm

**Definition 3.2.3** (1-norm ANN representations). *We denote by $(\mathbb{L}_d)_{d \in \mathbb{N}} \subseteq \mathbf{N}$ the neural networks which satisfy that*

*(i) it holds that*

$$\mathbb{L}_1 = \left(\left(\begin{pmatrix} 1 \\ -1 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \end{pmatrix}\right), \left((1 \quad 1), (0)\right)\right) \in (\mathbb{R}^{2 \times 1} \times \mathbb{R}^2) \times (\mathbb{R}^{1 \times 2} \times \mathbb{R}^1) \tag{3.68}$$

*and*

*(ii) it holds for all $d \in \{2,3,4,\ldots\}$ that $\mathbb{L}_d = \mathbb{S}_{1,d} \bullet \mathbf{P}_d(\mathbb{L}_1, \mathbb{L}_1, \ldots, \mathbb{L}_1)$*

*(cf. Definitions 2.2.1, 2.2.5, 2.2.11, and 2.2.25).*

**Proposition 3.2.4.** *Let $d \in \mathbb{N}$. Then*

*(i) it holds that $\mathcal{D}(\mathbb{L}_d) = (d, 2d, 1)$,*

*(ii) it holds that $\mathcal{R}_{\mathfrak{r}}(\mathbb{L}_d) \in C(\mathbb{R}^d, \mathbb{R})$, and*

*(iii) it holds for all $x \in \mathbb{R}^d$ that $(\mathcal{R}_{\mathfrak{r}}(\mathbb{L}_d))(x) = \|x\|_1$*

*(cf. Definitions 2.1.6, 2.2.1, 2.2.3, 3.1.16, and 3.2.3).*

*Proof of Proposition 3.2.4.* Note that the fact that $\mathcal{D}(\mathbb{L}_1) = (1, 2, 1)$ and Lemma 2.2.12 show that for all $\mathfrak{d} \in \{2, 3, 4, \ldots\}$ it holds that $\mathcal{D}(\mathbf{P}_{\mathfrak{d}}(\mathbb{L}_1, \mathbb{L}_1, \ldots, \mathbb{L}_1)) = (\mathfrak{d}, 2\mathfrak{d}, \mathfrak{d})$ (cf. Definitions 2.2.1, 2.2.11, and 3.2.3). Combining this, Proposition 2.2.7, and Lemma 2.2.21 ensures that for all $\mathfrak{d} \in \{2, 3, 4, \ldots\}$ it holds that $\mathcal{D}(\mathbb{S}_{1,\mathfrak{d}} \bullet \mathbf{P}_{\mathfrak{d}}(\mathbb{L}_1, \mathbb{L}_1, \ldots, \mathbb{L}_1)) = (\mathfrak{d}, 2\mathfrak{d}, 1)$ (cf. Definitions 2.2.5 and 2.2.25). This establishes item (i). Furthermore, observe that (3.68) assures that for all $x \in \mathbb{R}$ it holds that

$$(\mathcal{R}_{\mathfrak{r}}(\mathbb{L}_1))(x) = \mathfrak{r}(x) + \mathfrak{r}(-x) = \max\{x, 0\} + \max\{-x, 0\} = |x| = \|x\|_1 \tag{3.69}$$

(cf. Definitions 2.1.6, 2.2.3, and 3.1.16). Combining this and Proposition 2.2.13 shows that for all $\mathfrak{d} \in \{2, 3, 4, \ldots\}$, $x = (x_1, x_2, \ldots, x_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ it holds that

$$\left(\mathcal{R}_{\mathfrak{r}}(\mathbf{P}_{\mathfrak{d}}(\mathbb{L}_1, \mathbb{L}_1, \ldots, \mathbb{L}_1))\right)(x) = (|x_1|, |x_2|, \ldots, |x_{\mathfrak{d}}|). \tag{3.70}$$

This and Lemma 2.2.26 demonstrate that for all $\mathfrak{d} \in \{2, 3, 4, \ldots\}$, $x = (x_1, x_2, \ldots, x_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ it holds that

$$
\begin{aligned}
(\mathcal{R}_{\mathfrak{r}}(\mathbb{L}_{\mathfrak{d}}))(x) &= \big(\mathcal{R}_{\mathfrak{r}}(\mathbb{S}_{1,\mathfrak{d}} \bullet \mathbf{P}_{\mathfrak{d}}(\mathbb{L}_1, \mathbb{L}_1, \ldots, \mathbb{L}_1))\big)(x) \\
&= \big(\mathcal{R}_{\mathfrak{r}}(\mathbb{S}_{1,\mathfrak{d}})\big)(|x_1|, |x_2|, \ldots, |x_{\mathfrak{d}}|) = \sum_{n=1}^{d} |x_n| = \|x\|_1.
\end{aligned}
\tag{3.71}
$$

This establishes items (ii)–(iii). The proof of Proposition 3.2.4 is thus complete. $\qquad\square$

**Lemma 3.2.5.** *Let $d \in \mathbb{N}$. Then*

*(i) it holds that $\mathcal{B}_{1,\mathbb{L}_d} = 0 \in \mathbb{R}^{2d}$,*

*(ii) it holds that $\mathcal{B}_{2,\mathbb{L}_d} = 0 \in \mathbb{R}$,*

*(iii) it holds that $\mathcal{W}_{1,\mathbb{L}_d} \in \{-1, 0, 1\}^{(2d) \times d}$,*

*(iv) it holds for all $x \in \mathbb{R}^d$ that $\|\mathcal{W}_{1,\mathbb{L}_d} x\|_\infty = \|x\|_\infty$, and*

*(v) it holds that $\mathcal{W}_{2,\mathbb{L}_d} = (1\ 1\ \cdots\ 1) \in \mathbb{R}^{1 \times (2d)}$*

*(cf. Definitions 2.2.1, 3.1.16, and 3.2.3).*

*Proof of Lemma 3.2.5.* Throughout this proof assume w.l.o.g. that $d > 1$. Note that the fact that $\mathcal{B}_{1,\mathbb{L}_1} = 0 \in \mathbb{R}^2$, the fact that $\mathcal{B}_{2,\mathbb{L}_1} = 0 \in \mathbb{R}$, the fact that $\mathcal{B}_{1,\mathbb{S}_{1,d}} = 0 \in \mathbb{R}$, and the fact that $\mathbb{L}_d = \mathbb{S}_{1,d} \bullet \mathbf{P}_d(\mathbb{L}_1, \mathbb{L}_1, \ldots, \mathbb{L}_1)$ establish items (i)–(ii) (cf. Definitions 2.2.1, 2.2.5, 2.2.11, 2.2.25, and 3.2.3). In addition, observe that the fact that

$$
\mathcal{W}_{1,\mathbb{L}_1} = \begin{pmatrix} 1 \\ -1 \end{pmatrix} \qquad \text{and} \qquad \mathcal{W}_{1,\mathbb{L}_d} = \begin{pmatrix} \mathcal{W}_{1,\mathbb{L}_1} & 0 & \cdots & 0 \\ 0 & \mathcal{W}_{1,\mathbb{L}_1} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathcal{W}_{1,\mathbb{L}_1} \end{pmatrix} \in \mathbb{R}^{(2d) \times d} \tag{3.72}
$$

proves item (iii). Next note that (3.72) implies item (iv). Moreover, note that the fact that $\mathcal{B}_{2,\mathbb{L}_1} = (1\ 1)$ and the fact that $\mathbb{L}_d = \mathbb{S}_{1,d} \bullet \mathbf{P}_d(\mathbb{L}_1, \mathbb{L}_1, \ldots, \mathbb{L}_1)$ show that

$$
\mathcal{W}_{2,\mathbb{L}_d} = \underbrace{\begin{pmatrix} 1 & 1 & \cdots & 1 \end{pmatrix}}_{\in \mathbb{R}^{1 \times d}} \underbrace{\begin{pmatrix} \mathcal{W}_{2,\mathbb{L}_1} & 0 & \cdots & 0 \\ 0 & \mathcal{W}_{2,\mathbb{L}_1} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathcal{W}_{2,\mathbb{L}_1} \end{pmatrix}}_{\in \mathbb{R}^{d \times (2d)}} = \begin{pmatrix} 1 & 1 & \cdots & 1 \end{pmatrix} \in \mathbb{R}^{1 \times (2d)}.
$$

$$
\tag{3.73}
$$

This establishes item (v). The proof of Lemma 3.2.5 is thus complete. $\qquad\square$

### 3.2.2.2   ANN representations for maxima

**Lemma 3.2.6.** *There exist unique $(\phi_d)_{d \in \mathbb{N}} \subseteq \mathbf{N}$ which satisfy that*

*(i) it holds for all $d \in \mathbb{N}$ that $\mathcal{I}(\phi_d) = d$,*

*(ii) it holds for all $d \in \mathbb{N}$ that $\mathcal{O}(\phi_d) = 1$,*

(iii) *it holds that $\phi_1 = \mathbf{A}_{1,0} \in \mathbb{R}^{1\times 1} \times \mathbb{R}^1$,*

(iv) *it holds that*

$$\phi_2 = \left(\left(\begin{pmatrix} 1 & -1 \\ 0 & 1 \\ 0 & -1 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}\right), \left(\begin{pmatrix} 1 & 1 & -1 \end{pmatrix}, (0)\right)\right) \in (\mathbb{R}^{3\times 2} \times \mathbb{R}^3) \times (\mathbb{R}^{1\times 3} \times \mathbb{R}^1),$$

(3.74)

(v) *it holds for all $d \in \{2,3,4,\ldots\}$ that $\phi_{2d} = \phi_d \bullet \left(\mathbf{P}_d(\phi_2, \phi_2, \ldots, \phi_2)\right)$, and*

(vi) *it holds for all $d \in \{2,3,4,\ldots\}$ that $\phi_{2d-1} = \phi_d \bullet \left(\mathbf{P}_d(\phi_2, \phi_2, \ldots, \phi_2, \mathfrak{I}_1)\right)$*

*(cf. Definitions 2.2.1, 2.2.5, 2.2.11, 2.2.18, and 2.2.20).*

*Proof of Lemma 3.2.6.* Throughout this proof let $\psi \in \mathbf{N}$ satisfy

$$\psi = \left(\left(\begin{pmatrix} 1 & -1 \\ 0 & 1 \\ 0 & -1 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}\right), \left(\begin{pmatrix} 1 & 1 & -1 \end{pmatrix}, (0)\right)\right) \in (\mathbb{R}^{3\times 2} \times \mathbb{R}^3) \times (\mathbb{R}^{1\times 3} \times \mathbb{R}^1) \quad (3.75)$$

(cf. Definition 2.2.1). Note that the fact that $\mathcal{I}(\psi) = 2$, the fact that $\mathcal{O}(\psi) = 1$, the fact that $\mathcal{L}(\psi) = \mathcal{L}(\mathfrak{I}_1) = 2$, Lemma 2.2.12, and Lemma 2.2.19 assure that for all $d \in \mathbb{N}$ it holds that $\mathcal{I}(\mathbf{P}_d(\psi, \psi, \ldots, \psi)) = 2d$, $\mathcal{O}(\mathbf{P}_d(\psi, \psi, \ldots, \psi)) = d$, $\mathcal{I}(\mathbf{P}_d(\psi, \psi, \ldots, \psi, \mathfrak{I}_1)) = 2d-1$, and $\mathcal{O}(\mathbf{P}_d(\psi, \psi, \ldots, \psi, \mathfrak{I}_1)) = d$ (cf. Definitions 2.2.11 and 2.2.18). This, Proposition 2.2.7, and induction establish that there exists unique $\phi_d \in \mathbf{N}$, $d \in \mathbb{N}$, which satisfy that for all $d \in \mathbb{N}$ it holds that $\mathcal{I}(\phi_d) = d$, $\mathcal{O}(\phi_d) = 1$, and

$$\phi_d = \begin{cases} \mathbf{A}_{1,0} & : d = 1 \\ \psi & : d = 2 \\ \phi_{d/2} \bullet \left(\mathbf{P}_{d/2}(\psi, \psi, \ldots, \psi)\right) & : d \in \{4, 6, 8, \ldots\} \\ \phi_{(d+1)/2} \bullet \left(\mathbf{P}_{(d+1)/2}(\psi, \psi, \ldots, \psi, \mathfrak{I}_1)\right) & : d \in \{3, 5, 7, \ldots\}. \end{cases} \quad (3.76)$$

The proof of Lemma 3.2.6 is thus complete. $\qquad\square$

**Definition 3.2.7** (Maxima ANN representations)**.** *We denote by $(\mathbb{M}_d)_{d\in\mathbb{N}} \subseteq \mathbf{N}$ the neural networks which satisfy that*

(i) *it holds for all $d \in \mathbb{N}$ that $\mathcal{I}(\mathbb{M}_d) = d$,*

(ii) *it holds for all $d \in \mathbb{N}$ that $\mathcal{O}(\mathbb{M}_d) = 1$,*

(iii) *it holds that $\mathbb{M}_1 = \mathbf{A}_{1,0} \in \mathbb{R}^{1\times 1} \times \mathbb{R}^1$,*

(iv) *it holds that*

$$\mathbb{M}_2 = \left(\left(\begin{pmatrix} 1 & -1 \\ 0 & 1 \\ 0 & -1 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}\right), \left(\begin{pmatrix} 1 & 1 & -1 \end{pmatrix}, (0)\right)\right) \in (\mathbb{R}^{3\times 2} \times \mathbb{R}^3) \times (\mathbb{R}^{1\times 3} \times \mathbb{R}^1),$$

(3.77)

(v) *it holds for all $d \in \{2,3,4,\ldots\}$ that $\mathbb{M}_{2d} = \mathbb{M}_d \bullet \left(\mathbf{P}_d(\mathbb{M}_2, \mathbb{M}_2, \ldots, \mathbb{M}_2)\right)$, and*

*(vi) it holds for all $d \in \{2, 3, 4, \ldots\}$ that $\mathbb{M}_{2d-1} = \mathbb{M}_d \bullet \big( \mathbf{P}_d(\mathbb{M}_2, \mathbb{M}_2, \ldots, \mathbb{M}_2, \mathfrak{I}_1) \big)$*

*(cf. Definitions 2.2.1, 2.2.5, 2.2.11, 2.2.18, and 2.2.20).*

**Definition 3.2.8** (Floor and ceiling of real numbers)**.** *We denote by $\lceil \cdot \rceil \colon \mathbb{R} \to \mathbb{Z}$ and $\lfloor \cdot \rfloor \colon \mathbb{R} \to \mathbb{Z}$ the functions which satisfy for all $x \in \mathbb{R}$ that $\lceil x \rceil = \min(\mathbb{Z} \cap [x, \infty))$ and $\lfloor x \rfloor = \max(\mathbb{Z} \cap (-\infty, x]).$*

**Proposition 3.2.9.** *Let $d \in \mathbb{N}$. Then*

*(i) it holds that $\mathcal{H}(\mathbb{M}_d) = \lceil \log_2(d) \rceil$,*

*(ii) it holds for all $i \in \mathbb{N}$ that $\mathbb{D}_i(\mathbb{M}_d) \leq 3 \lceil \frac{d}{2^i} \rceil$,*

*(iii) it holds that $\mathcal{R}_{\mathfrak{r}}(\mathbb{M}_d) \in C(\mathbb{R}^d, \mathbb{R})$, and*

*(iv) it holds for all $x = (x_1, x_2, \ldots, x_d) \in \mathbb{R}^d$ that $(\mathcal{R}_{\mathfrak{r}}(\mathbb{M}_d))(x) = \max\{x_1, x_2, \ldots, x_d\}$*

*(cf. Definitions 2.1.6, 2.2.1, 2.2.3, 3.2.7, and 3.2.8).*

*Proof of Proposition 3.2.9.* Throughout this proof assume w.l.o.g. that $d > 1$. Note that (3.77) ensures that $\mathcal{H}(\mathbb{M}_2) = 1$ (cf. Definitions 2.2.1 and 3.2.7). This and (2.111) demonstrate that for all $\mathfrak{d} \in \{2, 3, 4, \ldots\}$ it holds that

$$\mathcal{H}(\mathbf{P}_{\mathfrak{d}}(\mathbb{M}_2, \mathbb{M}_2, \ldots, \mathbb{M}_2)) = \mathcal{H}(\mathbf{P}_{\mathfrak{d}}(\mathbb{M}_2, \mathbb{M}_2, \ldots, \mathbb{M}_2, \mathfrak{I}_1)) = \mathcal{H}(\mathbb{M}_2) = 1 \qquad (3.78)$$

(cf. Definitions 2.2.11 and 2.2.18). Combining this with Proposition 2.2.7 establishes that for all $\mathfrak{d} \in \{3, 4, 5, \ldots\}$ it holds that

$$\mathcal{H}(\mathbb{M}_{\mathfrak{d}}) = \mathcal{H}(\mathbb{M}_{\lceil \mathfrak{d}/2 \rceil}) + 1 \qquad (3.79)$$

(cf. Definition 3.2.8). This assures that for all $\mathfrak{d} \in \{4, 6, 8, \ldots\}$ with $\mathcal{H}(\mathbb{M}_{\mathfrak{d}/2}) = \lceil \log_2(\mathfrak{d}/2) \rceil$ it holds that

$$\mathcal{H}(\mathbb{M}_{\mathfrak{d}}) = \lceil \log_2(\mathfrak{d}/2) \rceil + 1 = \lceil \log_2(\mathfrak{d}) - 1 \rceil + 1 = \lceil \log_2(\mathfrak{d}) \rceil. \qquad (3.80)$$

Moreover, note that (3.79) and the fact that for all $\mathfrak{d} \in \{3, 5, 7, \ldots\}$ it holds that $\lceil \log_2(\mathfrak{d} + 1) \rceil = \lceil \log_2(\mathfrak{d}) \rceil$ ensure that for all $\mathfrak{d} \in \{3, 5, 7, \ldots\}$ with $\mathcal{H}(\mathbb{M}_{\lceil \mathfrak{d}/2 \rceil}) = \lceil \log_2(\lceil \mathfrak{d}/2 \rceil) \rceil$ it holds that

$$\begin{aligned}
\mathcal{H}(\mathbb{M}_{\mathfrak{d}}) &= \lceil \log_2(\lceil \mathfrak{d}/2 \rceil) \rceil + 1 = \lceil \log_2((\mathfrak{d}+1)/2) \rceil + 1 \\
&= \lceil \log_2(\mathfrak{d} + 1) - 1 \rceil + 1 = \lceil \log_2(\mathfrak{d} + 1) \rceil = \lceil \log_2(\mathfrak{d}) \rceil.
\end{aligned} \qquad (3.81)$$

Combining this and (3.80) demonstrates that for all $\mathfrak{d} \in \{3, 4, 5, \ldots\}$ with $\forall\, k \in \{2, 3, \ldots, \mathfrak{d} - 1\} \colon \mathcal{H}(\mathbb{M}_k) = \lceil \log_2(k) \rceil$ it holds that $\mathcal{H}(\mathbb{M}_{\mathfrak{d}}) = \lceil \log_2(\mathfrak{d}) \rceil$. The fact that $\mathcal{H}(\mathbb{M}_2) = 1$ and induction hence establish item (i). Next note that the fact that $\mathcal{D}(\mathbb{M}_2) = (2, 3, 1)$ assure that for all $i \in \mathbb{N}$ it holds that

$$\mathbb{D}_i(\mathbb{M}_2) \leq 3 = 3 \lceil \tfrac{2}{2^i} \rceil. \qquad (3.82)$$

Moreover, observe that Proposition 2.2.7 and Lemma 2.2.12 imply that for all $\mathfrak{d} \in \{2, 3, 4, \ldots\}$, $i \in \mathbb{N}$ it holds that

$$\mathbb{D}_i(\mathbb{M}_{2\mathfrak{d}}) = \begin{cases} 3\mathfrak{d} & : i = 1 \\ \mathbb{D}_{i-1}(\mathbb{M}_{\mathfrak{d}}) & : i \geq 2 \end{cases} \qquad (3.83)$$

and

$$\mathbb{D}_i(\mathbb{M}_{2\mathfrak{d}-1}) = \begin{cases} 3\mathfrak{d} - 1 & : i = 1 \\ \mathbb{D}_{i-1}(\mathbb{M}_{\mathfrak{d}}) & : i \geq 2. \end{cases} \tag{3.84}$$

This assures that for all $\mathfrak{d} \in \{2, 4, 6, \ldots\}$ it holds that

$$\mathbb{D}_1(\mathbb{M}_{\mathfrak{d}}) = 3(\tfrac{\mathfrak{d}}{2}) \leq 3\lceil\tfrac{\mathfrak{d}}{2}\rceil. \tag{3.85}$$

Moreover, note that (3.84) ensures that for all $\mathfrak{d} \in \{3, 5, 7, \ldots\}$ it holds that

$$\mathbb{D}_1(\mathbb{M}_{\mathfrak{d}}) = 3\lceil\tfrac{\mathfrak{d}}{2}\rceil - 1 \leq 3\lceil\tfrac{\mathfrak{d}}{2}\rceil. \tag{3.86}$$

This and (3.85) show that for all $\mathfrak{d} \in \{2, 3, \ldots\}$ it holds that

$$\mathbb{D}_1(\mathbb{M}_{\mathfrak{d}}) \leq 3\lceil\tfrac{\mathfrak{d}}{2}\rceil. \tag{3.87}$$

In addition, observe that (3.83) demonstrates that for all $\mathfrak{d} \in \{4, 6, 8, \ldots\}$, $i \in \{2, 3, \ldots\}$ with $\mathbb{D}_{i-1}(\mathbb{M}_{\mathfrak{d}/2}) \leq 3\lceil(\mathfrak{d}/2)\frac{1}{2^{i-1}}\rceil$ it holds that

$$\mathbb{D}_i(\mathbb{M}_{\mathfrak{d}}) = \mathbb{D}_{i-1}(\mathbb{M}_{\mathfrak{d}/2}) \leq 3\lceil(\mathfrak{d}/2)\tfrac{1}{2^{i-1}}\rceil = 3\lceil\tfrac{\mathfrak{d}}{2^i}\rceil. \tag{3.88}$$

Furthermore, note that the fact that for all $\mathfrak{d} \in \{3, 5, 7, \ldots\}$, $i \in \mathbb{N}$ it holds that $\lceil\frac{\mathfrak{d}+1}{2^i}\rceil = \lceil\frac{\mathfrak{d}}{2^i}\rceil$ and (3.84) assure that for all $\mathfrak{d} \in \{3, 5, 7, \ldots\}$, $i \in \{2, 3, \ldots\}$ with $\mathbb{D}_{i-1}(\mathbb{M}_{\lceil\mathfrak{d}/2\rceil}) \leq 3\lceil\lceil\mathfrak{d}/2\rceil\frac{1}{2^{i-1}}\rceil$ it holds that

$$\mathbb{D}_i(\mathbb{M}_{\mathfrak{d}}) = \mathbb{D}_{i-1}(\mathbb{M}_{\lceil\mathfrak{d}/2\rceil}) \leq 3\lceil\lceil\mathfrak{d}/2\rceil\tfrac{1}{2^{i-1}}\rceil = 3\lceil\tfrac{\mathfrak{d}+1}{2^i}\rceil = 3\lceil\tfrac{\mathfrak{d}}{2^i}\rceil. \tag{3.89}$$

This and (3.88) ensure that for all $\mathfrak{d} \in \{3, 4, \ldots\}$, $i \in \{2, 3, \ldots\}$ with $\forall\, k \in \{2, 3, \ldots, \mathfrak{d} - 1\}, j \in \{1, 2, \ldots, i-1\}\colon \mathbb{D}_j(\mathbb{M}_k) \leq 3\lceil\frac{k}{2^j}\rceil$ it holds that

$$\mathbb{D}_i(\mathbb{M}_{\mathfrak{d}}) \leq 3\lceil\tfrac{\mathfrak{d}}{2^i}\rceil. \tag{3.90}$$

Combining this, (3.82), and (3.87) with induction establishes item (ii). Next observe that (3.77) ensures that for all $x = (x_1, x_2) \in \mathbb{R}^2$ it holds that

$$\begin{aligned}(\mathcal{R}_{\mathfrak{r}}(\mathbb{M}_2))(x) &= \max\{x_1 - x_2, 0\} + \max\{x_2, 0\} - \max\{-x_2, 0\} \\ &= \max\{x_1 - x_2, 0\} + x_2 = \max\{x_1, x_2\}\end{aligned} \tag{3.91}$$

(cf. Definitions 2.1.6 and 2.2.3). Proposition 2.2.13, Proposition 2.2.7, Lemma 2.2.19, and induction hence imply that for all $\mathfrak{d} \in \{2, 3, 4, \ldots\}$, $x = (x_1, x_2, \ldots, x_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ it holds that $\mathcal{R}_{\mathfrak{r}}(\mathbb{M}_{\mathfrak{d}}) \in C(\mathbb{R}^{\mathfrak{d}}, \mathbb{R})$ and $(\mathcal{R}_{\mathfrak{r}}(\mathbb{M}_{\mathfrak{d}}))(x) = \max\{x_1, x_2, \ldots, x_{\mathfrak{d}}\}$. This establishes items (iii)–(iv). The proof of Proposition 3.2.9 is thus complete. $\qquad\square$

**Lemma 3.2.10.** *Let $d \in \mathbb{N}$, $i \in \{1, 2, \ldots, \mathcal{L}(\mathbb{M}_d)\}$ (cf. Definitions 2.2.1 and 3.2.7). Then*

*(i) it holds that $\mathcal{B}_{i,\mathbb{M}_d} = 0 \in \mathbb{R}^{\mathbb{D}_i(\mathbb{M}_d)}$,*

*(ii) it holds that $\mathcal{W}_{i,\mathbb{M}_d} \in \{-1, 0, 1\}^{\mathbb{D}_i(\mathbb{M}_d) \times \mathbb{D}_{i-1}(\mathbb{M}_d)}$, and*

*(iii) it holds for all $x \in \mathbb{R}^d$ that $\|\mathcal{W}_{1,\mathbb{M}_d} x\|_\infty \leq 2\|x\|_\infty$*

*(cf. Definition 3.1.16).*

*Proof of Lemma 3.2.10.* Throughout this proof assume w.l.o.g. that $d > 2$ (cf. items (iii)–(iv) in Definition 3.2.7) and let $A_1 \in \mathbb{R}^{3 \times 2}$, $A_2 \in \mathbb{R}^{1 \times 3}$, $C_1 \in \mathbb{R}^{2 \times 1}$, $C_2 \in \mathbb{R}^{1 \times 2}$ satisfy

$$A_1 = \begin{pmatrix} 1 & -1 \\ 0 & 1 \\ 0 & -1 \end{pmatrix}, \qquad A_2 = \begin{pmatrix} 1 & 1 & -1 \end{pmatrix}, \qquad C_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \qquad \text{and} \qquad C_2 = \begin{pmatrix} 1 & -1 \end{pmatrix}. \tag{3.92}$$

Note that items (iv)–(vi) in Definition 3.2.7 assure that for all $\mathfrak{d} \in \{2, 3, 4, \ldots\}$ it holds that

$$\mathcal{W}_{1,\mathbb{M}_{2\mathfrak{d}-1}} = \underbrace{\begin{pmatrix} A_1 & 0 & \cdots & 0 & 0 \\ 0 & A_1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & A_1 & 0 \\ 0 & 0 & \cdots & 0 & C_1 \end{pmatrix}}_{\in \mathbb{R}^{(3\mathfrak{d}-1) \times (2\mathfrak{d}-1)}}, \qquad \mathcal{W}_{1,\mathbb{M}_{2\mathfrak{d}}} = \underbrace{\begin{pmatrix} A_1 & 0 & \cdots & 0 \\ 0 & A_1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & A_1 \end{pmatrix}}_{\in \mathbb{R}^{(3\mathfrak{d}) \times (2\mathfrak{d})}}, \tag{3.93}$$

$$\mathcal{B}_{1,\mathbb{M}_{2\mathfrak{d}-1}} = 0 \in \mathbb{R}^{3\mathfrak{d}-1}, \qquad \text{and} \qquad \mathcal{B}_{1,\mathbb{M}_{2\mathfrak{d}}} = 0 \in \mathbb{R}^{3\mathfrak{d}}.$$

This and (3.92) proves item (iii). Furthermore, note that (3.93) and item (iv) in Definition 3.2.7 imply that for all $\mathfrak{d} \in \{2, 3, 4, \ldots\}$ it holds that $\mathcal{B}_{1,\mathbb{M}_{\mathfrak{d}}} = 0$. Items (iv)–(vi) in Definition 3.2.7 hence ensures that for all $\mathfrak{d} \in \{2, 3, 4, \ldots\}$ it holds that

$$\mathcal{W}_{2,\mathbb{M}_{2\mathfrak{d}-1}} = \mathcal{W}_{1,\mathbb{M}_{\mathfrak{d}}} \underbrace{\begin{pmatrix} A_2 & 0 & \cdots & 0 & 0 \\ 0 & A_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & A_2 & 0 \\ 0 & 0 & \cdots & 0 & C_2 \end{pmatrix}}_{\in \mathbb{R}^{\mathfrak{d} \times (3\mathfrak{d}-1)}}, \qquad \mathcal{W}_{2,\mathbb{M}_{2\mathfrak{d}}} = \mathcal{W}_{1,\mathbb{M}_{\mathfrak{d}}} \underbrace{\begin{pmatrix} A_2 & 0 & \cdots & 0 \\ 0 & A_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & A_2 \end{pmatrix}}_{\in \mathbb{R}^{\mathfrak{d} \times (3\mathfrak{d})}}, \tag{3.94}$$

$$\mathcal{B}_{2,\mathbb{M}_{2\mathfrak{d}-1}} = \mathcal{B}_{1,\mathbb{M}_{\mathfrak{d}}} = 0, \qquad \text{and} \qquad \mathcal{B}_{2,\mathbb{M}_{2\mathfrak{d}}} = \mathcal{B}_{1,\mathbb{M}_{\mathfrak{d}}} = 0.$$

Combining this and item (iv) in Definition 3.2.7 shows that for all $\mathfrak{d} \in \{2, 3, 4, \ldots\}$ it holds that $\mathcal{B}_{2,\mathbb{M}_{\mathfrak{d}}} = 0$. Moreover, note that (2.59) demonstrates that for all $\mathfrak{d} \in \{2, 3, 4, \ldots, \}$, $i \in \{3, 4, \ldots, \mathcal{L}(\mathbb{M}_{\mathfrak{d}}) + 1\}$ it holds that

$$\mathcal{W}_{i,\mathbb{M}_{2\mathfrak{d}-1}} = \mathcal{W}_{i,\mathbb{M}_{2\mathfrak{d}}} = \mathcal{W}_{i-1,\mathbb{M}_{\mathfrak{d}}} \qquad \text{and} \qquad \mathcal{B}_{i,\mathbb{M}_{2\mathfrak{d}-1}} = \mathcal{B}_{i,\mathbb{M}_{2\mathfrak{d}}} = \mathcal{B}_{i-1,\mathbb{M}_{\mathfrak{d}}}. \tag{3.95}$$

This, (3.92), (3.93), (3.94), the fact that for all $\mathfrak{d} \in \{2, 3, 4, \ldots\}$ it holds that $\mathcal{B}_{2,\mathbb{M}_{\mathfrak{d}}} = 0$, and induction establish items (i)–(ii). The proof of Lemma 3.2.10 is thus complete. $\square$

### 3.2.2.3 ANN representations for maximum convolutions

**Lemma 3.2.11.** *Let* $d, K \in \mathbb{N}$, $L \in [0, \infty)$, $\mathfrak{x}_1, \mathfrak{x}_2, \ldots, \mathfrak{x}_K \in \mathbb{R}^d$, $\mathfrak{y} = (\mathfrak{y}_1, \mathfrak{y}_2, \ldots, \mathfrak{y}_K) \in \mathbb{R}^K$, $\Phi \in \mathbf{N}$ *satisfy*

$$\Phi = \mathbb{M}_K \bullet \mathbf{A}_{-L\,\mathrm{I}_K,\mathfrak{y}} \bullet \mathbf{P}_K\big(\mathbb{L}_d \bullet \mathbf{A}_{I_d,-\mathfrak{x}_1}, \mathbb{L}_d \bullet \mathbf{A}_{I_d,-\mathfrak{x}_2}, \ldots, \mathbb{L}_d \bullet \mathbf{A}_{I_d,-\mathfrak{x}_K}\big) \bullet \mathbb{T}_{K,d} \tag{3.96}$$

*(cf. Definitions 2.2.1, 2.2.5, 2.2.9, 2.2.11, 2.2.20, 2.2.30, 3.2.3, and 3.2.7). Then*

(i) it holds that $\mathcal{I}(\Phi) = d$,

(ii) it holds that $\mathcal{O}(\Phi) = 1$,

(iii) it holds that $\mathcal{H}(\Phi) = \lceil \log_2(K) \rceil + 1$,

(iv) it holds that $\mathbb{D}_1(\Phi) = 2dK$,

(v) it holds for all $i \in \{2, 3, \ldots\}$ that $\mathbb{D}_i(\Phi) \leq 3\lceil \frac{K}{2^{i-1}} \rceil$,

(vi) it holds that $\|\mathcal{T}(\Phi)\|_\infty \leq \max\{1, L, \max_{k \in \{1,2,\ldots,K\}}\|\mathfrak{x}_k\|_\infty, 2\|\mathfrak{y}\|_\infty\}$, and

(vii) it holds for all $x \in \mathbb{R}^d$ that $(\mathcal{R}_\mathfrak{r}(\Phi))(x) = \max_{k \in \{1,2,\ldots,K\}}(\mathfrak{y}_k - L\|x - \mathfrak{x}_k\|_1)$

*(cf. Definitions 2.1.6, 2.2.3, 2.2.36, 3.1.16, and 3.2.8).*

*Proof of Lemma 3.2.11.* Throughout this proof let $\Psi_k \in \mathbf{N}$, $k \in \{1, 2, \ldots, K\}$, satisfy for all $k \in \{1, 2, \ldots, K\}$ that $\Psi_k = \mathbb{L}_d \bullet \mathbf{A}_{I_d, -\mathfrak{x}_k}$, let $\Xi \in \mathbf{N}$ satisfy

$$\Xi = \mathbf{A}_{-L\,\mathrm{I}_K, \mathfrak{y}} \bullet \mathbf{P}_K\big(\Psi_1, \Psi_2, \ldots, \Psi_K\big) \bullet \mathbb{T}_{K,d}, \tag{3.97}$$

and let $\|\!|\cdot|\!\|\colon \bigcup_{m,n \in \mathbb{N}} \mathbb{R}^{m \times n} \to [0, \infty)$ satisfy for all $m, n \in \mathbb{N}$, $M = (M_{i,j})_{i \in \{1,\ldots,m\}, j \in \{1,\ldots,n\}} \in \mathbb{R}^{m \times n}$ that $\|\!|M|\!\| = \max_{i \in \{1,\ldots,m\}, j \in \{1,\ldots,n\}}|M_{i,j}|$. Observe that (3.96) and Proposition 2.2.7 ensure that $\mathcal{O}(\Phi) = \mathcal{O}(\mathbb{M}_K) = 1$ and $\mathcal{I}(\Phi) = \mathcal{I}(\mathbb{T}_{K,d}) = d$. This proves items (i)–(ii). Moreover, observe that the fact that for all $m, n \in \mathbb{N}$, $\mathfrak{W} \in \mathbb{R}^{m \times n}$, $\mathfrak{B} \in \mathbb{R}^m$ it holds that $\mathcal{H}(\mathbf{A}_{\mathfrak{W},\mathfrak{B}}) = 0 = \mathcal{H}(\mathbb{T}_{K,d})$, the fact that $\mathcal{H}(\mathbb{L}_d) = 1$, and Proposition 2.2.7 assure that

$$\mathcal{H}(\Xi) = \mathcal{H}(\mathbf{A}_{-L\,\mathrm{I}_K, \mathfrak{y}}) + \mathcal{H}(\mathbf{P}_K(\Psi_1, \Psi_2, \ldots, \Psi_K)) + \mathcal{H}(\mathbb{T}_{K,d}) = \mathcal{H}(\Psi_1) = \mathcal{H}(\mathbb{L}_d) = 1. \tag{3.98}$$

Proposition 2.2.7 and Proposition 3.2.9 hence ensure that

$$\mathcal{H}(\Phi) = \mathcal{H}(\mathbb{M}_K \bullet \Xi) = \mathcal{H}(\mathbb{M}_K) + \mathcal{H}(\Xi) = \lceil \log_2(K) \rceil + 1 \tag{3.99}$$

(cf. Definition 3.2.8). This establishes item (iii). Next observe that the fact that $\mathcal{H}(\Xi) = 1$, Proposition 2.2.7, and Proposition 3.2.9 assure that for all $i \in \{2, 3, \ldots\}$ it holds that

$$\mathbb{D}_i(\Phi) = \mathbb{D}_{i-1}(\mathbb{M}_K) \leq 3\lceil \tfrac{K}{2^{i-1}} \rceil. \tag{3.100}$$

This proves item (v). Furthermore, note that Proposition 2.2.7, Proposition 2.2.14, and Proposition 3.2.4 assure that

$$\mathbb{D}_1(\Phi) = \mathbb{D}_1(\Xi) = \mathbb{D}_1(\mathbf{P}_K(\Psi_1, \Psi_2, \ldots, \Psi_K)) = \sum_{i=1}^{K} \mathbb{D}_1(\Psi_i) = \sum_{i=1}^{K} \mathbb{D}_1(\mathbb{L}_d) = 2dK. \tag{3.101}$$

This establishes item (iv). Next observe that (2.59) and Lemma 3.2.10 imply that

$$\begin{aligned}\Phi = \big(&(\mathcal{W}_{1,\Xi}, \mathcal{B}_{1,\Xi}), (\mathcal{W}_{1,\mathbb{M}_K}\mathcal{W}_{2,\Xi}, \mathcal{W}_{1,\mathbb{M}_K}\mathcal{B}_{2,\Xi}), \\ &(\mathcal{W}_{2,\mathbb{M}_K}, 0), \ldots, (\mathcal{W}_{\mathcal{L}(\mathbb{M}_K),\mathbb{M}_K}, 0)\big).\end{aligned} \tag{3.102}$$

Moreover, note that the fact that for all $k \in \{1, 2, \ldots, K\}$ it holds that $\mathcal{W}_{1,\Psi_k} = \mathcal{W}_{1,\mathbf{A}_{I_d,-\mathfrak{r}_k}} \mathcal{W}_{1,\mathbb{L}_d} = \mathcal{W}_{1,\mathbb{L}_d}$ assures that

$$
\mathcal{W}_{1,\Xi} = \mathcal{W}_{1,\mathbf{P}_K(\Psi_1,\Psi_2,\ldots,\Psi_K)} \mathcal{W}_{1,\mathbb{T}_{K,d}} = \begin{pmatrix} \mathcal{W}_{1,\Psi_1} & 0 & \cdots & 0 \\ 0 & \mathcal{W}_{1,\Psi_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathcal{W}_{1,\Psi_K} \end{pmatrix} \begin{pmatrix} \mathrm{I}_d \\ \mathrm{I}_d \\ \vdots \\ \mathrm{I}_d \end{pmatrix}
$$

$$
= \begin{pmatrix} \mathcal{W}_{1,\Psi_1} \\ \mathcal{W}_{1,\Psi_2} \\ \vdots \\ \mathcal{W}_{1,\Psi_K} \end{pmatrix} = \begin{pmatrix} \mathcal{W}_{1,\mathbb{L}_d} \\ \mathcal{W}_{1,\mathbb{L}_d} \\ \vdots \\ \mathcal{W}_{1,\mathbb{L}_d} \end{pmatrix}. \tag{3.103}
$$

Lemma 3.2.5 hence demonstrates that $\|\mathcal{W}_{1,\Xi}\| = 1$. In addition, note that (2.59) implies that

$$
\mathcal{B}_{1,\Xi} = \mathcal{W}_{1,\mathbf{P}_K(\Psi_1,\Psi_2,\ldots,\Psi_K)} \mathcal{B}_{1,\mathbb{T}_{K,d}} + \mathcal{B}_{1,\mathbf{P}_K(\Psi_1,\Psi_2,\ldots,\Psi_K)} = \mathcal{B}_{1,\mathbf{P}_K(\Psi_1,\Psi_2,\ldots,\Psi_K)} = \begin{pmatrix} \mathcal{B}_{1,\Psi_1} \\ \mathcal{B}_{1,\Psi_2} \\ \vdots \\ \mathcal{B}_{1,\Psi_K} \end{pmatrix}. \tag{3.104}
$$

Furthermore, observe that Lemma 3.2.5 implies that for all $k \in \{1, 2, \ldots, K\}$ it holds that

$$
\mathcal{B}_{1,\Psi_k} = \mathcal{W}_{1,\mathbb{L}_d} \mathcal{B}_{1,\mathbf{A}_{I_d,-\mathfrak{r}_k}} + \mathcal{B}_{1,\mathbb{L}_d} = -\mathcal{W}_{1,\mathbb{L}_d} \mathfrak{r}_k. \tag{3.105}
$$

This, (3.104), and Lemma 3.2.5 show that

$$
\|\mathcal{B}_{1,\Xi}\|_\infty = \max_{k \in \{1,2,\ldots,K\}} \|\mathcal{B}_{1,\Psi_k}\|_\infty = \max_{k \in \{1,2,\ldots,K\}} \|\mathcal{W}_{1,\mathbb{L}_d} \mathfrak{r}_k\|_\infty = \max_{k \in \{1,2,\ldots,K\}} \|\mathfrak{r}_k\|_\infty \tag{3.106}
$$

(cf. Definition 3.1.16). Combining this, (3.102), Lemma 3.2.10, and the fact that $\|\mathcal{W}_{1,\Xi}\| = 1$ shows that

$$
\|\mathcal{T}(\Phi)\|_\infty = \max\{\|\mathcal{W}_{1,\Xi}\|, \|\mathcal{B}_{1,\Xi}\|_\infty, \|\mathcal{W}_{1,\mathbb{M}_K}\mathcal{W}_{2,\Xi}\|, \|\mathcal{W}_{1,\mathbb{M}_K}\mathcal{B}_{2,\Xi}\|_\infty, 1\}
$$
$$
\leq \max\{1, \max_{k \in \{1,2,\ldots,K\}}\|\mathfrak{r}_k\|_\infty, \|\mathcal{W}_{1,\mathbb{M}_K}\mathcal{W}_{2,\Xi}\|, \|\mathcal{W}_{1,\mathbb{M}_K}\mathcal{B}_{2,\Xi}\|_\infty\} \tag{3.107}
$$

(cf. Definition 2.2.36). Next note that Lemma 3.2.5 ensures that for all $k \in \{1, 2, \ldots, K\}$ it holds that $\mathcal{B}_{2,\Psi_k} = \mathcal{B}_{2,\mathbb{L}_d} = 0$. Hence, we obtain that $\mathcal{B}_{2,\mathbf{P}_K(\Psi_1,\Psi_2,\ldots,\Psi_K)} = 0$. This implies that

$$
\mathcal{B}_{2,\Xi} = \mathcal{W}_{1,\mathbf{A}_{-L\,\mathrm{I}_K,\mathfrak{y}}} \mathcal{B}_{2,\mathbf{P}_K(\Psi_1,\Psi_2,\ldots,\Psi_K)} + \mathcal{B}_{1,\mathbf{A}_{-L\,\mathrm{I}_K,\mathfrak{y}}} = \mathcal{B}_{1,\mathbf{A}_{-L\,\mathrm{I}_K,\mathfrak{y}}} = \mathfrak{y}. \tag{3.108}
$$

In addition, observe that the fact that for all $k \in \{1, 2, \ldots, K\}$ it holds that $\mathcal{W}_{2,\Psi_k} = \mathcal{W}_{2,\mathbb{L}_d}$ assures that

$$
\mathcal{W}_{2,\Xi} = \mathcal{W}_{1,\mathbf{A}_{-L\,\mathrm{I}_K,\mathfrak{y}}} \mathcal{W}_{2,\mathbf{P}_K(\Psi_1,\Psi_2,\ldots,\Psi_K)} = -L\mathcal{W}_{2,\mathbf{P}_K(\Psi_1,\Psi_2,\ldots,\Psi_K)}
$$

$$
= -L \begin{pmatrix} \mathcal{W}_{2,\Psi_1} & 0 & \cdots & 0 \\ 0 & \mathcal{W}_{2,\Psi_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathcal{W}_{2,\Psi_K} \end{pmatrix} = \begin{pmatrix} -L\mathcal{W}_{2,\mathbb{L}_d} & 0 & \cdots & 0 \\ 0 & -L\mathcal{W}_{2,\mathbb{L}_d} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & -L\mathcal{W}_{2,\mathbb{L}_d} \end{pmatrix}. \tag{3.109}
$$

Item (v) in Lemma 3.2.5 and Lemma 3.2.10 hence imply that

$$\||\mathcal{W}_{1,\mathbb{M}_K}\mathcal{W}_{2,\Xi}\|| \le L\||\mathcal{W}_{1,\mathbb{M}_K}\|| \le L. \tag{3.110}$$

Moreover, observe that (3.108), (3.109), and Lemma 3.2.10 assure that

$$\|\mathcal{W}_{1,\mathbb{M}_K}\mathcal{B}_{2,\Xi}\|_\infty \le 2\|\mathcal{B}_{2,\Xi}\|_\infty = 2\|\mathfrak{y}\|_\infty. \tag{3.111}$$

Combining this with (3.107) and (3.110) establishes item (vi). Next observe that Proposition 3.2.4 and Lemma 2.2.22 show that for all $x \in \mathbb{R}^d$, $k \in \{1, 2, \dots, K\}$ it holds that

$$(\mathcal{R}_{\mathfrak{r}}(\Psi_k))(x) = \big(\mathcal{R}_{\mathfrak{r}}(\mathbb{L}_d) \circ \mathcal{R}_{\mathfrak{r}}(\mathbf{A}_{I_d,-\mathfrak{x}_k})\big)(x) = \|x - \mathfrak{x}_k\|_1. \tag{3.112}$$

This, Proposition 2.2.13, and Proposition 2.2.7 imply that for all $x \in \mathbb{R}^d$ it holds that

$$\big(\mathcal{R}_{\mathfrak{r}}(\mathbf{P}_K(\Psi_1, \Psi_2, \dots, \Psi_K) \bullet \mathbb{T}_{K,d})\big)(x) = \big(\|x - \mathfrak{x}_1\|_1, \|x - \mathfrak{x}_2\|_1, \dots, \|x - \mathfrak{x}_K\|_1\big). \tag{3.113}$$

(cf. Definitions 2.1.6 and 2.2.3). Combining this and Lemma 2.2.22 establishes that for all $x \in \mathbb{R}^d$ it holds that

$$\begin{aligned}
(\mathcal{R}_{\mathfrak{r}}(\Xi))(x) &= \big(\mathcal{R}_{\mathfrak{r}}(\mathbf{A}_{-L\,I_K,\mathfrak{y}}) \circ \mathcal{R}_{\mathfrak{r}}(\mathbf{P}_K(\Psi_1, \Psi_2, \dots, \Psi_K) \bullet \mathbb{T}_{K,d})\big)(x) \\
&= \big(\mathfrak{y}_1 - L\|x - \mathfrak{x}_1\|_1, \mathfrak{y}_2 - L\|x - \mathfrak{x}_2\|_1, \dots, \mathfrak{y}_K - L\|x - \mathfrak{x}_K\|_1\big).
\end{aligned} \tag{3.114}$$

Proposition 2.2.7 and Proposition 3.2.9 hence demonstrate that for all $x \in \mathbb{R}^d$ it holds that

$$\begin{aligned}
(\mathcal{R}_{\mathfrak{r}}(\Phi))(x) &= \big(\mathcal{R}_{\mathfrak{r}}(\mathbb{M}_K) \circ \mathcal{R}_{\mathfrak{r}}(\Xi)\big)(x) \\
&= (\mathcal{R}_{\mathfrak{r}}(\mathbb{M}_K))\big(\mathfrak{y}_1 - L\|x - \mathfrak{x}_1\|_1, \mathfrak{y}_2 - L\|x - \mathfrak{x}_2\|_1, \dots, \mathfrak{y}_K - L\|x - \mathfrak{x}_K\|_1\big) \\
&= \max_{k \in \{1,2,\dots,K\}}(\mathfrak{y}_k - L\|x - \mathfrak{x}_k\|_1).
\end{aligned} \tag{3.115}$$

This establishes item (vii). The proof of Lemma 3.2.11 is thus complete. $\square$

### 3.2.3   Explicit approximations through ANNs

**Proposition 3.2.12.** *Let $d, K \in \mathbb{N}$, $L \in [0, \infty)$, let $E \subseteq \mathbb{R}^d$ be a set, let $\mathfrak{x}_1, \mathfrak{x}_2, \dots, \mathfrak{x}_K \in E$, let $f\colon E \to \mathbb{R}$ satisfy for all $x, y \in E$ that $|f(x) - f(y)| \le L\|x - y\|_1$, and let $\mathfrak{y} \in \mathbb{R}^K$, $\Phi \in \mathbf{N}$ satisfy $\mathfrak{y} = (f(\mathfrak{x}_1), f(\mathfrak{x}_2), \dots, f(\mathfrak{x}_K))$ and*

$$\Phi = \mathbb{M}_K \bullet \mathbf{A}_{-L\,I_K,\mathfrak{y}} \bullet \mathbf{P}_K\big(\mathbb{L}_d \bullet \mathbf{A}_{I_d,-\mathfrak{x}_1}, \mathbb{L}_d \bullet \mathbf{A}_{I_d,-\mathfrak{x}_2}, \dots, \mathbb{L}_d \bullet \mathbf{A}_{I_d,-\mathfrak{x}_K}\big) \bullet \mathbb{T}_{K,d} \tag{3.116}$$

*(cf. Definitions 2.2.1, 2.2.5, 2.2.9, 2.2.11, 2.2.20, 2.2.30, 3.1.16, 3.2.3, and 3.2.7). Then*

$$\sup_{x \in E} |(\mathcal{R}_{\mathfrak{r}}(\Phi))(x) - f(x)| \le 2L\big[\sup_{x \in E}\big(\min_{k \in \{1,2,\dots,K\}}\|x - \mathfrak{x}_k\|_1\big)\big] \tag{3.117}$$

*(cf. Definitions 2.1.6 and 2.2.3).*

*Proof of Proposition 3.2.12.* Throughout this proof let $F\colon \mathbb{R}^d \to \mathbb{R}$ satisfy for all $x \in \mathbb{R}^d$ that

$$F(x) = \max_{k \in \{1,2,\dots,K\}}(f(\mathfrak{x}_k) - L\|x - \mathfrak{x}_k\|_1). \tag{3.118}$$

Observe that Corollary 3.2.2, (3.118), and the assumption that for all $x, y \in E$ it holds that $|f(x) - f(y)| \leq L\|x - y\|_1$ assure that

$$\sup_{x \in E}|F(x) - f(x)| \leq 2L\big[\sup_{x \in E}\big(\min_{k \in \{1,2,\ldots,K\}}\|x - \mathfrak{x}_k\|_1\big)\big]. \tag{3.119}$$

Moreover, note that Lemma 3.2.11 ensures that for all $x \in E$ it holds that $F(x) = (\mathcal{R}_{\mathfrak{r}}(\Phi))(x)$. Combining this and (3.119) establishes (3.117). The proof of Proposition 3.2.12 is thus complete. $\qquad\square$

**Exercise 3.2.2.** *Prove or disprove the following statement: There exists $\Phi \in \mathbf{N}$ such that $\mathcal{I}(\Phi) = 2$, $\mathcal{O}(\Phi) = 1$, $\mathcal{P}(\Phi) < 20$, and*

$$\sup_{v=(x,y)\in[0,2]^2}\big|x^2 + y^2 - 2x - 2y + 2 - (\mathcal{R}_{\mathfrak{r}}(\Phi))(v)\big| \leq \tfrac{3}{8}. \tag{3.120}$$

**Exercise 3.2.3.** *Prove or disprove the following statement: For all $n \in \{3, 5, 7, \ldots\}$ it holds that $\lceil \log_2(n+1) \rceil = \lceil \log_2(n) \rceil$.*

### 3.2.4 Analysis of the approximation error

#### 3.2.4.1 Covering number estimates

**Definition 3.2.13** (Covering number). *Let $(E, \delta)$ be a metric space and let $r \in [0, \infty]$. Then we denote by $\mathcal{C}^{(E,\delta),r} \in \mathbb{N}_0 \cup \{\infty\}$ (we denote by $\mathcal{C}^{E,r} \in \mathbb{N}_0 \cup \{\infty\}$) the extended real number given by*

$$\mathcal{C}^{(E,\delta),r} = \min\left(\left\{n \in \mathbb{N}_0 \colon \left[\exists\, A \subseteq E \colon \left(\begin{array}{c}(|A| \leq n) \wedge (\forall\, x \in E\colon \\ \exists\, a \in A\colon \delta(a,x) \leq r)\end{array}\right)\right]\right\} \cup \{\infty\}\right). \tag{3.121}$$

**Exercise 3.2.4.** *Prove or disprove the following statement: For every metric space $(X, d)$, every $Y \subseteq X$, and every $r \in [0, \infty]$ it holds that $\mathcal{C}^{(Y,d|_{Y \times Y}),r} \leq \mathcal{C}^{(X,d),r}$.*

**Exercise 3.2.5.** *Prove or disprove the following statement: For every metric space $(E, \delta)$ it holds that $\mathcal{C}^{(E,\delta),\infty} = 1$.*

**Exercise 3.2.6.** *Prove or disprove the following statement: For every metric space $(E, \delta)$ and every $r \in [0, \infty)$ with $\mathcal{C}^{(E,\delta),r} < \infty$ it holds that $E$ is bounded.* (Note: A metric space $(E, \delta)$ is bounded if and only if there exists $r \in [0, \infty)$ such that it holds for all $x, y \in E$ that $\delta(x, y) \leq r$.)

**Exercise 3.2.7.** *Prove or disprove the following statement: For every bounded metric space $(E, \delta)$ and every $r \in [0, \infty]$ it holds that $\mathcal{C}^{(E,\delta),r} < \infty$.*

**Lemma 3.2.14.** *Let $d \in \mathbb{N}$, $a \in \mathbb{R}$, $b \in (a, \infty)$, $r \in (0, \infty)$ and for every $p \in [1, \infty]$ let $\delta_p \colon ([a,b]^d) \times ([a,b]^d) \to [0, \infty)$ satisfy for all $x, y \in [a,b]^d$ that $\delta_p(x,y) = \|x-y\|_p$ (cf. Definition 3.1.16). Then*

*(i) it holds for all $p \in [1, \infty)$ that*

$$\mathcal{C}^{([a,b]^d,\delta_p),r} \leq \left(\left\lceil \tfrac{d^{1/p}(b-a)}{2r} \right\rceil\right)^d \leq \begin{cases} 1 & : r \geq {}^{d(b-a)}/_2 \\ \left(\tfrac{d(b-a)}{r}\right)^d & : r < {}^{d(b-a)}/_2 \end{cases} \tag{3.122}$$

*and*

*(ii)* it holds that

$$\mathcal{C}^{([a,b]^d,\delta_\infty),r} \leq \left(\left\lceil \tfrac{b-a}{2r} \right\rceil\right)^d \leq \begin{cases} 1 & : r \geq {(b-a)}/{2} \\ \left(\tfrac{b-a}{r}\right)^d & : r < {(b-a)}/{2} \end{cases} \tag{3.123}$$

*(cf. Definitions 3.2.8 and 3.2.13).*

*Proof of Lemma 3.2.14.* Throughout this proof let $(\mathfrak{N}_p)_{p\in[1,\infty]} \subseteq \mathbb{N}$ satisfy for all $p \in [1,\infty)$ that

$$\mathfrak{N}_p = \left\lceil \tfrac{d^{1/p}(b-a)}{2r} \right\rceil \qquad \text{and} \qquad \mathfrak{N}_\infty = \left\lceil \tfrac{b-a}{2r} \right\rceil, \tag{3.124}$$

for every $N \in \mathbb{N}$, $i \in \{1, 2, \ldots, N\}$ let $g_{N,i} \in [a,b]$ be given by $g_{N,i} = a + {(i-1/2)(b-a)}/{N}$, and for every $p \in [1,\infty]$ let $A_p \subseteq [a,b]^d$ be given by $A_p = \{g_{\mathfrak{N}_p,1}, g_{\mathfrak{N}_p,2}, \ldots, g_{\mathfrak{N}_p,\mathfrak{N}_p}\}^d$ (cf. Definition 3.2.8). Observe that it holds for all $N \in \mathbb{N}$, $i \in \{1, 2, \ldots, N\}$, $x \in [a + {(i-1)(b-a)}/{N}, g_{N,i}]$ that

$$|x - g_{N,i}| = a + \tfrac{(i-1/2)(b-a)}{N} - x \leq a + \tfrac{(i-1/2)(b-a)}{N} - \left(a + \tfrac{(i-1)(b-a)}{N}\right) = \tfrac{b-a}{2N}. \tag{3.125}$$

In addition, note that it holds for all $N \in \mathbb{N}$, $i \in \{1, 2, \ldots, N\}$, $x \in [g_{N,i}, a + {i(b-a)}/{N}]$ that

$$|x - g_{N,i}| = x - \left(a + \tfrac{(i-1/2)(b-a)}{N}\right) \leq a + \tfrac{i(b-a)}{N} - \left(a + \tfrac{(i-1/2)(b-a)}{N}\right) = \tfrac{b-a}{2N}. \tag{3.126}$$

Combining (3.125) and (3.126) implies for all $N \in \mathbb{N}$, $i \in \{1, 2, \ldots, N\}$, $x \in [a + {(i-1)(b-a)}/{N}, a + {i(b-a)}/{N}]$ that $|x - g_{N,i}| \leq {(b-a)}/{(2N)}$. This proves that for every $N \in \mathbb{N}$, $x \in [a,b]$ there exists $y \in \{g_{N,1}, g_{N,2}, \ldots, g_{N,N}\}$ such that

$$|x - y| \leq \tfrac{b-a}{2N}. \tag{3.127}$$

This establishes that for every $p \in [1,\infty)$, $x = (x_1, x_2, \ldots, x_d) \in [a,b]^d$ there exists $y = (y_1, y_2, \ldots, y_d) \in A_p$ such that

$$\delta_p(x,y) = \|x - y\|_p = \left(\sum_{i=1}^d |x_i - y_i|^p\right)^{1/p} \leq \left(\sum_{i=1}^d \tfrac{(b-a)^p}{(2\mathfrak{N}_p)^p}\right)^{1/p} = \tfrac{d^{1/p}(b-a)}{2\mathfrak{N}_p} \leq \tfrac{d^{1/p}(b-a)2r}{2d^{1/p}(b-a)} = r. \tag{3.128}$$

Furthermore, (3.127) shows that for every $x = (x_1, x_2, \ldots, x_d) \in [a,b]^d$ there exists $y = (y_1, y_2, \ldots, y_d) \in A_\infty$ such that

$$\delta_\infty(x,y) = \|x - y\|_\infty = \max_{i\in\{1,2,\ldots,d\}} |x_i - y_i| \leq \tfrac{b-a}{2\mathfrak{N}_\infty} \leq \tfrac{(b-a)2r}{2(b-a)} = r. \tag{3.129}$$

Note that (3.128), (3.124), and the fact that $\forall\, x \in [0,\infty): \lceil x \rceil \leq \mathbb{1}_{(0,1]}(x) + 2x\mathbb{1}_{(1,\infty)}(x) = \mathbb{1}_{(0,r]}(rx) + 2x\mathbb{1}_{(r,\infty)}(rx)$ yield that for all $p \in [1,\infty)$ it holds that

$$\begin{aligned} \mathcal{C}^{([a,b]^d,\delta_p),r} &\leq |A_p| = (\mathfrak{N}_p)^d = \left(\left\lceil \tfrac{d^{1/p}(b-a)}{2r} \right\rceil\right)^d \leq \left(\left\lceil \tfrac{d(b-a)}{2r} \right\rceil\right)^d \\ &\leq \left(\mathbb{1}_{(0,r]}\left(\tfrac{d(b-a)}{2}\right) + \tfrac{2d(b-a)}{2r}\mathbb{1}_{(r,\infty)}\left(\tfrac{d(b-a)}{2}\right)\right)^d \\ &= \mathbb{1}_{(0,r]}\left(\tfrac{d(b-a)}{2}\right) + \left(\tfrac{d(b-a)}{r}\right)^d \mathbb{1}_{(r,\infty)}\left(\tfrac{d(b-a)}{2}\right) \end{aligned} \tag{3.130}$$

(cf. Definition 3.2.13). This proves item (i). In addition, (3.129), (3.124), and the fact that $\forall\, x \in [0,\infty): \lceil x \rceil \leq \mathbb{1}_{(0,r]}(rx) + 2x\mathbb{1}_{(r,\infty)}(rx)$ demonstrate that

$$\mathcal{C}^{([a,b]^d,\delta_\infty),r} \leq |A_\infty| = (\mathfrak{N}_\infty)^d = \left(\left\lceil \tfrac{b-a}{2r} \right\rceil\right)^d \leq \mathbb{1}_{(0,r]}\left(\tfrac{b-a}{2}\right) + \left(\tfrac{b-a}{r}\right)^d \mathbb{1}_{(r,\infty)}\left(\tfrac{b-a}{2}\right). \tag{3.131}$$

This implies item (ii). and thus completes the proof of Lemma 3.2.14. $\square$

### 3.2.4.2 Convergence rates for the approximation error

**Lemma 3.2.15.** *Let $d \in \mathbb{N}$, $L, a \in \mathbb{R}$, $b \in (a, \infty)$, let $f \colon [a,b]^d \to \mathbb{R}$ satisfy for all $x, y \in [a,b]^d$ that $|f(x) - f(y)| \le L\|x - y\|_1$, and let $\mathbf{F} = \mathbf{A}_{0, f((a+b)/2, (a+b)/2, \dots, (a+b)/2)} \in \mathbb{R}^{1 \times d} \times \mathbb{R}^1$ (cf. Definitions 2.2.20 and 3.1.16). Then*

*(i) it holds that $\mathcal{I}(\mathbf{F}) = d$,*

*(ii) it holds that $\mathcal{O}(\mathbf{F}) = 1$,*

*(iii) it holds that $\mathcal{H}(\mathbf{F}) = 0$,*

*(iv) it holds that $\mathcal{P}(\mathbf{F}) = d + 1$,*

*(v) it holds that $\|\mathcal{T}(\mathbf{F})\|_\infty \le \sup_{x \in [a,b]^d} |f(x)|$, and*

*(vi) it holds that $\sup_{x \in [a,b]^d} |(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - f(x)| \le \frac{dL(b-a)}{2}$*

*(cf. Definitions 2.1.6, 2.2.1, 2.2.3, and 2.2.36).*

*Proof of Lemma 3.2.15.* Note that the assumption that for all $x, y \in [a,b]^d$ it holds that $|f(x) - f(y)| \le L\|x - y\|_1$ assures that $L \ge 0$. Next observe that Lemma 2.2.21 assures that for all $x \in \mathbb{R}^d$ it holds that

$$(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) = f\big((a+b)/2, (a+b)/2, \dots, (a+b)/2\big). \tag{3.132}$$

The fact that for all $x \in [a,b]$ it holds that $|x - (a+b)/2| \le (a+b)/2$ and the assumption that for all $x, y \in [a,b]^d$ it holds that $|f(x) - f(y)| \le L\|x - y\|_1$ hence ensure that for all $x = (x_1, x_2, \dots, x_d) \in [a,b]^d$ it holds that

$$
\begin{aligned}
|(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - f(x)| &= |f\big((a+b)/2, (a+b)/2, \dots, (a+b)/2\big) - f(x)| \\
&\le L\big\|\big((a+b)/2, (a+b)/2, \dots, (a+b)/2\big) - x\big\|_1 \\
&= L\sum_{i=1}^{d} |(a+b)/2 - x_i| \le \sum_{i=1}^{d} \tfrac{L(b-a)}{2} = \tfrac{dL(b-a)}{2}.
\end{aligned}
\tag{3.133}
$$

This and the fact that $\|\mathcal{T}(\mathbf{F})\|_\infty = |f((a+b)/2, (a+b)/2, \dots, (a+b)/2)| \le \sup_{x \in [a,b]^d} |f(x)|$ complete the proof of Lemma 3.2.15. $\qquad\square$

**Proposition 3.2.16.** *Let $d \in \mathbb{N}$, $L, a \in \mathbb{R}$, $b \in (a, \infty)$, $r \in (0, d/4)$, let $f \colon [a,b]^d \to \mathbb{R}$ and $\delta \colon [a,b]^d \times [a,b]^d \to \mathbb{R}$ satisfy for all $x, y \in [a,b]^d$ that $|f(x) - f(y)| \le L\|x - y\|_1$ and $\delta(x, y) = \|x - y\|_1$, and let $K \in \mathbb{N}$, $\mathfrak{x}_1, \mathfrak{x}_2, \dots, \mathfrak{x}_K \in [a,b]^d$, $\mathfrak{y} \in \mathbb{R}^K$, $\mathbf{F} \in \mathbf{N}$ satisfy $K = \mathcal{C}^{([a,b]^d, \delta), (b-a)r}$, $\sup_{x \in [a,b]^d} \big[\min_{k \in \{1,2,\dots,K\}} \delta(x, \mathfrak{x}_k)\big] \le (b - a)r$, $\mathfrak{y} = (f(\mathfrak{x}_1), f(\mathfrak{x}_2), \dots, f(\mathfrak{x}_K))$, and*

$$\mathbf{F} = \mathbb{M}_K \bullet \mathbf{A}_{-L\,\mathrm{I}_K, \mathfrak{y}} \bullet \mathbf{P}_K\big(\mathbb{L}_d \bullet \mathbf{A}_{\mathrm{I}_d, -\mathfrak{x}_1}, \mathbb{L}_d \bullet \mathbf{A}_{\mathrm{I}_d, -\mathfrak{x}_2}, \dots, \mathbb{L}_d \bullet \mathbf{A}_{\mathrm{I}_d, -\mathfrak{x}_K}\big) \bullet \mathbb{T}_{K,d} \tag{3.134}$$

*(cf. Definitions 2.2.1, 2.2.5, 2.2.9, 2.2.11, 2.2.20, 2.2.30, 3.1.16, 3.2.3, 3.2.7, and 3.2.13). Then*

*(i) it holds that $\mathcal{I}(\mathbf{F}) = d$,*

*(ii) it holds that $\mathcal{O}(\mathbf{F}) = 1$,*

*(iii)  it holds that $\mathcal{H}(\mathbf{F}) \leq \lceil d\log_2\left(\frac{3d}{4r}\right)\rceil + 1$,*

*(iv)  it holds that $\mathbb{D}_1(\mathbf{F}) \leq 2d\left(\frac{3d}{4r}\right)^d$,*

*(v)  it holds for all $i \in \{2,3,\ldots\}$ that $\mathbb{D}_i(\mathbf{F}) \leq 3\lceil\left(\frac{3d}{4r}\right)^d\frac{1}{2^{i-1}}\rceil$,*

*(vi)  it holds that $\mathcal{P}(\mathbf{F}) \leq 35\left(\frac{3d}{4r}\right)^{2d}d^2$,*

*(vii)  it holds that $\|\mathcal{T}(\mathbf{F})\|_\infty \leq \max\{1, L, |a|, |b|, 2[\sup_{x\in[a,b]^d}|f(x)|]\}$, and*

*(viii)  it holds that $\sup_{x\in[a,b]^d}|(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - f(x)| \leq 2L(b-a)r$*

*(cf. Definitions 2.1.6, 2.2.3, 2.2.36, and 3.2.8).*

*Proof of Proposition 3.2.16.* Note that the assumption that for all $x, y \in [a,b]^d$ it holds that $|f(x) - f(y)| \leq L\|x-y\|_1$ assures that $L \geq 0$. Next observe that Lemma 3.2.11, (3.134), and Proposition 3.2.12 demonstrate that

(I)  it holds that $\mathcal{I}(\mathbf{F}) = d$,

(II)  it holds that $\mathcal{O}(\mathbf{F}) = 1$,

(III)  it holds that $\mathcal{H}(\mathbf{F}) = \lceil\log_2(K)\rceil + 1$,

(IV)  it holds that $\mathbb{D}_1(\mathbf{F}) = 2dK$,

(V)  it holds for all $i \in \{2,3,\ldots\}$ that $\mathbb{D}_i(\mathbf{F}) \leq 3\lceil\frac{K}{2^{i-1}}\rceil$,

(VI)  it holds that $\|\mathcal{T}(\mathbf{F})\|_\infty \leq \max\{1, L, \max_{k\in\{1,2,\ldots,K\}}\|\mathfrak{x}_k\|_\infty, 2[\max_{k\in\{1,2,\ldots,K\}}|f(\mathfrak{x}_k)|]\}$, and

(VII)  it holds that $\sup_{x\in[a,b]^d}|(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - f(x)| \leq 2L\left[\sup_{x\in[a,b]^d}\left(\min_{k\in\{1,2,\ldots,K\}}\delta(x,\mathfrak{x}_k)\right)\right]$

(cf. Definitions 2.1.6, 2.2.3, 2.2.36, and 3.2.8). Note that items (I) and (II) establish items (i) and (ii). Next observe that item (i) in Lemma 3.2.14 and the fact that $\frac{d}{2r} \geq 2$ imply that

$$K = \mathcal{C}^{([a,b]^d,\delta),(b-a)r} \leq \left(\left\lceil\frac{d(b-a)}{2(b-a)r}\right\rceil\right)^d = \left(\lceil\frac{d}{2r}\rceil\right)^d \leq \left(\frac{3}{2}\left(\frac{d}{2r}\right)\right)^d = \left(\frac{3d}{4r}\right)^d. \tag{3.135}$$

Combining this with item (III) assures that

$$\mathcal{H}(\mathbf{F}) = \lceil\log_2(K)\rceil + 1 \leq \left\lceil\log_2\left(\left(\frac{3d}{4r}\right)^d\right)\right\rceil + 1 = \lceil d\log_2\left(\frac{3d}{4r}\right)\rceil + 1. \tag{3.136}$$

This establishes item (iii). Moreover, note that (3.135) and item (IV) imply that

$$\mathbb{D}_1(\mathbf{F}) = 2dK \leq 2d\left(\frac{3d}{4r}\right)^d. \tag{3.137}$$

This establishes item (iv). In addition, observe that item (V) and (3.135) establish item (v). Next note that item (III) ensures that for all $i \in \mathbb{N} \cap (1, \mathcal{H}(\mathbf{F})]$ it holds that

$$\frac{K}{2^{i-1}} \geq \frac{K}{2^{\mathcal{H}(\mathbf{F})-1}} = \frac{K}{2^{\lceil\log_2(K)\rceil}} \geq \frac{K}{2^{\log_2(K)+1}} = \frac{K}{2K} = \frac{1}{2}. \tag{3.138}$$

Item (V) and (3.135) hence show that for all $i \in \mathbb{N} \cap (1, \mathcal{H}(\mathbf{F})]$ it holds that

$$\mathbb{D}_i(\mathbf{F}) \leq 3\left\lceil \tfrac{K}{2^{i-1}} \right\rceil \leq \tfrac{3K}{2^{i-2}} \leq \left(\tfrac{3d}{4r}\right)^d \tfrac{3}{2^{i-2}}. \tag{3.139}$$

Furthermore, note that the fact that for all $x \in [a,b]^d$ it holds that $\|x\|_\infty \leq \max\{|a|,|b|\}$ and item (VI) imply that

$$\|\mathcal{T}(\mathbf{F})\|_\infty \leq \max\{1, L, \max_{k \in \{1,2,\ldots,K\}}\|\mathfrak{x}_k\|_\infty, 2[\max_{k \in \{1,2,\ldots,K\}}|f(\mathfrak{x}_k)|]\} \\ \leq \max\{1, L, |a|, |b|, 2[\sup_{x \in [a,b]^d}|f(x)|]\}. \tag{3.140}$$

This establishes item (vii). Moreover, observe that the assumption that for all $x \in [a,b]^d$ it holds that $\min_{k \in \{1,2,\ldots,K\}} \delta(x, \mathfrak{x}_k) \leq (b-a)r$ and item (VII) demonstrate that

$$\sup_{x \in [a,b]^d}|(\mathcal{R}_\mathfrak{r}(\mathbf{F}))(x) - f(x)| \leq 2L\left[\sup_{x \in [a,b]^d}\left(\min_{k \in \{1,2,\ldots,K\}} \delta(x, \mathfrak{x}_k)\right)\right] \\ \leq 2L(b-a)r. \tag{3.141}$$

This establishes item (viii). It thus remains to prove item (vi). For this note that items (I) and (II), (3.137), and (3.139) assure that

$$\begin{aligned} \mathcal{P}(\mathbf{F}) &= \sum_{i=1}^{\mathcal{L}(\mathbf{F})} \mathbb{D}_i(\mathbf{F})(\mathbb{D}_{i-1}(\mathbf{F}) + 1) \\ &\leq 2d\left(\tfrac{3d}{4r}\right)^d(d+1) + \left(\tfrac{3d}{4r}\right)^d 3\left(2d\left(\tfrac{3d}{4r}\right)^d + 1\right) \\ &\quad + \left[\sum_{i=3}^{\mathcal{L}(\mathbf{F})-1}\left(\tfrac{3d}{4r}\right)^d \tfrac{3}{2^{i-2}}\left(\left(\tfrac{3d}{4r}\right)^d \tfrac{3}{2^{i-3}} + 1\right)\right] + \left(\tfrac{3d}{4r}\right)^d \tfrac{3}{2^{\mathcal{L}(\mathbf{F})-3}} + 1. \end{aligned} \tag{3.142}$$

Next note that the fact that $\tfrac{3d}{4r} \geq 3$ ensures that

$$\begin{aligned} &2d\left(\tfrac{3d}{4r}\right)^d(d+1) + \left(\tfrac{3d}{4r}\right)^d 3\left(2d\left(\tfrac{3d}{4r}\right)^d + 1\right) + \left(\tfrac{3d}{4r}\right)^d \tfrac{3}{2^{\mathcal{L}(\mathbf{F})-3}} + 1 \\ &\leq \left(\tfrac{3d}{4r}\right)^{2d}\left(2d(d+1) + 3(2d+1) + \tfrac{3}{2^{1-3}} + 1\right) \\ &\leq \left(\tfrac{3d}{4r}\right)^{2d} d^2(4 + 9 + 12 + 1) = 26\left(\tfrac{3d}{4r}\right)^{2d} d^2. \end{aligned} \tag{3.143}$$

Moreover, observe that the fact that $\tfrac{3d}{4r} \geq 3$ implies that

$$\begin{aligned} \sum_{i=3}^{\mathcal{L}(\mathbf{F})-1}\left(\tfrac{3d}{4r}\right)^d \tfrac{3}{2^{i-2}}\left(\left(\tfrac{3d}{4r}\right)^d \tfrac{3}{2^{i-3}} + 1\right) &\leq \left(\tfrac{3d}{4r}\right)^{2d} \sum_{i=3}^{\mathcal{L}(\mathbf{F})-1} \tfrac{3}{2^{i-2}}\left(\tfrac{3}{2^{i-3}} + 1\right) \\ &= \left(\tfrac{3d}{4r}\right)^{2d} \sum_{i=3}^{\mathcal{L}(\mathbf{F})-1}\left[\tfrac{9}{2^{2i-5}} + \tfrac{3}{2^{i-2}}\right] \\ &= \left(\tfrac{3d}{4r}\right)^{2d} \sum_{i=0}^{\mathcal{L}(\mathbf{F})-4}\left[\tfrac{9}{2}(4^{-i}) + \tfrac{3}{2}(2^{-i})\right] \\ &\leq \left(\tfrac{3d}{4r}\right)^{2d}\left(\tfrac{9}{2}\left(\tfrac{1}{1-4^{-1}}\right) + \tfrac{3}{2}\left(\tfrac{1}{1-2^{-1}}\right)\right) = 9\left(\tfrac{3d}{4r}\right)^{2d}. \end{aligned} \tag{3.144}$$

Combining this, (3.142), and (3.143) demonstrates that

$$\mathcal{P}(\mathbf{F}) \leq 26\left(\tfrac{3d}{4r}\right)^{2d} d^2 + 9\left(\tfrac{3d}{4r}\right)^{2d} \leq 35\left(\tfrac{3d}{4r}\right)^{2d} d^2. \tag{3.145}$$

This establishes item (vi). The proof of Proposition 3.2.16 is thus complete. □

**Proposition 3.2.17.** *Let $d \in \mathbb{N}$, $L, a \in \mathbb{R}$, $b \in (a, \infty)$, $r \in (0, \infty)$ and let $f \colon [a,b]^d \to \mathbb{R}$ satisfy for all $x, y \in [a,b]^d$ that $|f(x) - f(y)| \leq L\|x - y\|_1$ (cf. Definition 3.1.16). Then there exists $\mathbf{F} \in \mathbf{N}$ such that*

(i) *it holds that $\mathcal{I}(\mathbf{F}) = d$,*

(ii) *it holds that $\mathcal{O}(\mathbf{F}) = 1$,*

(iii) *it holds that $\mathcal{H}(\mathbf{F}) \leq \left( \left\lceil d \log_2\left(\frac{3d}{4r}\right) \right\rceil + 1 \right) \mathbb{1}_{(0, d/4)}(r)$,*

(iv) *it holds that $\mathbb{D}_1(\mathbf{F}) \leq 2d\left(\frac{3d}{4r}\right)^d \mathbb{1}_{(0, d/4)}(r) + \mathbb{1}_{[d/4, \infty)}(r)$,*

(v) *it holds for all $i \in \{2, 3, \ldots\}$ that $\mathbb{D}_i(\mathbf{F}) \leq 3\left\lceil \left(\frac{3d}{4r}\right)^d \frac{1}{2^{i-1}} \right\rceil$,*

(vi) *it holds that $\mathcal{P}(\mathbf{F}) \leq 35\left(\frac{3d}{4r}\right)^{2d} d^2 \mathbb{1}_{(0, d/4)}(r) + (d+1)\mathbb{1}_{[d/4, \infty)}(r)$,*

(vii) *it holds that $\|\mathcal{T}(\mathbf{F})\|_\infty \leq \max\{1, L, |a|, |b|, 2[\sup_{x \in [a,b]^d}|f(x)|]\}$, and*

(viii) *it holds that $\sup_{x \in [a,b]^d}|(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - f(x)| \leq 2L(b-a)r$*

*(cf. Definitions 2.1.6, 2.2.1, 2.2.3, 2.2.36, and 3.2.8).*

*Proof of Proposition 3.2.17.* Throughout this proof assume w.l.o.g. that $r < d/4$ (cf. Lemma 3.2.15), let $\delta \colon [a,b]^d \times [a,b]^d \to \mathbb{R}$ satisfy for all $x, y \in [a,b]^d$ that $\delta(x,y) = \|x-y\|_1$, and let $K \in \mathbb{N} \cup \{\infty\}$ satisfy $K = \mathcal{C}^{([a,b]^d, \delta), (b-a)r}$. Note that item (i) in Lemma 3.2.14 assures that $K < \infty$. This and (3.121) ensure that there exist $\mathfrak{x}_1, \mathfrak{x}_2, \ldots, \mathfrak{x}_K \in [a,b]^d$ such that $\sup_{x \in [a,b]^d}\left[\min_{k \in \{1,2,\ldots,K\}} \delta(x, \mathfrak{x}_k)\right] \leq (b-a)r$. Combining this with Proposition 3.2.16 establishes items (i)–(viii). The proof of Proposition 3.2.17 is thus complete. $\qquad\square$

**Proposition 3.2.18.** *Let $d \in \mathbb{N}$, $L, a \in \mathbb{R}$, $b \in (a, \infty)$, $\varepsilon \in (0, 1]$ and let $f \colon [a,b]^d \to \mathbb{R}$ satisfy for all $x, y \in [a,b]^d$ that $|f(x) - f(y)| \leq L\|x - y\|_1$ (cf. Definition 3.1.16). Then there exists $\mathbf{F} \in \mathbf{N}$ such that*

(i) *it holds that $\mathcal{I}(\mathbf{F}) = d$,*

(ii) *it holds that $\mathcal{O}(\mathbf{F}) = 1$,*

(iii) *it holds that $\mathcal{H}(\mathbf{F}) \leq d\left(\max\left\{\log_2\left(\frac{3dL(b-a)}{2}\right), 0\right\} + \log_2(\varepsilon^{-1})\right) + 2$,*

(iv) *it holds that $\mathbb{D}_1(\mathbf{F}) \leq \varepsilon^{-d} d(3d\max\{L(b-a), 1\})^d$,*

(v) *it holds for all $i \in \{2, 3, \ldots\}$ that $\mathbb{D}_i(\mathbf{F}) \leq \varepsilon^{-d} 3\left(\frac{(3dL(b-a))^d}{2^i} + 1\right)$,*

(vi) *it holds that $\mathcal{P}(\mathbf{F}) \leq \varepsilon^{-2d} 9\left(3d\max\{L(b-a), 1\}\right)^{2d} d^2$,*

(vii) *it holds that $\|\mathcal{T}(\mathbf{F})\|_\infty \leq \max\{1, L, |a|, |b|, 2[\sup_{x \in [a,b]^d}|f(x)|]\}$, and*

(viii) *it holds that $\sup_{x \in [a,b]^d}|(\mathcal{R}_{\mathfrak{r}}(\mathbf{F}))(x) - f(x)| \leq \varepsilon$*

*(cf. Definitions 2.1.6, 2.2.1, 2.2.3, and 2.2.36).*

*Proof of Proposition 3.2.18.* Throughout this proof assume w.l.o.g. that $L \neq 0$. Observe that the assumption that for all $x, y \in [a,b]^d$ it holds that $|f(x) - f(y)| \leq L\|x-y\|_1$ and the assumption that $L \neq 0$ ensure that $L > 0$. Note that Proposition 3.2.17 shows that there exists $\mathbf{F} \in \mathbf{N}$ which satisfies that

(I) it holds that $\mathcal{I}(\mathbf{F}) = d$,

(II) it holds that $\mathcal{O}(\mathbf{F}) = 1$,

(III) it holds that $\mathcal{H}(\mathbf{F}) \leq \left(\left\lceil d\log_2\left(\frac{3dL(b-a)}{2\varepsilon}\right)\right\rceil + 1\right)\mathbb{1}_{(0,d/4)}\left(\frac{\varepsilon}{2L(b-a)}\right)$,

(IV) it holds that $\mathbb{D}_1(\mathbf{F}) \leq 2d\left(\frac{3dL(b-a)}{2\varepsilon}\right)^d\mathbb{1}_{(0,d/4)}\left(\frac{\varepsilon}{2L(b-a)}\right) + \mathbb{1}_{[d/4,\infty)}\left(\frac{\varepsilon}{2L(b-a)}\right)$,

(V) it holds for all $i \in \{2,3,\ldots\}$ that $\mathbb{D}_i(\mathbf{F}) \leq 3\left\lceil\left(\frac{3dL(b-a)}{2\varepsilon}\right)^d\frac{1}{2^{i-1}}\right\rceil$,

(VI) it holds that $\mathcal{P}(\mathbf{F}) \leq 35\left(\frac{3dL(b-a)}{2\varepsilon}\right)^{2d}d^2\mathbb{1}_{(0,d/4)}\left(\frac{\varepsilon}{2L(b-a)}\right) + (d+1)\mathbb{1}_{[d/4,\infty)}\left(\frac{\varepsilon}{2L(b-a)}\right)$,

(VII) it holds that $\|\mathcal{T}(\mathbf{F})\|_\infty \leq \max\{1, L, |a|, |b|, 2[\sup_{x\in[a,b]^d}|f(x)|]\}$, and

(VIII) it holds that $\sup_{x\in[a,b]^d}|(\mathcal{R}_\mathfrak{r}(\mathbf{F}))(x) - f(x)| \leq \varepsilon$

(cf. Definitions 2.1.6, 2.2.1, 2.2.3, 2.2.36, and 3.2.8). Moreover, note that item (III) assures that

$$\begin{aligned}
\mathcal{H}(\mathbf{F}) &\leq \left(d\left(\log_2\left(\tfrac{3dL(b-a)}{2}\right) + \log_2(\varepsilon^{-1})\right) + 2\right)\mathbb{1}_{(0,d/4)}\left(\tfrac{\varepsilon}{2L(b-a)}\right)\\
&\leq d\left(\max\left\{\log_2\left(\tfrac{3dL(b-a)}{2}\right), 0\right\} + \log_2(\varepsilon^{-1})\right) + 2.
\end{aligned}$$
(3.146)

In addition, observe that item (IV) implies that

$$\begin{aligned}
\mathbb{D}_1(\mathbf{F}) &\leq d\left(\tfrac{3d\max\{L(b-a),1\}}{\varepsilon}\right)^d\mathbb{1}_{(0,d/4)}\left(\tfrac{\varepsilon}{2L(b-a)}\right) + \mathbb{1}_{[d/4,\infty)}\left(\tfrac{\varepsilon}{2L(b-a)}\right)\\
&\leq \varepsilon^{-d}d(3d\max\{L(b-a),1\})^d.
\end{aligned}$$
(3.147)

Furthermore, note that item (V) ensures that for all $i \in \{2,3,\ldots\}$ it holds that

$$\mathbb{D}_i(\mathbf{F}) \leq 3\left(\left(\tfrac{3dL(b-a)}{2\varepsilon}\right)^d\tfrac{1}{2^{i-1}} + 1\right) \leq \varepsilon^{-d}3\left(\tfrac{(3dL(b-a))^d}{2^i} + 1\right).$$
(3.148)

Moreover, observe that item (VI) ensures that

$$\begin{aligned}
\mathcal{P}(\mathbf{F}) &\leq 9\left(\tfrac{3d\max\{L(b-a),1\}}{\varepsilon}\right)^{2d}d^2\mathbb{1}_{(0,d/4)}\left(\tfrac{\varepsilon}{2L(b-a)}\right) + (d+1)\mathbb{1}_{[d/4,\infty)}\left(\tfrac{\varepsilon}{2L(b-a)}\right)\\
&\leq \varepsilon^{-2d}9\left(3d\max\{L(b-a),1\}\right)^{2d}d^2.
\end{aligned}$$
(3.149)

Combining this, (3.146), (3.147), (3.148), and items (I), (II), (VII), and (VIII) establishes items (i), (ii), (iii), (iv), (v), (vi), (vii), and (viii). The proof of Proposition 3.2.18 is thus complete. $\qquad\square$

**Corollary 3.2.19.** *Let $d \in \mathbb{N}$, $L, a \in \mathbb{R}$, $b \in (a,\infty)$ and let $f: [a,b]^d \to \mathbb{R}$ satisfy for all $x, y \in [a,b]^d$ that $|f(x)-f(y)| \leq L\|x-y\|_1$ (cf. Definition 3.1.16). Then there exist $C \in \mathbb{R}$ and $\mathbf{F} = (\mathbf{F}_\varepsilon)_{\varepsilon\in(0,1]}: (0,1] \to \mathbf{N}$ such that for all $\varepsilon \in (0,1]$ it holds that $\mathcal{R}_\mathfrak{r}(\mathbf{F}_\varepsilon) \in C(\mathbb{R}^d, \mathbb{R})$, $\sup_{x\in[a,b]^d}|(\mathcal{R}_\mathfrak{r}(\mathbf{F}_\varepsilon))(x) - f(x)| \leq \varepsilon$, and $\mathcal{P}(\mathbf{F}_\varepsilon) \leq C\varepsilon^{-2d}$ (cf. Definitions 2.1.6, 2.2.1, and 2.2.3).*

*Proof of Corollary 3.2.19.* Throughout this proof let $C = 9\left(3d\max\{L(b-a),1\}\right)^{2d}d^2$. Note that items (i), (ii), (vi), and (viii) in Proposition 3.2.18 imply that for every $\varepsilon \in (0,1]$ there exists $\mathbf{F}_\varepsilon \in \mathbf{N}$ such that $\mathcal{R}_\mathfrak{r}(\mathbf{F}_\varepsilon) \in C(\mathbb{R}^d, \mathbb{R})$, $\sup_{x\in[a,b]^d}|(\mathcal{R}_\mathfrak{r}(\mathbf{F}_\varepsilon))(x) - f(x)| \leq \varepsilon$, and $\mathcal{P}(\mathbf{F}_\varepsilon) \leq C\varepsilon^{-2d}$. The proof of Corollary 3.2.19 is thus complete. $\qquad\square$

### 3.2.5 Implicit approximations through ANNs

#### 3.2.5.1 Embedding ANNs in larger architectures

**Lemma 3.2.20.** *Let* $a \in C(\mathbb{R}, \mathbb{R})$, $L \in \mathbb{N}$, $l_0, l_1, \ldots, l_L, \mathfrak{l}_0, \mathfrak{l}_1, \ldots, \mathfrak{l}_L \in \mathbb{N}$ *satisfy for all* $k \in \{1, 2, \ldots, L\}$ *that* $\mathfrak{l}_0 = l_0$, $\mathfrak{l}_L = l_L$, *and* $\mathfrak{l}_k \geq l_k$, *for every* $k \in \{1, 2, \ldots, L\}$ *let* $W_k = (W_{k,i,j})_{(i,j) \in \{1,2,\ldots,l_k\} \times \{1,2,\ldots,l_{k-1}\}} \in \mathbb{R}^{l_k \times l_{k-1}}$, $\mathscr{W}_k = (\mathscr{W}_{k,i,j})_{(i,j) \in \{1,2,\ldots,\mathfrak{l}_k\} \times \{1,2,\ldots,\mathfrak{l}_{k-1}\}} \in \mathbb{R}^{\mathfrak{l}_k \times \mathfrak{l}_{k-1}}$, $B_k = (B_{k,i})_{i \in \{1,2,\ldots,l_k\}} \in \mathbb{R}^{l_k}$, $\mathscr{B}_k = (\mathscr{B}_{k,i})_{i \in \{1,2,\ldots,\mathfrak{l}_k\}} \in \mathbb{R}^{\mathfrak{l}_k}$, *assume for all* $k \in \{1, 2, \ldots, L\}$, $i \in \{1, 2, \ldots, l_k\}$, $j \in \mathbb{N} \cap (0, l_{k-1}]$ *that* $\mathscr{W}_{k,i,j} = W_{k,i,j}$ *and* $\mathscr{B}_{k,i} = B_{k,i}$, *and assume for all* $k \in \{1, 2, \ldots, L\}$, $i \in \{1, 2, \ldots, l_k\}$, $j \in \mathbb{N} \cap (l_{k-1}, \mathfrak{l}_{k-1} + 1)$ *that* $\mathscr{W}_{k,i,j} = 0$. *Then*

$$\mathcal{R}_a\big(((W_1, B_1), (W_2, B_2), \ldots, (W_L, B_L))\big) = \mathcal{R}_a\big(((\mathscr{W}_1, \mathscr{B}_1), (\mathscr{W}_2, \mathscr{B}_2), \ldots, (\mathscr{W}_L, \mathscr{B}_L))\big) \tag{3.150}$$

*(cf. Definition 2.2.3).*

*Proof of Lemma 3.2.20.* Throughout this proof let $\pi_k \colon \mathbb{R}^{\mathfrak{l}_k} \to \mathbb{R}^{l_k}$, $k \in \{0, 1, \ldots, L\}$, satisfy for all $k \in \{0, 1, \ldots, L\}$, $x = (x_1, x_2, \ldots, x_{\mathfrak{l}_k})$ that

$$\pi_k(x) = (x_1, x_2, \ldots, x_{l_k}). \tag{3.151}$$

Observe that the assumption that $\mathfrak{l}_0 = l_0$ and $\mathfrak{l}_L = l_L$ shows that

$$\mathcal{R}_a\big(((W_1, B_1), (W_2, B_2), \ldots, (W_L, B_L))\big) \in C(\mathbb{R}^{\mathfrak{l}_0}, \mathbb{R}^{\mathfrak{l}_L}) \tag{3.152}$$

(cf. Definition 2.2.3). Furthermore, note that the assumption that for all $k \in \{1, 2, \ldots, l\}$, $i \in \{1, 2, \ldots, l_k\}$, $j \in \mathbb{N} \cap (l_{k-1}, \mathfrak{l}_{k-1} + 1)$ it holds that $\mathscr{W}_{k,i,j} = 0$ ensures that for all $k \in \{1, 2, \ldots, L\}$, $x = (x_1, x_2, \ldots, x_{\mathfrak{l}_{k-1}}) \in \mathbb{R}^{\mathfrak{l}_{k-1}}$ it holds that

$$\pi_k(\mathscr{W}_k x + \mathscr{B}_k) = \left(\left[\sum_{i=1}^{\mathfrak{l}_{k-1}} \mathscr{W}_{k,1,i} x_i\right] + \mathscr{B}_{k,1}, \left[\sum_{i=1}^{\mathfrak{l}_{k-1}} \mathscr{W}_{k,2,i} x_i\right] + \mathscr{B}_{k,2}, \ldots, \left[\sum_{i=1}^{\mathfrak{l}_{k-1}} \mathscr{W}_{k,l_k,i} x_i\right] + \mathscr{B}_{k,l_k}\right)$$
$$= \left(\left[\sum_{i=1}^{l_{k-1}} \mathscr{W}_{k,1,i} x_i\right] + \mathscr{B}_{k,1}, \left[\sum_{i=1}^{l_{k-1}} \mathscr{W}_{k,2,i} x_i\right] + \mathscr{B}_{k,2}, \ldots, \left[\sum_{i=1}^{l_{k-1}} \mathscr{W}_{k,l_k,i} x_i\right] + \mathscr{B}_{k,l_k}\right). \tag{3.153}$$

Combining this with the assumption that for all $k \in \{1, 2, \ldots, L\}$, $i \in \{1, 2, \ldots, l_k\}$, $j \in \mathbb{N} \cap (0, l_{k-1}]$ that $\mathscr{W}_{k,i,j} = W_{k,i,j}$ and $\mathscr{B}_{k,i} = B_{k,i}$ shows that for all $k \in \{1, 2, \ldots, L\}$, $x = (x_1, x_2, \ldots, x_{\mathfrak{l}_{k-1}}) \in \mathbb{R}^{\mathfrak{l}_{k-1}}$ it holds that

$$\pi_k(\mathscr{W}_k x + \mathscr{B}_k) = \left(\left[\sum_{i=1}^{l_{k-1}} W_{k,1,i} x_i\right] + B_{k,1}, \left[\sum_{i=1}^{l_{k-1}} W_{k,2,i} x_i\right] + B_{k,2}, \ldots, \left[\sum_{i=1}^{l_{k-1}} W_{k,l_k,i} x_i\right] + B_{k,l_k}\right)$$
$$= W_k \pi_{k-1}(x) + B_k. \tag{3.154}$$

Hence, we obtain that for all $x_0 \in \mathbb{R}^{\mathfrak{l}_0}$, $x_1 \in \mathbb{R}^{\mathfrak{l}_1}, \ldots, x_{L-1} \in \mathbb{R}^{\mathfrak{l}_{L-1}}$, $k \in \mathbb{N} \cap (0, L)$ with $\forall\, m \in \mathbb{N} \cap (0, L) \colon x_m = \mathfrak{M}_{a,\mathfrak{l}_m}(\mathscr{W}_m x_{m-1} + \mathscr{B}_m)$ it holds that

$$\pi_k(x_k) = \mathfrak{M}_{a,l_k}(\pi_k(\mathscr{W}_k x_{k-1} + \mathscr{B}_k)) = \mathfrak{M}_{a,l_k}(W_k \pi_{k-1}(x_{k-1}) + B_k) \tag{3.155}$$

(cf. Definition 2.1.4). Induction, the assumption that $l_0 = \mathfrak{l}_0$ and $l_L = \mathfrak{l}_L$, and (3.154) therefore prove that for all $x_0 \in \mathbb{R}^{\mathfrak{l}_0}$, $x_1 \in \mathbb{R}^{\mathfrak{l}_1}, \ldots, x_{L-1} \in \mathbb{R}^{\mathfrak{l}_{L-1}}$ with $\forall\, k \in \mathbb{N} \cap (0, L) \colon x_k = \mathfrak{M}_{a, \mathfrak{l}_k}(\mathscr{W}_k x_{k-1} + \mathscr{B}_k)$ it holds that

$$
\begin{aligned}
\big(\mathcal{R}_a\big(((W_1, B_1), (W_2, B_2), \ldots, (W_L, B_L))\big)\big)(x_0) &= \big(\mathcal{R}_a\big(((W_1, B_1), (W_2, B_2), \ldots, (W_L, B_L))\big)\big)(\pi_0(x_0)) \\
&= W_L \pi_{L-1}(x_{L-1}) + B_L \\
&= \pi_L(\mathscr{W}_L x_{L-1} + \mathscr{B}_L) = \mathscr{W}_L x_{L-1} + \mathscr{B}_L \\
&= \big(\mathcal{R}_a\big(((\mathscr{W}_1, \mathscr{B}_1), (\mathscr{W}_2, \mathscr{B}_2), \ldots, (\mathscr{W}_L, \mathscr{B}_L))\big)\big)(x_0)
\end{aligned}
\tag{3.156}
$$

(cf. Definition 2.2.3). The proof of Lemma 3.2.20 is thus complete. $\qquad\square$

**Lemma 3.2.21.** *Let $a \in C(\mathbb{R}, \mathbb{R})$, $L \in \mathbb{N}$, $l_0, l_1, \ldots, l_L, \mathfrak{l}_0, \mathfrak{l}_1, \ldots, \mathfrak{l}_L \in \mathbb{N}$ satisfy for all $k \in \{1, 2, \ldots, L\}$ that $\mathfrak{l}_0 = l_0$, $\mathfrak{l}_L = l_L$, and $\mathfrak{l}_k \geq l_k$ and let $\Phi \in \mathbf{N}$ satisfy $\mathcal{D}(\Phi) = (l_0, l_1, \ldots, l_L)$ (cf. Definition 2.2.1). Then there exists $\Psi \in \mathbf{N}$ such that*

$$
\mathcal{D}(\Psi) = (\mathfrak{l}_0, \mathfrak{l}_1, \ldots, \mathfrak{l}_L), \qquad \|\mathcal{T}(\Psi)\|_\infty = \|\mathcal{T}(\Phi)\|_\infty, \qquad \text{and} \qquad \mathcal{R}_a(\Psi) = \mathcal{R}_a(\Phi)
\tag{3.157}
$$

*(cf. Definitions 2.2.3, 2.2.36, and 3.1.16).*

*Proof of Lemma 3.2.21.* Throughout this proof let $B_k = (B_{k,i})_{i \in \{1, 2, \ldots, l_k\}} \in \mathbb{R}^{l_k}$, $k \in \{1, 2, \ldots, L\}$, and $W_k = (W_{k,i,j})_{(i,j) \in \{1, 2, \ldots, l_k\} \times \{1, 2, \ldots, l_{k-1}\}} \in \mathbb{R}^{l_k \times l_{k-1}}$, $k \in \{1, 2, \ldots, L\}$, satisfy $\Phi = ((W_1, B_1), (W_2, B_2), \ldots, (W_L, B_L))$ and let $\mathfrak{W}_k = (\mathfrak{W}_{k,i,j})_{(i,j) \in \{1, 2, \ldots, \mathfrak{l}_k\} \times \{1, 2, \ldots, \mathfrak{l}_{k-1}\}} \in \mathbb{R}^{\mathfrak{l}_k \times \mathfrak{l}_{k-1}}$, $k \in \{1, 2, \ldots, L\}$, and $\mathfrak{B}_k = (\mathfrak{B}_{k,i})_{i \in \{1, 2, \ldots, \mathfrak{l}_k\}} \in \mathbb{R}^{\mathfrak{l}_k}$, $k \in \{1, 2, \ldots, L\}$, satisfy for all $k \in \{1, 2, \ldots, L\}$, $i \in \{1, 2, \ldots, \mathfrak{l}_k\}$, $j \in \{1, 2, \ldots, \mathfrak{l}_{k-1}\}$ that

$$
\mathfrak{W}_{k,i,j} = \begin{cases} W_{k,i,j} & : (i \leq l_k) \wedge (j \leq l_{k-1}) \\ 0 & : (i > l_k) \vee (j > l_{k-1}) \end{cases} \qquad \text{and} \qquad \mathfrak{B}_{k,i} = \begin{cases} B_{k,i} & : i \leq l_k \\ 0 & : i > l_k. \end{cases}
\tag{3.158}
$$

Note that (2.51) ensures that $((\mathfrak{W}_1, \mathfrak{B}_1), (\mathfrak{W}_2, \mathfrak{B}_2), \ldots, (\mathfrak{W}_L, \mathfrak{B}_L)) \in \big(\bigtimes_{i=1}^L (\mathbb{R}^{\mathfrak{l}_i \times \mathfrak{l}_{i-1}} \times \mathbb{R}^{\mathfrak{l}_i})\big) \subseteq \mathbf{N}$ and

$$
\mathcal{D}\big(((\mathfrak{W}_1, \mathfrak{B}_1), (\mathfrak{W}_2, \mathfrak{B}_2), \ldots, (\mathfrak{W}_L, \mathfrak{B}_L))\big) = (\mathfrak{l}_0, \mathfrak{l}_1, \ldots, \mathfrak{l}_L).
\tag{3.159}
$$

Furthermore, observe that Lemma 2.2.38 and (3.158) show that

$$
\|\mathcal{T}\big(((\mathfrak{W}_1, \mathfrak{B}_1), (\mathfrak{W}_2, \mathfrak{B}_2), \ldots, (\mathfrak{W}_L, \mathfrak{B}_L))\big)\|_\infty = \|\mathcal{T}(\Phi)\|_\infty
\tag{3.160}
$$

(cf. Definitions 2.2.36 and 3.1.16). In addition, note that Lemma 3.2.20 establishes that

$$
\mathcal{R}_a(\Phi) = \mathcal{R}_a\big(((W_1, B_1), (W_2, B_2), \ldots, (W_L, B_L))\big) = \mathcal{R}_a\big(((\mathfrak{W}_1, \mathfrak{B}_1), (\mathfrak{W}_2, \mathfrak{B}_2), \ldots, (\mathfrak{W}_L, \mathfrak{B}_L))\big)
\tag{3.161}
$$

(cf. Definition 2.2.3). The proof of Lemma 3.2.21 is thus complete. $\qquad\square$

**Lemma 3.2.22.** *Let $L, \mathfrak{L} \in \mathbb{N}$, $l_0, l_1, \ldots, l_L, \mathfrak{l}_0, \mathfrak{l}_1, \ldots, \mathfrak{l}_\mathfrak{L} \in \mathbb{N}$, $\Phi_1 = ((W_1, B_1), (W_2, B_2), \ldots, (W_L, B_L)) \in \big(\bigtimes_{k=1}^L (\mathbb{R}^{l_k \times l_{k-1}} \times \mathbb{R}^{l_k})\big)$, $\Phi_2 = ((\mathfrak{W}_1, \mathfrak{B}_1), (\mathfrak{W}_2, \mathfrak{B}_2), \ldots, (\mathfrak{W}_\mathfrak{L}, \mathfrak{B}_\mathfrak{L})) \in \big(\bigtimes_{k=1}^\mathfrak{L} (\mathbb{R}^{\mathfrak{l}_k \times \mathfrak{l}_{k-1}} \times \mathbb{R}^{\mathfrak{l}_k})\big)$. Then*

$$
\|\mathcal{T}(\Phi_1 \bullet \Phi_2)\|_\infty \leq \max\big\{\|\mathcal{T}(\Phi_1)\|_\infty, \|\mathcal{T}(\Phi_2)\|_\infty, \|\mathcal{T}\big(((W_1 \mathfrak{W}_\mathfrak{L}, W_1 \mathfrak{B}_\mathfrak{L} + B_1))\big)\|_\infty\big\}
\tag{3.162}
$$

*(cf. Definitions 2.2.5, 2.2.36, and 3.1.16).*

*Proof of Lemma 3.2.22.* Note that (2.59) and Lemma 2.2.38 establish (3.162). The proof of Lemma 3.2.22 is thus complete. □

**Lemma 3.2.23.** *Let $d, L \in \mathbb{N}$, $\Phi \in \mathbf{N}$ satisfy $L \geq \mathcal{L}(\Phi)$ and $d = \mathcal{O}(\Phi)$ (cf. Definition 2.2.1). Then*

$$\|\mathcal{T}(\mathcal{E}_{L,\mathfrak{I}_d}(\Phi))\|_\infty \leq \max\{1, \|\mathcal{T}(\Phi)\|_\infty\} \tag{3.163}$$

*(cf. Definitions 2.2.18, 2.2.36, 3.1.16, and 16.2.1).*

*Proof of Lemma 3.2.23.* Throughout this proof assume w.l.o.g. that $L > \mathcal{L}(\Phi)$ and let $l_0, l_1, \ldots, l_{L-\mathcal{L}(\Phi)+1} \in \mathbb{N}$ satisfy $(l_0, l_1, \ldots, l_{L-\mathcal{L}(\Phi)+1}) = (d, 2d, 2d, \ldots, 2d, d)$. Note that Lemma 2.2.19 ensures that $\mathcal{D}(\mathfrak{I}_d) = (d, 2d, d) \in \mathbb{N}^3$ (cf. Definition 2.2.18). Item (i) in Lemma 16.2.2 hence establishes that

$$\mathcal{L}((\mathfrak{I}_d)^{\bullet(L-\mathcal{L}(\Phi))}) = L-\mathcal{L}(\Phi)+1 \qquad \text{and} \qquad \mathcal{D}((\mathfrak{I}_d)^{\bullet(L-\mathcal{L}(\Phi))}) = (l_0, l_1, \ldots, l_{L-\mathcal{L}(\Phi)+1}) \in \mathbb{N}^{L-\mathcal{L}(\Phi)+2} \tag{3.164}$$

(cf. Definition 2.2.10). This shows that there exist $W_k \in \mathbb{R}^{l_k \times l_{k-1}}$, $k \in \{1, 2, \ldots, L-\mathcal{L}(\Phi)+1\}$, and $B_k \in \mathbb{R}^{l_k}$, $k \in \{1, 2, \ldots, L-\mathcal{L}(\Phi)+1\}$, which satisfy

$$(\mathfrak{I}_d)^{\bullet(L-\mathcal{L}(\Phi))} = ((W_1, B_1), (W_2, B_2), \ldots, (W_{L-\mathcal{L}(\Phi)+1}, B_{L-\mathcal{L}(\Phi)+1})). \tag{3.165}$$

Next observe that (2.111), (2.136), (2.137), (2.59), and (2.109) demonstrate that

$$W_1 = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ -1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ 0 & -1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \\ 0 & 0 & \cdots & -1 \end{pmatrix} \in \mathbb{R}^{(2d)\times d}$$

$$\text{and} \qquad W_{L-\mathcal{L}(\Phi)+1} = \begin{pmatrix} 1 & -1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & -1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 & -1 \end{pmatrix} \in \mathbb{R}^{d\times(2d)}. \tag{3.166}$$

Moreover, note that (2.111), (2.136), (2.137), (2.59), and (2.109) prove that for all $k \in$

$\mathbb{N} \cap (1, L - \mathcal{L}(\Phi) + 1)$ it holds that

$$
W_k = \underbrace{\begin{pmatrix} 1 & 0 & \cdots & 0 \\ -1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ 0 & -1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \\ 0 & 0 & \cdots & -1 \end{pmatrix}}_{\in \mathbb{R}^{(2d) \times d}} \underbrace{\begin{pmatrix} 1 & -1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & -1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 & -1 \end{pmatrix}}_{\in \mathbb{R}^{d \times (2d)}}
$$

$$
= \begin{pmatrix} 1 & -1 & 0 & 0 & \cdots & 0 & 0 \\ -1 & 1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & -1 & \cdots & 0 & 0 \\ 0 & 0 & -1 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 & -1 \\ 0 & 0 & 0 & 0 & \cdots & -1 & 1 \end{pmatrix} \in \mathbb{R}^{(2d) \times (2d)}. \tag{3.167}
$$

In addition, observe that (2.136), (2.137), (2.111), (2.109), and (2.59) show that for all $k \in \mathbb{N} \cap [1, L - \mathcal{L}(\Phi)]$ it holds that

$$
B_k = 0 \in \mathbb{R}^{2d} \qquad \text{and} \qquad B_{L-\mathcal{L}(\Phi)+1} = 0 \in \mathbb{R}^d. \tag{3.168}
$$

Combining this, (3.166), and (3.167) establishes that

$$
\left\| \mathcal{T}\big((\mathfrak{I}_d)^{\bullet(L-\mathcal{L}(\Phi))}\big) \right\|_\infty = 1 \tag{3.169}
$$

(cf. Definitions 2.2.36 and 3.1.16). Furthermore, note that (3.166) demonstrates that for all $k \in \mathbb{N}$, $\mathfrak{W} = (w_{i,j})_{(i,j) \in \{1,2,\ldots,d\} \times \{1,2,\ldots,k\}} \in \mathbb{R}^{d \times k}$ it holds that

$$
W_1 \mathfrak{W} = \begin{pmatrix} w_{1,1} & w_{1,2} & \cdots & w_{1,k} \\ -w_{1,1} & -w_{1,2} & \cdots & -w_{1,k} \\ w_{2,1} & w_{2,2} & \cdots & w_{2,k} \\ -w_{2,1} & -w_{2,2} & \cdots & -w_{2,k} \\ \vdots & \vdots & \ddots & \vdots \\ w_{d,1} & w_{d,2} & \cdots & w_{d,k} \\ -w_{d,1} & -w_{d,2} & \cdots & -w_{d,k} \end{pmatrix} \in \mathbb{R}^{(2d) \times k}. \tag{3.170}
$$

Next observe that (3.166) and (3.168) show that for all $\mathfrak{B} = (b_1, b_2, \ldots, b_d) \in \mathbb{R}^d$ it holds that

$$
W_1 \mathfrak{B} + B_1 = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ -1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ 0 & -1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \\ 0 & 0 & \cdots & -1 \end{pmatrix} \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_d \end{pmatrix} = \begin{pmatrix} b_1 \\ -b_1 \\ b_2 \\ -b_2 \\ \vdots \\ b_d \\ -b_d \end{pmatrix} \in \mathbb{R}^{2d}. \tag{3.171}
$$

Combining this with (3.170) proves that for all $k \in \mathbb{N}$, $\mathfrak{W} \in \mathbb{R}^{d \times k}$, $\mathfrak{B} \in \mathbb{R}^d$ it holds that

$$\left\| \mathcal{T}\big(((W_1 \mathfrak{W}, W_1 \mathfrak{B} + B_1))\big) \right\|_\infty = \left\| \mathcal{T}\big(((\mathfrak{W}, \mathfrak{B}))\big) \right\|_\infty. \tag{3.172}$$

This, Lemma 3.2.22, and (3.169) establish that

$$\begin{aligned}
\left\| \mathcal{T}(\mathcal{E}_{L,\mathfrak{I}_d}(\Phi)) \right\|_\infty &= \left\| \mathcal{T}\big(((\mathfrak{I}_d)^{\bullet(L-\mathcal{L}(\Phi))}) \bullet \Phi\big) \right\|_\infty \\
&\leq \max\left\{ \left\| \mathcal{T}\big(((\mathfrak{I}_d)^{\bullet(L-\mathcal{L}(\Phi))})\big) \right\|_\infty, \left\| \mathcal{T}(\Phi) \right\|_\infty \right\} = \max\{1, \left\| \mathcal{T}(\Phi) \right\|_\infty\}
\end{aligned} \tag{3.173}$$

(cf. Definition 16.2.1). The proof of Lemma 3.2.23 is thus complete. $\qquad\square$

**Lemma 3.2.24.** *Let $L, \mathfrak{L} \in \mathbb{N}$, $l_0, l_1, \ldots, l_L, \mathfrak{l}_0, \mathfrak{l}_1, \ldots, \mathfrak{l}_{\mathfrak{L}} \in \mathbb{N}$ satisfy for all $i \in \mathbb{N} \cap [0, L)$ that $\mathfrak{L} \geq L$ $\mathfrak{l}_0 = l_0$, $\mathfrak{l}_{\mathfrak{L}} = l_L$, and $\mathfrak{l}_i \geq l_i$, assume for all $i \in \mathbb{N} \cap (L-1, \mathfrak{L})$ that $\mathfrak{l}_i \geq 2l_L$, and let $\Phi \in \mathbf{N}$ satisfy $\mathcal{D}(\Phi) = (l_0, l_1, \ldots, l_L)$ (cf. Definition 2.2.1). Then there exists $\Psi \in \mathbf{N}$ such that*

$$\mathcal{D}(\Psi) = (\mathfrak{l}_0, \mathfrak{l}_1, \ldots, \mathfrak{l}_{\mathfrak{L}}), \qquad \left\| \mathcal{T}(\Psi) \right\|_\infty \leq \max\{1, \left\| \mathcal{T}(\Phi) \right\|_\infty\}, \qquad and \qquad \mathcal{R}_{\mathfrak{r}}(\Psi) = \mathcal{R}_{\mathfrak{r}}(\Phi) \tag{3.174}$$

*(cf. Definitions 2.1.6, 2.2.3, 2.2.36, and 3.1.16).*

*Proof of Lemma 3.2.24.* Throughout this proof let $\Xi \in \mathbf{N}$ satisfy $\Xi = \mathcal{E}_{\mathfrak{L},\mathfrak{I}_{l_L}}(\Phi)$ (cf. Definitions 2.2.18 and 16.2.1). Note that item (i) in Lemma 2.2.19 demonstrates that $\mathcal{D}(\mathfrak{I}_{l_L}) = (l_L, 2l_L, l_L) \in \mathbb{N}^3$. Combining this with Lemma 16.2.4 shows that $\mathcal{D}(\Xi) \in \mathbb{N}^{\mathfrak{L}+1}$ and

$$\mathcal{D}(\Xi) = \begin{cases} (l_0, l_1, \ldots, l_L) & : \mathfrak{L} = L \\ (l_0, l_1, \ldots, l_{L-1}, 2l_L, 2l_L, \ldots, 2l_L, l_L) & : \mathfrak{L} > L. \end{cases} \tag{3.175}$$

Moreover, observe that Lemma 3.2.23 (applied with $d \curvearrowleft l_L$, $L \curvearrowleft \mathfrak{L}$, $\Phi \curvearrowleft \Phi$ in the notation of Lemma 3.2.23) establishes that

$$\left\| \mathcal{T}(\Xi) \right\|_\infty \leq \max\{1, \left\| \mathcal{T}(\Phi) \right\|_\infty\} \tag{3.176}$$

(cf. Definitions 2.2.36 and 3.1.16). In addition, note that item (iii) in Lemma 2.2.19 ensures that for all $x \in \mathbb{R}^{l_L}$ it holds that

$$(\mathcal{R}_{\mathfrak{r}}(\mathfrak{I}_{l_L}))(x) = x \tag{3.177}$$

(cf. Definitions 2.1.6 and 2.2.3). This and item (ii) in Lemma 16.2.3 prove that

$$\mathcal{R}_{\mathfrak{r}}(\Xi) = \mathcal{R}_{\mathfrak{r}}(\Phi). \tag{3.178}$$

In the next step, we observe that (3.175), the assumption that for all $i \in [0, L)$ it holds that $\mathfrak{l}_0 = l_0$, $\mathfrak{l}_{\mathfrak{L}} = l_L$, and $\mathfrak{l}_i \leq l_i$, the assumption that for all $i \in \mathbb{N} \cap (L-1, \mathfrak{L})$ it holds that $\mathfrak{l}_i \geq 2l_L$, and Lemma 3.2.21 (applied with $a \curvearrowleft \mathfrak{r}$, $L \curvearrowleft \mathfrak{L}$, $(l_0, l_1, \ldots, l_L) \curvearrowleft \mathcal{D}(\Xi)$, $(\mathfrak{l}_0, \mathfrak{l}_1, \ldots, \mathfrak{l}_{\mathfrak{L}}) \curvearrowleft (\mathfrak{l}_0, \mathfrak{l}_1, \ldots, \mathfrak{l}_{\mathfrak{L}})$, $\Phi \curvearrowleft \Xi$ in the notation of Lemma 3.2.21) ensure that there exists $\Psi \in \mathbf{N}$ such that

$$\mathcal{D}(\Psi) = (\mathfrak{l}_0, \mathfrak{l}_1, \ldots, \mathfrak{l}_{\mathfrak{L}}), \qquad \left\| \mathcal{T}(\Psi) \right\|_\infty = \left\| \mathcal{T}(\Xi) \right\|_\infty, \qquad and \qquad \mathcal{R}_{\mathfrak{r}}(\Psi) = \mathcal{R}_{\mathfrak{r}}(\Xi). \tag{3.179}$$

Combining this with (3.176) and (3.178) establishes (3.174). The proof of Lemma 3.2.24 is thus complete. $\qquad\square$

**Lemma 3.2.25.** *Let $u \in [-\infty, \infty)$, $v \in (u, \infty]$, $L, \mathfrak{L}, d, \mathfrak{d} \in \mathbb{N}$, $\theta \in \mathbb{R}^d$, $l_0, l_1, \dots, l_L, \mathfrak{l}_0, \mathfrak{l}_1, \dots, \mathfrak{l}_{\mathfrak{L}} \in \mathbb{N}$ satisfy for all $i \in \mathbb{N} \cap [0, L)$ that $d \geq \sum_{i=1}^{L} l_i(l_{i-1} + 1)$, $\mathfrak{d} \geq \sum_{i=1}^{\mathfrak{L}} \mathfrak{l}_i(\mathfrak{l}_{i-1} + 1)$, $\mathfrak{L} \geq L$, $\mathfrak{l}_0 = l_0$, $\mathfrak{l}_{\mathfrak{L}} = l_L$, and $\mathfrak{l}_i \geq l_i$ and assume for all $i \in \mathbb{N} \cap (L - 1, \mathfrak{L})$ that $\mathfrak{l}_i \geq 2l_L$. Then there exists $\vartheta \in \mathbb{R}^{\mathfrak{d}}$ such that*

$$\|\vartheta\|_\infty \leq \max\{1, \|\theta\|_\infty\} \qquad and \qquad \mathcal{N}_{u,v}^{\vartheta,(\mathfrak{l}_0,\mathfrak{l}_1,\dots,\mathfrak{l}_{\mathfrak{L}})} = \mathcal{N}_{u,v}^{\theta,(l_0,l_1,\dots,l_L)} \tag{3.180}$$

*(cf. Definitions 2.1.27 and 3.1.16).*

*Proof of Lemma 3.2.25.* Throughout this proof let $\eta_1, \eta_2, \dots, \eta_d \in \mathbb{R}$ satisfy

$$\theta = (\eta_1, \eta_2, \dots, \eta_d) \tag{3.181}$$

and let $\Phi \in \left( \bigtimes_{i=1}^{L} \mathbb{R}^{l_i \times l_{i-1}} \times \mathbb{R}^{l_i} \right)$ satisfy

$$\mathcal{T}(\Phi) = (\eta_1, \eta_2, \dots, \eta_{\mathcal{P}(\Phi)}) \tag{3.182}$$

(cf. Definition 2.2.36). Note that Lemma 3.2.24 establishes that there exists $\Psi \in \mathbf{N}$ which satisfies

$$\mathcal{D}(\Psi) = (\mathfrak{l}_0, \mathfrak{l}_1, \dots, \mathfrak{l}_{\mathfrak{L}}), \qquad \|\mathcal{T}(\Psi)\|_\infty \leq \max\{1, \|\mathcal{T}(\Phi)\|_\infty\}, \qquad and \qquad \mathcal{R}_{\mathfrak{r}}(\Psi) = \mathcal{R}_{\mathfrak{r}}(\Phi) \tag{3.183}$$

(cf. Definitions 2.1.6, 2.2.1, 2.2.3, and 3.1.16). Next let $\vartheta = (\vartheta_1, \vartheta_2, \dots, \vartheta_{\mathfrak{d}}) \in \mathbb{R}^{\mathfrak{d}}$ satisfy

$$(\vartheta_1, \vartheta_2, \dots, \vartheta_{\mathcal{P}(\Psi)}) = \mathcal{T}(\Psi) \qquad and \qquad \forall i \in \mathbb{N} \cap (\mathcal{P}(\Psi), \mathfrak{d}+1) \colon \vartheta_i = 0. \tag{3.184}$$

Note that (3.181), (3.182), (3.183), and (3.184) show that

$$\|\vartheta\|_\infty = \|\mathcal{T}(\Psi)\|_\infty \leq \max\{1, \|\mathcal{T}(\Phi)\|_\infty\} \leq \max\{1, \|\theta\|_\infty\}. \tag{3.185}$$

Next observe that Corollary 2.2.40 and (3.182) establish that for all $x \in \mathbb{R}^{l_0}$ it holds that

$$\left(\mathcal{N}_{-\infty,\infty}^{\theta,(l_0,l_1,\dots,l_L)}\right)(x) = \left(\mathcal{N}_{-\infty,\infty}^{\mathcal{T}(\Phi),\mathcal{D}(\Phi)}\right)(x) = (\mathcal{R}_{\mathfrak{r}}(\Phi))(x). \tag{3.186}$$

In addition, observe that Corollary 2.2.40, (3.183), and (3.184) prove that for all $x \in \mathbb{R}^{\mathfrak{l}_0}$ it holds that

$$\left(\mathcal{N}_{-\infty,\infty}^{\vartheta,(\mathfrak{l}_0,\mathfrak{l}_1,\dots,\mathfrak{l}_{\mathfrak{L}})}\right)(x) = \left(\mathcal{N}_{-\infty,\infty}^{\mathcal{T}(\Psi),\mathcal{D}(\Psi)}\right)(x) = (\mathcal{R}_{\mathfrak{r}}(\Psi))(x). \tag{3.187}$$

Combining this and (3.186) with (3.183) and the assumption that $\mathfrak{l}_0 = l_0$ and $\mathfrak{l}_{\mathfrak{L}} = l_L$ demonstrates that

$$\mathcal{N}_{-\infty,\infty}^{\theta,(l_0,l_1,\dots,l_L)} = \mathcal{N}_{-\infty,\infty}^{\vartheta,(\mathfrak{l}_0,\mathfrak{l}_1,\dots,\mathfrak{l}_{\mathfrak{L}})}. \tag{3.188}$$

Hence, we obtain that

$$\mathcal{N}_{u,v}^{\theta,(l_0,l_1,\dots,l_L)} = \mathfrak{C}_{u,v,l_L} \circ \mathcal{N}_{-\infty,\infty}^{\theta,(l_0,l_1,\dots,l_L)} = \mathfrak{C}_{u,v,\mathfrak{l}_{\mathfrak{L}}} \circ \mathcal{N}_{-\infty,\infty}^{\vartheta,(\mathfrak{l}_0,\mathfrak{l}_1,\dots,\mathfrak{l}_{\mathfrak{L}})} = \mathcal{N}_{u,v}^{\vartheta,(\mathfrak{l}_0,\mathfrak{l}_1,\dots,\mathfrak{l}_{\mathfrak{L}})} \tag{3.189}$$

(cf. Definition 2.1.12). This and (3.185) establish (3.180). The proof of Lemma 3.2.25 is thus complete. $\qquad \square$

### 3.2.5.2 Implicit approximation through ANNs with variable architectures

**Corollary 3.2.26.** *Let $d, K, \mathbf{d}, \mathbf{L} \in \mathbb{N}$, $\mathbf{l} = (\mathbf{l_0}, \mathbf{l_1}, \ldots, \mathbf{l_L}) \in \mathbb{N}^{\mathbf{L}+1}$, $L \in [0, \infty)$ satisfy for all $i \in \{2, 3, \ldots, \mathbf{L} - 1\}$ that $\mathbf{L} \geq \lceil \log_2(K) \rceil + 2$, $\mathbf{l_0} = d$, $\mathbf{l_L} = 1$, $\mathbf{l_1} \geq 2dK$, $\mathbf{l_i} \geq 3\lceil \frac{K}{2^{i-1}} \rceil$, and $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l_i}(\mathbf{l_{i-1}} + 1)$, let $E \subseteq \mathbb{R}^d$ be a set, let $\mathfrak{x}_1, \mathfrak{x}_2, \ldots, \mathfrak{x}_K \in E$, and let $f \colon E \to \mathbb{R}$ satisfy for all $x, y \in E$ that $|f(x) - f(y)| \leq L\|x - y\|_1$ (cf. Definitions 3.1.16 and 3.2.8). Then there exists $\theta \in \mathbb{R}^{\mathbf{d}}$ such that*

$$\|\theta\|_\infty \leq \max\{1, L, \max_{k \in \{1,2,\ldots,K\}} \|\mathfrak{x}_k\|_\infty, 2\max_{k \in \{1,2,\ldots,K\}} |f(\mathfrak{x}_k)|\} \tag{3.190}$$

*and*

$$\sup_{x \in E} |f(x) - \mathcal{N}_{-\infty,\infty}^{\theta,\mathbf{l}}(x)| \leq 2L\big[\sup_{x \in E}\big(\inf_{k \in \{1,2,\ldots,K\}} \|x - \mathfrak{x}_k\|_1\big)\big] \tag{3.191}$$

*(cf. Definition 2.1.27).*

*Proof of Corollary 3.2.26.* Throughout this proof let let $\mathfrak{y} \in \mathbb{R}^K$, $\Phi \in \mathbf{N}$ satisfy $\mathfrak{y} = (f(\mathfrak{x}_1), f(\mathfrak{x}_2), \ldots, f(\mathfrak{x}_K))$ and

$$\Phi = \mathbb{M}_K \bullet \mathbf{A}_{-L\,\mathrm{I}_K,\mathfrak{y}} \bullet \mathbf{P}_K\big(\mathbb{L}_d \bullet \mathbf{A}_{\mathrm{I}_d,-\mathfrak{x}_1}, \mathbb{L}_d \bullet \mathbf{A}_{\mathrm{I}_d,-\mathfrak{x}_2}, \ldots, \mathbb{L}_d \bullet \mathbf{A}_{\mathrm{I}_d,-\mathfrak{x}_K}\big) \bullet \mathbb{T}_{K,d} \tag{3.192}$$

(cf. Definitions 2.2.1, 2.2.5, 2.2.9, 2.2.11, 2.2.20, 2.2.30, 3.2.3, and 3.2.7). Observe that Lemma 3.2.11 and Proposition 3.2.12 establish that

(I) it holds that $\mathcal{L}(\Phi) = \lceil \log_2(K) \rceil + 2$,

(II) it holds that $\mathcal{I}(\Phi) = d$,

(III) it holds that $\mathcal{O}(\Phi) = 1$,

(IV) it holds that $\mathbb{D}_1(\Phi) = 2dK$,

(V) it holds for all $i \in \{2, 3, \ldots, \mathcal{L}(\Phi) - 1\}$ that $\mathbb{D}_i(\Phi) \leq 3\lceil \frac{K}{2^{i-1}} \rceil$,

(VI) it holds that $\|\mathcal{T}(\Phi)\|_\infty \leq \max\{1, L, \max_{k \in \{1,2,\ldots,K\}} \|\mathfrak{x}_k\|_\infty, 2\max_{k \in \{1,2,\ldots,K\}} |f(\mathfrak{x}_k)|\}$, and

(VII) it holds that $\sup_{x \in E} |f(x) - (\mathcal{R}_{\mathfrak{r}}(\Phi))(x)| \leq 2L\big[\sup_{x \in E}\big(\inf_{k \in \{1,2,\ldots,K\}} \|x - \mathfrak{x}_k\|_1\big)\big]$

(cf. Definitions 2.1.6, 2.2.3, and 2.2.36). In addition note that the fact that $\mathbf{L} \geq \lceil \log_2(K) \rceil + 2 = \mathcal{L}(\Phi)$, the fact that $\mathbf{l_0} = d = \mathbb{D}_0(\Phi)$, the fact that $\mathbf{l_1} \geq 2dK = \mathbb{D}_1(\Phi)$, the fact that for all $i \in \{2, 3, \ldots, \mathcal{L}(\Phi) - 1\}$ it holds that $\mathbf{l_i} \geq 3\lceil \frac{K}{2^{i-1}} \rceil \geq \mathbb{D}_i(\Phi)$, the fact that for all $i \in \{\mathcal{L}(\Phi), \mathcal{L}(\Phi) + 1, \ldots, \mathbf{L} - 1\}$ it holds that $\mathbf{l_i} \geq 3\lceil \frac{K}{2^{i-1}} \rceil \geq 2 = 2\mathbb{D}_{\mathcal{L}(\Phi)}(\Phi)$, and the fact that $\mathbf{l_L} = 1 = \mathbb{D}_{\mathcal{L}(\Phi)}(\Phi)$ with Lemma 3.2.25 establishes that there exists $\theta \in \mathbb{R}^{\mathbf{d}}$ which satisfies that

$$\|\theta\|_\infty \leq \max\{1, \|\mathcal{T}(\Phi)\|_\infty\} \qquad \text{and} \qquad \mathcal{N}_{-\infty,\infty}^{\theta,(\mathbf{l_0},\mathbf{l_1},\ldots,\mathbf{l_L})} = \mathcal{N}_{-\infty,\infty}^{\mathcal{T}(\Phi),\mathcal{D}(\Phi)}. \tag{3.193}$$

This and item (VI) ensure that

$$\|\theta\|_\infty \leq \max\{1, L, \max_{k \in \{1,2,\ldots,K\}} \|\mathfrak{x}_k\|_\infty, 2\max_{k \in \{1,2,\ldots,K\}} |f(\mathfrak{x}_k)|\}. \tag{3.194}$$

Moreover, note that (3.193), Corollary 2.2.40, and item (VII) assure that

$$\begin{aligned}
\sup_{x \in E} |f(x) - \mathcal{N}_{-\infty,\infty}^{\theta,(\mathbf{l_0},\mathbf{l_1},\ldots,\mathbf{l_L})}(x)| &= \sup_{x \in E} |f(x) - \mathcal{N}_{-\infty,\infty}^{\mathcal{T}(\Phi),\mathcal{D}(\Phi)}(x)| \\
&= \sup_{x \in E} |f(x) - (\mathcal{R}_{\mathfrak{r}}(\Phi))(x)| \\
&\leq 2L\big[\sup_{x \in E}\big(\inf_{k \in \{1,2,\ldots,K\}} \|x - \mathfrak{x}_k\|_1\big)\big]
\end{aligned} \tag{3.195}$$

(cf. Definition 2.1.27). The proof of Corollary 3.2.26 is thus complete. $\square$

**Corollary 3.2.27.** *Let $d, K, \mathbf{d}, \mathbf{L} \in \mathbb{N}$, $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \ldots, \mathbf{l_L}) \in \mathbb{N}^{\mathbf{L}+1}$, $L \in [0, \infty)$, $u \in [-\infty, \infty)$, $v \in (u, \infty]$ satisfy for all $i \in \{2, 3, \ldots, \mathbf{L}-1\}$ that $\mathbf{L} \geq \lceil \log_2 K \rceil + 2$, $\mathbf{l}_0 = d$, $\mathbf{l_L} = 1$, $\mathbf{l}_1 \geq 2dK$, $\mathbf{l}_i \geq 3\lceil \frac{K}{2^{i-1}} \rceil$, and $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1}+1)$, let $E \subseteq \mathbb{R}^d$ be a set, let $\mathfrak{x}_1, \mathfrak{x}_2, \ldots, \mathfrak{x}_K \in E$, and let $f \colon E \to ([u,v] \cap \mathbb{R})$ satisfy for all $x, y \in E$ that $|f(x) - f(y)| \leq L\|x - y\|_1$ (cf. Definitions 3.1.16 and 3.2.8). Then there exists $\theta \in \mathbb{R}^{\mathbf{d}}$ such that*

$$\|\theta\|_\infty \leq \max\{1, L, \max_{k \in \{1,2,\ldots,K\}}\|\mathfrak{x}_k\|_\infty, 2\max_{k \in \{1,2,\ldots,K\}}|f(\mathfrak{x}_k)|\} \tag{3.196}$$

*and*

$$\sup_{x \in E}|f(x) - \mathscr{N}_{u,v}^{\theta,\mathbf{l}}(x)| \leq 2L\big[\sup_{x \in E}\big(\inf_{k \in \{1,2,\ldots,K\}}\|x - \mathfrak{x}_k\|_1\big)\big]. \tag{3.197}$$

*(cf. Definition 2.1.27).*

*Proof of Corollary 3.2.27.* Observe that Corollary 3.2.26 implies that there exists $\theta \in \mathbb{R}^{\mathbf{d}}$ such that

$$\|\theta\|_\infty \leq \max\{1, L, \max_{k \in \{1,2,\ldots,K\}}\|\mathfrak{x}_k\|_\infty, 2\max_{k \in \{1,2,\ldots,K\}}|f(\mathfrak{x}_k)|\} \tag{3.198}$$

and

$$\sup_{x \in E}|f(x) - \mathscr{N}_{-\infty,\infty}^{\theta,\mathbf{l}}(x)| \leq 2L\big[\sup_{x \in E}\big(\inf_{k \in \{1,2,\ldots,K\}}\|x - \mathfrak{x}_k\|_1\big)\big]. \tag{3.199}$$

Moreover, observe that the assumption that $f(E) \subseteq [u,v]$ shows that for all $x \in E$ it holds that $f(x) = \mathfrak{c}_{u,v}(f(x))$ (cf. Definitions 2.1.11 and 2.1.27). The fact that for all $x, y \in \mathbb{R}$ it holds that $|\mathfrak{c}_{u,v}(x) - \mathfrak{c}_{u,v}(y)| \leq |x - y|$ and (3.199) hence establish that

$$\begin{aligned}
\sup_{x \in E}|f(x) - \mathscr{N}_{u,v}^{\theta,\mathbf{l}}(x)| &= \sup_{x \in E}|\mathfrak{c}_{u,v}(f(x)) - \mathfrak{c}_{u,v}(\mathscr{N}_{-\infty,\infty}^{\theta,\mathbf{l}}(x))| \\
&\leq \sup_{x \in E}|f(x) - \mathscr{N}_{-\infty,\infty}^{\theta,\mathbf{l}}(x)| \leq 2L\big[\sup_{x \in E}\big(\inf_{k \in \{1,2,\ldots,K\}}\|x - \mathfrak{x}_k\|_1\big)\big].
\end{aligned} \tag{3.200}$$

The proof of Corollary 3.2.27 is thus complete. $\square$

### 3.2.5.3 OLD Convergence rates for the approximation error

**Lemma 3.2.28.** *Let $d, \mathbf{d}, \mathbf{L} \in \mathbb{N}$, $L, a \in \mathbb{R}$, $b \in (a, \infty)$, $u \in [-\infty, \infty)$, $v \in (u, \infty]$, $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \ldots, \mathbf{l_L}) \in \mathbb{N}^{\mathbf{L}+1}$, assume $\mathbf{l}_0 = d$, $\mathbf{l_L} = 1$, and $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1}+1)$, and let $f \colon [a,b]^d \to ([u,v] \cap \mathbb{R})$ satisfy for all $x, y \in [a,b]^d$ that $|f(x) - f(y)| \leq L\|x - y\|_1$ (cf. Definition 3.1.16). Then there exists $\vartheta \in \mathbb{R}^{\mathbf{d}}$ such that $\|\vartheta\|_\infty \leq \sup_{x \in [a,b]^d}|f(x)|$ and*

$$\sup_{x \in [a,b]^d}|\mathscr{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - f(x)| \leq \frac{dL(b-a)}{2} \tag{3.201}$$

*(cf. Definition 2.1.27).*

*Proof of Lemma 3.2.28.* Throughout this proof let $\mathfrak{d} \in \mathbb{N}$ be given by $\mathfrak{d} = \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1}+1)$, let $\mathbf{m} = (\mathbf{m}_1, \mathbf{m}_2, \ldots, \mathbf{m}_d) \in [a,b]^d$ satisfy for all $i \in \{1, 2, \ldots, d\}$ that $\mathbf{m}_i = (a+b)/2$, and let $\vartheta = (\vartheta_1, \vartheta_2, \ldots, \vartheta_{\mathbf{d}}) \in \mathbb{R}^{\mathbf{d}}$ satisfy for all $i \in \{1, 2, \ldots, \mathbf{d}\} \setminus \{\mathfrak{d}\}$ that $\vartheta_i = 0$ and $\vartheta_{\mathfrak{d}} = f(\mathbf{m})$. Observe that the assumption that $\mathbf{l_L} = 1$ and the fact that $\forall i \in \{1, 2, \ldots, \mathfrak{d}-1\}\colon \vartheta_i = 0$ show that for all $x = (x_1, x_2, \ldots, x_{\mathbf{l_{L-1}}}) \in \mathbb{R}^{\mathbf{l_{L-1}}}$ it holds that

$$\begin{aligned}
\mathcal{A}_{1,\mathbf{l_{L-1}}}^{\vartheta,\sum_{i=1}^{\mathbf{L}-1}\mathbf{l}_i(\mathbf{l}_{i-1}+1)}(x) &= \left[\sum_{i=1}^{\mathbf{l_{L-1}}}\vartheta_{[\sum_{i=1}^{\mathbf{L}-1}\mathbf{l}_i(\mathbf{l}_{i-1}+1)]+i}x_i\right] + \vartheta_{[\sum_{i=1}^{\mathbf{L}-1}\mathbf{l}_i(\mathbf{l}_{i-1}+1)]+\mathbf{l_{L-1}}+1} \\
&= \left[\sum_{i=1}^{\mathbf{l_{L-1}}}\vartheta_{[\sum_{i=1}^{\mathbf{L}}\mathbf{l}_i(\mathbf{l}_{i-1}+1)]-(\mathbf{l_{L-1}}-i+1)}x_i\right] + \vartheta_{\sum_{i=1}^{\mathbf{L}}\mathbf{l}_i(\mathbf{l}_{i-1}+1)} \\
&= \left[\sum_{i=1}^{\mathbf{l_{L-1}}}\vartheta_{\mathfrak{d}-(\mathbf{l_{L-1}}-i+1)}x_i\right] + \vartheta_{\mathfrak{d}} = \vartheta_{\mathfrak{d}} = f(\mathbf{m})
\end{aligned} \tag{3.202}$$

(cf. Definition 2.1.1). Combining this with the fact that $f(\mathbf{m}) \in [u, v]$ ensures that for all $x \in \mathbb{R}^{\mathbf{l_{L-1}}}$ it holds that

$$
\begin{aligned}
\left(\mathfrak{C}_{u,v,\mathbf{l_L}} \circ \mathcal{A}_{\mathbf{l_L},\mathbf{l_{L-1}}}^{\vartheta, \sum_{i=1}^{\mathbf{L}-1} \mathbf{l}_i(\mathbf{l}_{i-1}+1)}\right)(x) &= \left(\mathfrak{C}_{u,v,1} \circ \mathcal{A}_{1,\mathbf{l_{L-1}}}^{\vartheta, \sum_{i=1}^{\mathbf{L}-1} \mathbf{l}_i(\mathbf{l}_{i-1}+1)}\right)(x) = \mathfrak{c}_{u,v}(f(\mathbf{m})) \\
&= \max\{u, \min\{f(\mathbf{m}), v\}\} = \max\{u, f(\mathbf{m})\} = f(\mathbf{m})
\end{aligned}
\tag{3.203}
$$

(cf. Definitions 2.1.11 and 2.1.12). This proves for all $x \in \mathbb{R}^d$ that

$$
\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) = f(\mathbf{m}). \tag{3.204}
$$

In addition, note that it holds for all $x \in [a, \mathbf{m}_1]$, $\mathfrak{x} \in [\mathbf{m}_1, b]$ that $|\mathbf{m}_1 - x| = \mathbf{m}_1 - x = {}^{(a+b)}/2 - x \leq {}^{(a+b)}/2 - a = {}^{(b-a)}/2$ and $|\mathbf{m}_1 - \mathfrak{x}| = \mathfrak{x} - \mathbf{m}_1 = \mathfrak{x} - {}^{(a+b)}/2 \leq b - {}^{(a+b)}/2 = {}^{(b-a)}/2$. The assumption that $\forall\, x, y \in [a, b]^d \colon |f(x) - f(y)| \leq L\|x - y\|_1$ and (3.204) hence demonstrate that for all $x = (x_1, x_2, \ldots, x_d) \in [a, b]^d$ it holds that

$$
\begin{aligned}
|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - f(x)| = |f(\mathbf{m}) - f(x)| \leq L\|\mathbf{m} - x\|_1 &= L\sum_{i=1}^{d}|\mathbf{m}_i - x_i| \\
= L\sum_{i=1}^{d}|\mathbf{m}_1 - x_i| \leq \sum_{i=1}^{d}\frac{L(b-a)}{2} &= \frac{dL(b-a)}{2}.
\end{aligned}
\tag{3.205}
$$

This and the fact that $\|\vartheta\|_\infty = \max_{i \in \{1,2,\ldots,\mathbf{d}\}}|\vartheta_i| = |f(\mathbf{m})| \leq \sup_{x \in [a,b]^d}|f(x)|$ complete the proof of Lemma 3.2.28. $\qquad\square$

**Proposition 3.2.29.** *Let $d, \mathbf{d}, \mathbf{L} \in \mathbb{N}$, $A \in (0, \infty)$, $L, a \in \mathbb{R}$, $b \in (a, \infty)$, $u \in [-\infty, \infty)$, $v \in (u, \infty]$, $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \ldots, \mathbf{l_L}) \in \mathbb{N}^{\mathbf{L}+1}$, assume $\mathbf{L} \geq 1 + (\lceil \log_2({}^A/(2d)) \rceil + 1)\mathbb{1}_{(6^d, \infty)}(A)$, $\mathbf{l}_0 = d$, $\mathbf{l}_1 \geq A\mathbb{1}_{(6^d, \infty)}(A)$, $\mathbf{l_L} = 1$, and $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1}+1)$, assume for all $i \in \{2, 3, \ldots, \mathbf{L}-1\}$ that $\mathbf{l}_i \geq 3\lceil {}^A/(2^i d) \rceil \mathbb{1}_{(6^d, \infty)}(A)$, and let $f\colon [a, b]^d \to ([u, v] \cap \mathbb{R})$ satisfy for all $x, y \in [a, b]^d$ that $|f(x) - f(y)| \leq L\|x - y\|_1$ (cf. Definitions 3.1.16 and 3.2.8). Then there exists $\vartheta \in \mathbb{R}^{\mathbf{d}}$ such that $\|\vartheta\|_\infty \leq \max\{1, L, |a|, |b|, 2[\sup_{x \in [a,b]^d}|f(x)|]\}$ and*

$$
\sup_{x \in [a,b]^d}|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - f(x)| \leq \frac{3dL(b-a)}{A^{1/d}} \tag{3.206}
$$

*(cf. Definition 2.1.27).*

*Proof of Proposition 3.2.29.* Throughout this proof assume w.l.o.g. that $A > 6^d$ (cf. Lemma 3.2.28), let $\mathfrak{Z} \in \mathbb{Z}$ satisfy $\mathfrak{Z} = \lfloor \left(\frac{A}{2d}\right)^{1/d} \rfloor$. Note that the fact that for all $k \in \mathbb{N}$ it holds that $2k \leq 2(2^{k-1}) = 2^k$ implies that $3^d = {}^{6^d}/2^d \leq {}^A/(2d)$. Therefore, we obtain that

$$
2 \leq \tfrac{2}{3}\left(\tfrac{A}{2d}\right)^{1/d} \leq \left(\tfrac{A}{2d}\right)^{1/d} - 1 < \mathfrak{Z}. \tag{3.207}
$$

In the next step let $r \in (0, \infty)$ satisfy $r = {}^{d(b-a)}/2\mathfrak{Z}$, let $\delta\colon [a, b]^d \times [a, b]^d \to \mathbb{R}$ satisfy for all $x, y \in [a, b]^d$ that $\delta(x, y) = \|x - y\|_1$, and let $K \in \mathbb{N} \cup \{\infty\}$ satisfy $K = \max(2, \mathcal{C}^{([a,b]^d, \delta), r})$ (cf. Definition 3.2.13). Observe that equation (3.207) and item (i) in Lemma 3.2.14 establish that

$$
K = \max\{2, \mathcal{C}^{([a,b]^d,\delta),r}\} \leq \max\left\{2, \left(\lceil \tfrac{d(b-a)}{2r} \rceil\right)^d\right\} = \max\{2, (\lceil \mathfrak{Z} \rceil)^d\} = \mathfrak{Z}^d < \infty. \tag{3.208}
$$

This implies that

$$
4 \leq 2dK \leq 2d\mathfrak{Z}^d \leq \tfrac{2dA}{2d} = A. \tag{3.209}
$$

Combining this and the fact that $\mathbf{L} \geq 1 + (\lceil \log_2(A/(2d)) \rceil + 1) \mathbb{1}_{(6^d,\infty)}(A) = \lceil \log_2(A/(2d)) \rceil + 2$ hence proves that $\lceil \log_2(K) \rceil \leq \lceil \log_2(A/(2d)) \rceil \leq \mathbf{L} - 2$. This, (3.209), the assumption that $\mathbf{l}_1 \geq A \mathbb{1}_{(6^d,\infty)}(A) = A$, and the assumption that $\forall\, i \in \{2, 3, \ldots, \mathbf{L}-1\} \colon \mathbf{l}_i \geq 3 \lceil A/(2^i d) \rceil \mathbb{1}_{(6^d,\infty)}(A) = 3\lceil A/(2^i d) \rceil$ imply that for all $i \in \{2, 3, \ldots, \mathbf{L}-1\}$ it holds that

$$\mathbf{L} \geq \lceil \log_2(K) \rceil + 2, \quad \mathbf{l}_1 \geq A \geq 2dK, \qquad \text{and} \qquad \mathbf{l}_i \geq 3\lceil \tfrac{A}{2^i d} \rceil \geq 3\lceil \tfrac{K}{2^{i-1}} \rceil. \qquad (3.210)$$

Let $\mathfrak{x}_1, \mathfrak{x}_2, \ldots, \mathfrak{x}_K \in [a,b]^d$ satisfy

$$\sup_{x \in [a,b]^d} \big[ \inf_{k \in \{1,2,\ldots,K\}} \delta(x, \mathfrak{x}_k) \big] \leq r. \qquad (3.211)$$

Observe that (3.210), the assumptions that $\mathbf{l}_0 = d$, $\mathbf{l_L} = 1$, $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1}+1)$, and $\forall\, x, y \in [a,b]^d \colon |f(x) - f(y)| \leq L\|x - y\|_1$, and Corollary 3.2.27 show that there exists $\vartheta \in \mathbb{R}^{\mathbf{d}}$ such that

$$\|\vartheta\|_\infty \leq \max\{1, L, \max_{k \in \{1,2,\ldots,K\}}\|\mathfrak{x}_k\|_\infty, 2\max_{k \in \{1,2,\ldots,K\}}|f(\mathfrak{x}_k)|\} \qquad (3.212)$$

and

$$\begin{aligned}
\sup_{x \in [a,b]^d}|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - f(x)| &\leq 2L\big[\sup_{x \in [a,b]^d}\big(\inf_{k \in \{1,2,\ldots,K\}}\|x - \mathfrak{x}_k\|_1\big)\big] \\
&= 2L\big[\sup_{x \in [a,b]^d}\big(\inf_{k \in \{1,2,\ldots,K\}}\delta(x, \mathfrak{x}_k)\big)\big].
\end{aligned} \qquad (3.213)$$

Note that (3.212) implies that

$$\|\vartheta\|_\infty \leq \max\{1, L, |a|, |b|, 2\sup_{x \in [a,b]^d}|f(x)|\}. \qquad (3.214)$$

Moreover, note that the fact that for all $k \in \mathbb{N}$ it holds that $2k \leq 2(2^{k-1}) = 2^k$, (3.213), (3.207), and (3.211) demonstrate that

$$\begin{aligned}
\sup_{x \in [a,b]^d}|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - f(x)| &\leq 2L\big[\sup_{x \in [a,b]^d}\big(\inf_{k \in \{1,2,\ldots,K\}}\delta(x, \mathfrak{x}_k)\big)\big] \leq 2Lr = \frac{dL(b-a)}{3} \\
&\leq \frac{dL(b-a)}{\frac{2}{3}\left(\frac{A}{2d}\right)^{1/d}} = \frac{(2d)^{1/d}3dL(b-a)}{2A^{1/d}} \leq \frac{3dL(b-a)}{A^{1/d}}.
\end{aligned}$$
$$(3.215)$$

Combining this with (3.214) completes the proof of Proposition 3.2.29. $\qquad\square$

**Corollary 3.2.30.** *Let $d \in \mathbb{N}$, $L, a \in \mathbb{R}$, $b \in (a, \infty)$ and let $f\colon [a,b]^d \to \mathbb{R}$ satisfy for all $x, y \in [a,b]^d$ that $|f(x) - f(y)| \leq L\|x - y\|_1$ (cf. Definition 3.1.16). Then there exist $C \in \mathbb{R}$ and $\Phi = (\Phi_\varepsilon)_{\varepsilon \in (0,1]} \colon (0,1] \to \mathbf{N}$ such that*

*(i) it holds for all $\varepsilon \in (0,1]$ that $\|\mathcal{T}(\Phi_\varepsilon)\|_\infty \leq \max\{1, L, |a|, |b|, 2[\sup_{x \in [a,b]^d}|f(x)|]\}$,*

*(ii) it holds for all $\varepsilon \in (0,1]$ that $\sup_{x \in [a,b]^d}|(\mathcal{R}_\mathfrak{r}(\Phi_\varepsilon))(x) - f(x)| \leq \varepsilon$,*

*(iii) it holds for all $\varepsilon \in (0,1]$ that $\mathcal{H}(\Phi_\varepsilon) \leq d(\log_2(\varepsilon^{-1}) + \log_2(d) + \log_2(3L(b-a)) + 1)$, and*

*(iv) it holds for all $\varepsilon \in (0,1]$ that $\mathcal{P}(\Phi_\varepsilon) \leq C\varepsilon^{-2d}$*

*(cf. Definitions 2.1.6, 2.2.1, 2.2.3, and 2.2.36).*

*Proof of Corollary 3.2.30.* Throughout this proof assume w.l.o.g. that $L > 0$, let $\varepsilon \in (0,1]$, $A \in (0,\infty)$ $\mathbf{L} \in \mathbb{N}$ satisfy $A = \left(\frac{3dL(b-a)}{\varepsilon}\right)^d$ and $\mathbf{L} = 2 + \lceil \log_2(A/(2d)) \rceil$, let $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \ldots, \mathbf{l_L}) \in \mathbb{N}^{\mathbf{L}+1}$ satisfy for all $i \in \{2, 3, \ldots, \mathbf{L}-1\}$ that $\mathbf{l}_0 = d$, $\mathbf{l_L} = 1$, $\mathbf{l}_1 = \lceil A \rceil$, and $\mathbf{l}_i = 3\lceil A/(2^i d) \rceil$, and let $c, C \in \mathbb{R}$ satisfy

$$c = \log_2(3L(b-a)) + 1 \qquad \text{and} \qquad C = \tfrac{9}{8}\big(3dL(b-a)\big)^{2d} + (d+19)\big(3dL(b-a)\big)^d + d + 11. \tag{3.216}$$

(cf. Definition 3.2.8). Observe that the fact that $\mathbf{L} \geq 1 + \big(\lceil \log_2(\frac{A}{2d}) \rceil + 1\big) \mathbb{1}_{(6^d,\infty)}(A)$, the fact that $\mathbf{l}_1 \geq A\mathbb{1}_{(6^d,\infty)}(A)$, the fact that $\mathbf{l}_0 = d$, the fact that for all $i \in \{2, 3, \ldots, \mathbf{L}-1\}$ it holds that $\mathbf{l}_i \geq 3\lceil \frac{A}{2^i d} \rceil \mathbb{1}_{(6^d,\infty)}(A)$, the fact that the fact that $\mathbf{l_L} = 1$, Proposition 3.2.29, and Corollary 2.2.40 prove that there exists $\Psi \in \big(\bigtimes_{i=1}^{\mathbf{L}}(\mathbb{R}^{l_i \times l_{i-1}} \times \mathbb{R}^{l_i})\big) \subseteq \mathbf{N}$ which satisfies $\|\mathcal{T}(\Psi)\|_\infty \leq \max\{1, L, |a|, |b|, 2[\sup_{x \in [a,b]^d} |f(x)|]\}$ and

$$\sup_{x \in [a,b]^d} |(\mathcal{R}_{\mathfrak{r}}(\Psi))(x) - f(x)| \leq \frac{3dL(b-a)}{A^{1/d}} = \varepsilon. \tag{3.217}$$

(cf. Definitions 2.1.6, 2.2.1, 2.2.3, and 2.2.36). Note that the fact that $d \geq 1$ implies that

$$\mathcal{H}(\Psi) = \mathbf{L} - 1 = 1 + \lceil \log_2(A/(2d)) \rceil = \lceil \log_2(\tfrac{A}{d}) \rceil \leq \lceil \log_2(A) \rceil = \lceil d \log_2\left(\tfrac{3dL(b-a)}{\varepsilon}\right) \rceil$$
$$\leq d\big(\log_2(\varepsilon^{-1}) + \log_2(d) + \log_2(3L(b-a))\big) + 1 \leq d(\log_2(\varepsilon^{-1}) + \log_2(d) + c). \tag{3.218}$$

Furthermore, observe that

$$\mathcal{P}(\Psi) = \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1}+1) = \lceil A \rceil(d+1) + 3\lceil \tfrac{A}{4d} \rceil(\lceil A \rceil+1) + \sum_{i=3}^{\mathbf{L}-1}\big[3\lceil \tfrac{A}{2^i d} \rceil\big(3\lceil \tfrac{A}{2^{i-1}d} \rceil + 1\big)\big] + 3\lceil \tfrac{A}{2^{\mathbf{L}-1}d} \rceil + 1. \tag{3.219}$$

Next note that the fact that $\mathbf{L} = 1 + \lceil \log_2(A/d) \rceil \geq 1 + \log_2(A/d)$, the fact that $d \geq 1$, and the fact that $\forall\, x \in \mathbb{R}\colon \lceil x \rceil \leq x + 1$ imply that

$$\lceil A \rceil(d+1) + 3\lceil \tfrac{A}{4d} \rceil(\lceil A \rceil+1) + 3\lceil \tfrac{A}{2^{\mathbf{L}-1}d} \rceil + 1 \leq (A+1)(d+1) + 3\big(\tfrac{A}{4d}+1\big)(A+2) + 4$$
$$= \tfrac{3}{4d}A^2 + (d+4+\tfrac{3}{2d})A + d + 11$$
$$\leq \tfrac{3}{4}A^2 + (d+\tfrac{11}{2})A + d + 11. \tag{3.220}$$

Moreover, observe that the fact that $\forall\, x \in (0,\infty)\colon \log_2(x) = \log_2(x/2) + 1 \leq x/2 + 1$ implies that
$$\mathbf{L} \leq 2 + \log_2(\tfrac{A}{d}) \leq 3 + \tfrac{A}{2d} \leq 3 + \tfrac{A}{2}. \tag{3.221}$$

This demonstrates that

$$\sum_{i=3}^{\mathbf{L}-1}\big[3\lceil \tfrac{A}{2^i d} \rceil\big(3\lceil \tfrac{A}{2^{i-1}d} \rceil + 1\big)\big] \leq 3\sum_{i=3}^{\mathbf{L}-1}\big(\tfrac{A}{2^i d}+1\big)\big(\tfrac{3A}{2^{i-1}d}+4\big)$$
$$= \tfrac{9A^2}{d^2}\sum_{i=3}^{\mathbf{L}-1}2^{1-2i} + \tfrac{12A}{d}\sum_{i=3}^{\mathbf{L}-1}2^{-i} + \tfrac{9A}{d}\sum_{i=3}^{\mathbf{L}-1}2^{1-i} + 12(\mathbf{L}-3)$$
$$\leq \tfrac{3}{8}A^2 + 3A + \tfrac{9}{2}A + 6A = \tfrac{3}{8}A^2 + \tfrac{27}{2}A. \tag{3.222}$$

Combining (3.216), (3.219), (3.220), and (3.222) with the fact that $\varepsilon \in (0, 1]$ shows that

$$
\begin{aligned}
\mathcal{P}(\Psi) &\leq \tfrac{9}{8}A^2 + (d + 19)A + d + 11 \\
&= \tfrac{9}{8}\big(3dL(b-a)\big)^{2d}\varepsilon^{-2d} + (d+19)\big(3dL(b-a)\big)^{d}\varepsilon^{-d} + d + 11 \\
&\leq \Big[\tfrac{9}{8}\big(3dL(b-a)\big)^{2d} + (d+19)\big(3dL(b-a)\big)^{d} + d + 11\Big]\varepsilon^{-2d} = C\varepsilon^{-2d}.
\end{aligned}
\tag{3.223}
$$

Combining this (3.218), and (3.223) completes the proof of Corollary 3.2.30.  $\square$

# Chapter 4

# Overall error analysis

## 4.1 Bias-variance decomposition

**Lemma 4.1.1** (Bias-variance decomposition)**.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $(S, \mathcal{S})$ be a measurable space, let $X\colon \Omega \to S$ and $Y\colon \Omega \to \mathbb{R}$ be random variables with $\mathbb{E}[|Y|^2] < \infty$, and let $\mathbf{r}\colon \mathcal{L}^2(\mathbb{P}_X; \mathbb{R}) \to [0, \infty)$ satisfy for all $f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ that $\mathbf{r}(f) = \mathbb{E}[|f(X) - Y|^2]$. Then*

*(i) it holds for all $f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ that*

$$\mathbf{r}(f) = \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] + \mathbb{E}\big[|Y - \mathbb{E}[Y|X]|^2\big], \tag{4.1}$$

*(ii) it holds for all $f, g \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ that*

$$\mathbf{r}(f) - \mathbf{r}(g) = \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] - \mathbb{E}\big[|g(X) - \mathbb{E}[Y|X]|^2\big], \tag{4.2}$$

*and*

*(iii) it holds for all $f, g \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ that*

$$\mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] = \mathbb{E}\big[|g(X) - \mathbb{E}[Y|X]|^2\big] + \big(\mathbf{r}(f) - \mathbf{r}(g)\big). \tag{4.3}$$

*Proof of Lemma 4.1.1.* First, observe that the assumption that for all $f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ it holds that $\mathbf{r}(f) = \mathbb{E}[|f(X) - Y|^2]$ shows that for all $f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ it holds that

$$
\begin{aligned}
\mathbf{r}(f) &= \mathbb{E}\big[|f(X) - Y|^2\big] = \mathbb{E}\big[|(f(X) - \mathbb{E}[Y|X]) + (\mathbb{E}[Y|X] - Y)|^2\big] \\
&= \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] + 2\,\mathbb{E}\big[(f(X) - \mathbb{E}[Y|X])(\mathbb{E}[Y|X] - Y)\big] + \mathbb{E}\big[|\mathbb{E}[Y|X] - Y|^2\big] \\
&= \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] + 2\,\mathbb{E}\big[\mathbb{E}\big[(f(X) - \mathbb{E}[Y|X])(\mathbb{E}[Y|X] - Y)\big|X\big]\big] + \mathbb{E}\big[|\mathbb{E}[Y|X] - Y|^2\big] \\
&= \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] + 2\,\mathbb{E}\big[(f(X) - \mathbb{E}[Y|X])\mathbb{E}\big[(\mathbb{E}[Y|X] - Y)\big|X\big]\big] + \mathbb{E}\big[|\mathbb{E}[Y|X] - Y|^2\big] \\
&= \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] + 2\,\mathbb{E}\big[(f(X) - \mathbb{E}[Y|X])(\mathbb{E}[Y|X] - \mathbb{E}[Y|X])\big] + \mathbb{E}\big[|\mathbb{E}[Y|X] - Y|^2\big] \\
&= \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] + \mathbb{E}\big[|\mathbb{E}[Y|X] - Y|^2\big].
\end{aligned}
\tag{4.4}
$$

This implies that for all $f, g \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ it holds that

$$\mathbf{r}(f) - \mathbf{r}(g) = \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] - \mathbb{E}\big[|g(X) - \mathbb{E}[Y|X]|^2\big]. \tag{4.5}$$

Hence, we obtain that for all $f, g \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ it holds that

$$\mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] = \mathbb{E}\big[|g(X) - \mathbb{E}[Y|X]|^2\big] + \mathbf{r}(f) - \mathbf{r}(g). \qquad (4.6)$$

Combining this with (4.4) and (4.5) establishes items (i), (ii), and (iii). The proof of Lemma 4.1.1 is thus complete. $\qquad\square$

## 4.2 Overall error decomposition

**Proposition 4.2.1.** *Let* $d, \mathbf{d}, \mathbf{L}, M, K, N \in \mathbb{N}$, $B \in [0, \infty)$, $u \in \mathbb{R}$, $v \in (u, \infty)$, $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \ldots, \mathbf{l}_\mathbf{L}) \in \mathbb{N}^{\mathbf{L}+1}$, $\mathbf{T} \subseteq \{0, 1, \ldots, N\}$, $D \subseteq \mathbb{R}^d$, *assume* $0 \in \mathbf{T}$, $\mathbf{l}_0 = d$, $\mathbf{l}_\mathbf{L} = 1$, *and* $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1} + 1)$, *let* $(\Omega, \mathcal{F}, \mathbb{P})$ *be a probability space, let* $X_j \colon \Omega \to D$, $j \in \{1, 2, \ldots, M\}$, *and* $Y_j \colon \Omega \to [u, v]$, $j \in \{1, 2, \ldots, M\}$, *be random variables, let* $\mathcal{E} \colon D \to [u, v]$ *be* $\mathcal{B}(D)/\mathcal{B}([u, v])$*-measurable, assume that it holds* $\mathbb{P}$*-a.s. that* $\mathcal{E}(X_1) = \mathbb{E}[Y_1|X_1]$, *let* $\Theta_{k,n} \colon \Omega \to \mathbb{R}^\mathbf{d}$, $k, n \in \mathbb{N}_0$, *satisfy* $\left(\bigcup_{k=1}^\infty \Theta_{k,0}(\Omega)\right) \subseteq [-B, B]^\mathbf{d}$, *let* $\mathbf{R} \colon \mathbb{R}^\mathbf{d} \to [0, \infty)$ *satisfy for all* $\theta \in \mathbb{R}^\mathbf{d}$ *that* $\mathbf{R}(\theta) = \mathbb{E}[|\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(X_1) - Y_1|^2]$, *and let* $\mathscr{R} \colon \mathbb{R}^\mathbf{d} \times \Omega \to [0, \infty)$ *and* $\mathbf{k} \colon \Omega \to (\mathbb{N}_0)^2$ *satisfy for all* $\theta \in \mathbb{R}^\mathbf{d}$, $\omega \in \Omega$ *that*

$$\mathscr{R}(\theta, \omega) = \frac{1}{M}\left[\sum_{j=1}^M |\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)|^2\right] \qquad and \qquad (4.7)$$

$$\mathbf{k}(\omega) \in \operatorname{argmin}_{(k,n)\in\{1,2,\ldots,K\}\times\mathbf{T}, \|\Theta_{k,n}(\omega)\|_\infty \leq B} \mathscr{R}(\Theta_{k,n}(\omega), \omega) \qquad (4.8)$$

*(cf. Definitions 2.1.27 and 3.1.16). Then it holds for all* $\vartheta \in [-B, B]^\mathbf{d}$ *that*

$$
\begin{aligned}
&\int_D |\mathcal{N}_{u,v}^{\Theta_\mathbf{k},\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) \\
&\leq \big[\sup_{x\in D}|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)|^2\big] + 2\big[\sup_{\theta\in[-B,B]^\mathbf{d}}|\mathscr{R}(\theta) - \mathbf{R}(\theta)|\big] \\
&\quad + \min_{(k,n)\in\{1,2,\ldots,K\}\times\mathbf{T}, \|\Theta_{k,n}\|_\infty \leq B}|\mathscr{R}(\Theta_{k,n}) - \mathscr{R}(\vartheta)|.
\end{aligned}
\qquad (4.9)
$$

*Proof of Proposition 4.2.1.* Throughout this proof let $\mathbf{r} \colon \mathcal{L}^2(\mathbb{P}_{X_1}; \mathbb{R}) \to [0, \infty)$ satisfy for all $f \in \mathcal{L}^2(\mathbb{P}_{X_1}; \mathbb{R})$ that $\mathbf{r}(f) = \mathbb{E}[|f(X_1) - Y_1|^2]$. Observe that the assumption that $\forall\, \omega \in \Omega \colon Y_1(\omega) \in [u, v]$ and the fact that $\forall\, \theta \in \mathbb{R}^\mathbf{d}$, $x \in \mathbb{R}^d \colon \mathcal{N}_{u,v}^{\theta,\mathbf{l}}(x) \in [u, v]$ ensure that for all $\theta \in \mathbb{R}^\mathbf{d}$ it holds that $\mathbb{E}[|Y_1|^2] \leq \max\{u^2, v^2\} < \infty$ and

$$\int_D |\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) = \mathbb{E}\big[|\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(X_1)|^2\big] \leq \max\{u^2, v^2\} < \infty. \qquad (4.10)$$

Item (iii) in Lemma 4.1.1 (applied with $(\Omega, \mathcal{F}, \mathbb{P}) \curvearrowleft (\Omega, \mathcal{F}, \mathbb{P})$, $(S, \mathcal{S}) \curvearrowleft (D, \mathcal{B}(D))$, $X \curvearrowleft X_1$, $Y \curvearrowleft (\Omega \ni \omega \mapsto Y_1(\omega) \in \mathbb{R})$, $\mathbf{r} \curvearrowleft \mathbf{r}$, $f \curvearrowleft \mathcal{N}_{u,v}^{\theta,\mathbf{l}}|_D$, $g \curvearrowleft \mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}|_D$ for $\theta, \vartheta \in \mathbb{R}^\mathbf{d}$ in the notation of item (iii) in Lemma 4.1.1) hence proves that for all $\theta, \vartheta \in \mathbb{R}^\mathbf{d}$ it holds that

$$
\begin{aligned}
&\int_D |\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) \\
&= \mathbb{E}\big[|\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(X_1) - \mathcal{E}(X_1)|^2\big] = \mathbb{E}\big[|\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] \\
&= \mathbb{E}\big[|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] + \mathbf{r}(\mathcal{N}_{u,v}^{\theta,\mathbf{l}}|_D) - \mathbf{r}(\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}|_D) \\
&= \mathbb{E}\big[|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(X_1) - \mathcal{E}(X_1)|^2\big] + \mathbb{E}\big[|\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(X_1) - Y_1|^2\big] - \mathbb{E}\big[|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(X_1) - Y_1|^2\big] \\
&= \int_D |\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) + \mathbf{R}(\theta) - \mathbf{R}(\vartheta).
\end{aligned}
\qquad (4.11)
$$

This implies that for all $\theta, \vartheta \in \mathbb{R}^{\mathbf{d}}$ it holds that

$$
\int_D |\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x)
$$

$$
= \int_D |\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) - [\mathscr{R}(\theta) - \mathbf{R}(\theta)] + \mathscr{R}(\vartheta) - \mathbf{R}(\vartheta) + \mathscr{R}(\theta) - \mathscr{R}(\vartheta)
$$

$$
\leq \int_D |\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) + |\mathscr{R}(\theta) - \mathbf{R}(\theta)| + |\mathscr{R}(\vartheta) - \mathbf{R}(\vartheta)| + \mathscr{R}(\theta) - \mathscr{R}(\vartheta)
$$

$$
\leq \int_D |\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) + 2\big[\max_{\eta \in \{\theta,\vartheta\}} |\mathscr{R}(\eta) - \mathbf{R}(\eta)|\big] + \mathscr{R}(\theta) - \mathscr{R}(\vartheta).
$$

$$
(4.12)
$$

Next note that the fact that $\forall \, \omega \in \Omega \colon \|\Theta_{\mathbf{k}(\omega)}(\omega)\|_\infty \leq B$ ensures that for all $\omega \in \Omega$ it holds that $\Theta_{\mathbf{k}(\omega)}(\omega) \in [-B, B]^{\mathbf{d}}$. Combining (4.12) with (4.8) hence establishes that for all $\vartheta \in [-B, B]^{\mathbf{d}}$ it holds that

$$
\int_D |\mathcal{N}_{u,v}^{\Theta_{\mathbf{k}},\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x)
$$

$$
\leq \int_D |\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) + 2\big[\sup_{\theta \in [-B,B]^{\mathbf{d}}} |\mathscr{R}(\theta) - \mathbf{R}(\theta)|\big] + \mathscr{R}(\Theta_{\mathbf{k}}) - \mathscr{R}(\vartheta)
$$

$$
= \int_D |\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) + 2\big[\sup_{\theta \in [-B,B]^{\mathbf{d}}} |\mathscr{R}(\theta) - \mathbf{R}(\theta)|\big]
$$

$$
\quad + \min_{(k,n) \in \{1,2,\ldots,K\} \times \mathbf{T}, \, \|\Theta_{k,n}\|_\infty \leq B} [\mathscr{R}(\Theta_{k,n}) - \mathscr{R}(\vartheta)]
$$

$$
\leq \big[\sup_{x \in D} |\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)|^2\big] + 2\big[\sup_{\theta \in [-B,B]^{\mathbf{d}}} |\mathscr{R}(\theta) - \mathbf{R}(\theta)|\big]
$$

$$
\quad + \min_{(k,n) \in \{1,2,\ldots,K\} \times \mathbf{T}, \, \|\Theta_{k,n}\|_\infty \leq B} |\mathscr{R}(\Theta_{k,n}) - \mathscr{R}(\vartheta)|.
$$

$$
(4.13)
$$

The proof of Proposition 4.2.1 is thus complete. $\qquad\square$

# Chapter 5

# Optimization through random initializations

## 5.1 Analysis of the optimization error

### 5.1.1 The complementary distribution function formula

**Lemma 5.1.1** (Complementary distribution function formula). *Let* $\mu\colon \mathcal{B}([0,\infty)) \to [0,\infty]$ *be a sigma-finite measure. Then*

$$\int_0^\infty x\,\mu(\mathrm{d}x) = \int_0^\infty \mu([x,\infty))\,\mathrm{d}x = \int_0^\infty \mu((x,\infty))\,\mathrm{d}x. \tag{5.1}$$

*Proof of Lemma 5.1.1.* First, observe that

$$\int_0^\infty x\,\mu(\mathrm{d}x) = \int_0^\infty \left[\int_0^x \mathrm{d}y\right]\mu(\mathrm{d}x) = \int_0^\infty \left[\int_0^\infty \mathbb{1}_{(-\infty,x]}(y)\,\mathrm{d}y\right]\mu(\mathrm{d}x)$$
$$= \int_0^\infty \int_0^\infty \mathbb{1}_{[y,\infty)}(x)\,\mathrm{d}y\,\mu(\mathrm{d}x). \tag{5.2}$$

Next note that the fact that $[0,\infty)^2 \ni (x,y) \mapsto \mathbb{1}_{[y,\infty)}(x) \in \mathbb{R}$ is $(\mathcal{B}([0,\infty))\otimes\mathcal{B}([0,\infty)))/\mathcal{B}(\mathbb{R})$-measurable, the assumption that $\mu$ is a sigma-finite measure, and Fubini's theorem show that

$$\int_0^\infty \int_0^\infty \mathbb{1}_{[y,\infty)}(x)\,\mathrm{d}y\,\mu(\mathrm{d}x) = \int_0^\infty \int_0^\infty \mathbb{1}_{[y,\infty)}(x)\,\mu(\mathrm{d}x)\,\mathrm{d}y = \int_0^\infty \mu([y,\infty))\,\mathrm{d}y. \tag{5.3}$$

Combining this with (5.2) demonstrates that for all $\varepsilon \in (0,\infty)$ it holds that

$$\int_0^\infty x\,\mu(\mathrm{d}x) = \int_0^\infty \mu([y,\infty))\,\mathrm{d}y \geq \int_0^\infty \mu((y,\infty))\,\mathrm{d}y$$
$$\geq \int_0^\infty \mu([y+\varepsilon,\infty))\,\mathrm{d}y = \int_\varepsilon^\infty \mu([y,\infty))\,\mathrm{d}y. \tag{5.4}$$

Beppo Levi's monotone convergence theorem hence establishes that

$$\int_0^\infty x\,\mu(\mathrm{d}x) = \int_0^\infty \mu([y,\infty))\,\mathrm{d}y \geq \int_0^\infty \mu((y,\infty))\,\mathrm{d}y$$
$$\geq \sup_{\varepsilon\in(0,\infty)}\left[\int_\varepsilon^\infty \mu([y,\infty))\,\mathrm{d}y\right] \tag{5.5}$$
$$= \sup_{\varepsilon\in(0,\infty)}\left[\int_0^\infty \mu([y,\infty))\,\mathbb{1}_{(\varepsilon,\infty)}(y)\,\mathrm{d}y\right] = \int_0^\infty \mu([y,\infty))\,\mathrm{d}y.$$

The proof of Lemma 5.1.1 is thus complete. □

## 5.1.2 Estimates for the optimization error involving complementary distribution functions

**Lemma 5.1.2.** *Let* $(E, \delta)$ *be a metric space, let* $x \in E$, $K \in \mathbb{N}$, $p, L \in (0, \infty)$, *let* $(\Omega, \mathcal{F}, \mathbb{P})$ *be a probability space, let* $\mathscr{R} \colon E \times \Omega \to \mathbb{R}$ *be* $(\mathcal{B}(E) \otimes \mathcal{F})/\mathcal{B}(\mathbb{R})$-*measurable, assume for all* $y \in E$, $\omega \in \Omega$ *that* $|\mathscr{R}(x, \omega) - \mathscr{R}(y, \omega)| \leq L\delta(x, y)$, *and let* $X_k \colon \Omega \to E$, $k \in \{1, 2, \ldots, K\}$, *be i.i.d. random variables. Then*

$$\mathbb{E}\big[\min_{k \in \{1,2,\ldots,K\}}|\mathscr{R}(X_k) - \mathscr{R}(x)|^p\big] \leq L^p \int_0^\infty [\mathbb{P}(\delta(X_1, x) > \varepsilon^{1/p})]^K \, \mathrm{d}\varepsilon. \tag{5.6}$$

*Proof of Lemma 5.1.2.* Throughout this proof let $Y \colon \Omega \to [0, \infty)$ satisfy for all $\omega \in \Omega$ that $Y(\omega) = \min_{k \in \{1,2,\ldots,K\}}[\delta(X_k(\omega), x)]^p$. Observe that the fact that $Y$ is a random variable, the assumption that $\forall\, y \in E, \omega \in \Omega \colon |\mathscr{R}(x, \omega) - \mathscr{R}(y, \omega)| \leq L\delta(x, y)$, and Lemma 5.1.1 demonstrate that

$$\begin{aligned}
\mathbb{E}\big[\min_{k \in \{1,2,\ldots,K\}}|\mathscr{R}(X_k) - \mathscr{R}(x)|^p\big] &\leq L^p\, \mathbb{E}\big[\min_{k \in \{1,2,\ldots,K\}}[\delta(X_k, x)]^p\big] \\
&= L^p\, \mathbb{E}[Y] = L^p \int_0^\infty y\, \mathbb{P}_Y(\mathrm{d}y) = L^p \int_0^\infty \mathbb{P}_Y((\varepsilon, \infty))\, \mathrm{d}\varepsilon \\
&= L^p \int_0^\infty \mathbb{P}(Y > \varepsilon)\, \mathrm{d}\varepsilon = L^p \int_0^\infty \mathbb{P}\big(\min_{k \in \{1,2,\ldots,K\}}[\delta(X_k, x)]^p > \varepsilon\big)\, \mathrm{d}\varepsilon.
\end{aligned} \tag{5.7}$$

Moreover, observe that the assumption that $X_k$, $k \in \{1, 2, \ldots, K\}$, are i.i.d. random variables shows that for all $\varepsilon \in (0, \infty)$ it holds that

$$\begin{aligned}
\mathbb{P}\big(\min_{k \in \{1,2,\ldots,K\}}[\delta(X_k, x)]^p > \varepsilon\big) &= \mathbb{P}\big(\forall\, k \in \{1, 2, \ldots, K\} \colon [\delta(X_k, x)]^p > \varepsilon\big) \\
&= \prod_{k=1}^K \mathbb{P}([\delta(X_k, x)]^p > \varepsilon) = [\mathbb{P}([\delta(X_1, x)]^p > \varepsilon)]^K = [\mathbb{P}(\delta(X_1, x) > \varepsilon^{1/p})]^K.
\end{aligned} \tag{5.8}$$

Combining this with (5.7) proves (5.6). The proof of Lemma 5.1.2 is thus complete. □

# 5.2 Strong convergences rates for the optimization error

## 5.2.1 Properties of the gamma and the beta function

**Lemma 5.2.1.** *Let* $\Gamma \colon (0, \infty) \to (0, \infty)$ *and* $\mathbb{B} \colon (0, \infty)^2 \to (0, \infty)$ *satisfy for all* $x, y \in (0, \infty)$ *that* $\Gamma(x) = \int_0^\infty t^{x-1} e^{-t}\, \mathrm{d}t$ *and* $\mathbb{B}(x, y) = \int_0^1 t^{x-1}(1-t)^{y-1}\, \mathrm{d}t$. *Then*

*(i) it holds for all* $x \in (0, \infty)$ *that* $\Gamma(x + 1) = x\, \Gamma(x)$,

*(ii) it holds that* $\Gamma(1) = \Gamma(2) = 1$, *and*

*(iii) it holds for all* $x, y \in (0, \infty)$ *that* $\mathbb{B}(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}$.

*Proof of Lemma 5.2.1.* Throughout this proof let $x, y \in (0, \infty)$, let $\Phi \colon (0, \infty) \times (0, 1) \to (0, \infty)^2$ satisfy for all $u \in (0, \infty)$, $v \in (0, 1)$ that

$$\Phi(u, v) = (u(1 - v), uv), \tag{5.9}$$

and let $f \colon (0, \infty)^2 \to (0, \infty)$ satisfy for all $s, t \in (0, \infty)$ that

$$f(s, t) = s^{(x-1)} t^{(y-1)} e^{-(s+t)}. \tag{5.10}$$

Observe that the integration by parts formula assures that for all $x \in (0, \infty)$ it holds that

$$\begin{aligned}
\Gamma(x + 1) &= \int_0^\infty t^{((x+1)-1)} e^{-t} \, dt = -\int_0^\infty t^x \left[ -e^{-t} \right] dt \\
&= -\left( \left[ t^x e^{-t} \right]_{t=0}^{t=\infty} - x \int_0^\infty t^{(x-1)} e^{-t} \, dt \right) = x \int_0^\infty t^{(x-1)} e^{-t} \, dt = x \cdot \Gamma(x).
\end{aligned} \tag{5.11}$$

This establishes item (i). Moreover, note that

$$\Gamma(1) = \int_0^\infty t^0 e^{-t} \, dt = [-e^{-t}]_{t=0}^{t=\infty} = 1. \tag{5.12}$$

This and item (i) prove item (ii). Next note that the integral transformation theorem with the diffeomorphism $(1, \infty) \ni t \mapsto \frac{1}{t} \in (0, 1)$ ensures that

$$\begin{aligned}
B(x, y) &= \int_0^1 t^{(x-1)} (1 - t)^{(y-1)} dt = \int_1^\infty \left[ \tfrac{1}{t} \right]^{(x-1)} \left[ 1 - \tfrac{1}{t} \right]^{(y-1)} \tfrac{1}{t^2} \, dt \\
&= \int_1^\infty t^{(-x-1)} \left[ \tfrac{t-1}{t} \right]^{(y-1)} dt = \int_1^\infty t^{(-x-y)} (t - 1)^{(y-1)} dt \\
&= \int_0^\infty (t + 1)^{(-x-y)} t^{(y-1)} dt = \int_0^\infty \frac{t^{(y-1)}}{(t + 1)^{(x+y)}} \, dt.
\end{aligned} \tag{5.13}$$

In addition, note that

$$\begin{aligned}
\Gamma(x) \cdot \Gamma(y) &= \left[ \int_0^\infty t^{(x-1)} e^{-t} \, dt \right] \left[ \int_0^\infty t^{(y-1)} e^{-t} \, dt \right] \\
&= \left[ \int_0^\infty s^{(x-1)} e^{-s} \, ds \right] \left[ \int_0^\infty t^{(y-1)} e^{-t} \, dt \right] \\
&= \int_0^\infty \int_0^\infty s^{(x-1)} t^{(y-1)} e^{-(s+t)} \, dt \, ds \\
&= \int_{(0,\infty)^2} f(s, t) \, d(s, t).
\end{aligned} \tag{5.14}$$

Moreover, observe that for all $(u, v) \in (0, \infty) \times (0, 1)$ it holds that

$$\Phi'(u, v) = \begin{bmatrix} 1 - v & -u \\ v & u \end{bmatrix}. \tag{5.15}$$

Hence, we obtain that for all $(u, v) \in (0, \infty) \times (0, 1)$ it holds that

$$\det(\Phi'(u, v)) = (1 - v)u - v(-u) = u - vu + vu = u \in (0, \infty). \tag{5.16}$$

This, (5.14), and the integral transformation theorem show that

$$
\begin{aligned}
\Gamma(x) \cdot \Gamma(y) &= \int_{(0,\infty) \times (0,1)} f(\Phi(u,v)) \, |\det(\Phi'(u,v))| \, d(u,v) \\
&= \int_0^\infty \int_0^1 (u(1-v))^{(x-1)} \, (uv)^{(y-1)} \, e^{-(u(1-v)+uv)} \, u \, dv \, du \\
&= \int_0^\infty \int_0^1 u^{(x+y-1)} \, e^{-u} \, v^{(y-1)} \, (1-v)^{(x-1)} \, dv \, du \\
&= \left[ \int_0^\infty u^{(x+y-1)} \, e^{-u} \, du \right] \left[ \int_0^1 v^{(y-1)} \, (1-v)^{(x-1)} \, dv \right] \\
&= \Gamma(x+y) \, B(y,x).
\end{aligned}
\tag{5.17}
$$

This establishes item (iii). The proof of Lemma 5.2.1 is thus complete.  □

**Lemma 5.2.2.** *It holds for all $\alpha, x \in [0,1]$ that $(1-x)^\alpha \le 1 - \alpha x$.*

*Proof of Lemma 5.2.2.* Note that the fact that for all $y \in [0,\infty)$ it holds that $[0,\infty) \ni z \mapsto y^z \in [0,\infty)$ is convex implies that for all $\alpha, x \in [0,1]$ it holds that

$$
\begin{aligned}
(1-x)^\alpha &\le \alpha(1-x)^1 + (1-\alpha)(1-x)^0 \\
&= \alpha - \alpha x + 1 - \alpha = 1 - \alpha x.
\end{aligned}
\tag{5.18}
$$

The proof of Lemma 5.2.2 is thus complete.  □

**Proposition 5.2.3.** *Let $\Gamma \colon (0,\infty) \to (0,\infty)$ and $\lfloor \cdot \rfloor \colon (0,\infty) \to \mathbb{N}_0$ satisfy for all $x \in (0,\infty)$ that $\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} \, dt$ and $\lfloor x \rfloor = \max([0,x) \cap \mathbb{N}_0)$. Then*

*(i) it holds that $\Gamma \colon (0,\infty) \to (0,\infty)$ is convex,*

*(ii) it holds for all $x \in (0,\infty)$ that $\Gamma(x+1) = x\,\Gamma(x) \le x^{\lfloor x \rfloor} \le \max\{1, x^x\}$,*

*(iii) it holds for all $x \in (0,\infty)$, $\alpha \in [0,1]$ that*

$$
(\max\{x+\alpha-1, 0\})^\alpha \le \frac{x}{(x+\alpha)^{1-\alpha}} \le \frac{\Gamma(x+\alpha)}{\Gamma(x)} \le x^\alpha,
\tag{5.19}
$$

*and*

*(iv) it holds for all $x \in (0,\infty)$, $\alpha \in [0,\infty)$ that*

$$
(\max\{x + \min\{\alpha-1, 0\}, 0\})^\alpha \le \frac{\Gamma(x+\alpha)}{\Gamma(x)} \le (x + \max\{\alpha-1, 0\})^\alpha.
\tag{5.20}
$$

*Proof of Proposition 5.2.3.* Throughout this proof let $\lfloor \cdot \rfloor \colon [0,\infty) \to \mathbb{N}_0$ satisfy for all $x \in [0,\infty)$ that $\lfloor x \rfloor = \max([0,x] \cap \mathbb{N}_0)$. Observe that the fact that for all $t \in (0,\infty)$ it holds that $\mathbb{R} \ni x \mapsto t^x \in (0,\infty)$ is convex implies that for all $x, y \in (0,\infty)$, $\alpha \in [0,1]$ it holds that

$$
\begin{aligned}
\Gamma(\alpha x + (1-\alpha)y) &= \int_0^\infty t^{\alpha x + (1-\alpha)y - 1} e^{-t} \, dt = \int_0^\infty t^{\alpha x + (1-\alpha)y} t^{-1} e^{-t} \, dt \\
&\le \int_0^\infty (\alpha t^x + (1-\alpha)t^y) t^{-1} e^{-t} \, dt \\
&= \alpha \int_0^\infty t^{x-1} e^{-t} \, dt + (1-\alpha) \int_0^\infty t^{y-1} e^{-t} \, dt \\
&= \alpha \, \Gamma(x) + (1-\alpha)\Gamma(y).
\end{aligned}
\tag{5.21}
$$

This shows item (i). Next note that item (ii) in Lemma 5.2.1 and item (i) establish that for all $\alpha \in [0, 1]$ it holds that

$$\Gamma(\alpha + 1) = \Gamma(\alpha \cdot 2 + (1 - \alpha) \cdot 1) \leq \alpha\,\Gamma(2) + (1 - \alpha)\Gamma(1) = \alpha + (1 - \alpha) = 1. \tag{5.22}$$

This yields for all $x \in (0, 1]$ that

$$\Gamma(x + 1) \leq 1 = x^{\lceil x \rceil} = \max\{1, x^x\}. \tag{5.23}$$

Induction, item (i) in Lemma 5.2.1, and the fact that $\forall\, x \in (0, \infty) \colon x - \lfloor x \rfloor \in (0, 1]$ hence ensure that for all $x \in [1, \infty)$ it holds that

$$\Gamma(x + 1) = \left[\prod_{i=1}^{\lceil x \rceil}(x - i + 1)\right]\Gamma(x - \lfloor x \rfloor + 1) \leq x^{\lceil x \rceil}\Gamma(x - \lfloor x \rfloor + 1) \leq x^{\lceil x \rceil} \leq x^x = \max\{1, x^x\}. \tag{5.24}$$

Combining this and (5.23) with item (i) in Lemma 5.2.1 establishes item (ii). Furthermore, note that Hölder's inequality and item (i) in Lemma 5.2.1 prove that for all $x \in (0, \infty)$, $\alpha \in [0, 1]$ it holds that

$$
\begin{aligned}
\Gamma(x + \alpha) &= \int_0^\infty t^{x + \alpha - 1}e^{-t}\,\mathrm{d}t = \int_0^\infty t^{\alpha x}e^{-\alpha t}t^{(1-\alpha)x - (1-\alpha)}e^{-(1-\alpha)t}\,\mathrm{d}t \\
&= \int_0^\infty [t^x e^{-t}]^\alpha [t^{x-1}e^{-t}]^{1-\alpha}\,\mathrm{d}t \\
&\leq \left(\int_0^\infty t^x e^{-t}\,\mathrm{d}t\right)^\alpha \left(\int_0^\infty t^{x-1}e^{-t}\,\mathrm{d}t\right)^{1-\alpha} \\
&= [\Gamma(x + 1)]^\alpha [\Gamma(x)]^{1-\alpha} = x^\alpha [\Gamma(x)]^\alpha [\Gamma(x)]^{1-\alpha} \\
&= x^\alpha \Gamma(x).
\end{aligned}
\tag{5.25}
$$

This and item (i) in Lemma 5.2.1 demonstrate that for all $x \in (0, \infty)$, $\alpha \in [0, 1]$ it holds that

$$x\,\Gamma(x) = \Gamma(x + 1) = \Gamma(x + \alpha + (1 - \alpha)) \leq (x + \alpha)^{1-\alpha}\Gamma(x + \alpha). \tag{5.26}$$

Combining (5.25) and (5.26) yields that for all $x \in (0, \infty)$, $\alpha \in [0, 1]$ it holds that

$$\frac{x}{(x + \alpha)^{1-\alpha}} \leq \frac{\Gamma(x + \alpha)}{\Gamma(x)} \leq x^\alpha. \tag{5.27}$$

Furthermore, observe that item (i) in Lemma 5.2.1 and (5.27) imply that for all $x \in (0, \infty)$, $\alpha \in [0, 1]$ it holds that

$$\frac{\Gamma(x + \alpha)}{\Gamma(x + 1)} = \frac{\Gamma(x + \alpha)}{x\,\Gamma(x)} \leq x^{\alpha - 1}. \tag{5.28}$$

This shows for all $\alpha \in [0, 1]$, $x \in (\alpha, \infty)$ that

$$\frac{\Gamma(x)}{\Gamma(x + (1 - \alpha))} = \frac{\Gamma((x - \alpha) + \alpha)}{\Gamma((x - \alpha) + 1)} \leq (x - \alpha)^{\alpha - 1} = \frac{1}{(x - \alpha)^{1-\alpha}}. \tag{5.29}$$

This, in turn, ensures for all $\alpha \in [0, 1]$, $x \in (1 - \alpha, \infty)$ that

$$(x + \alpha - 1)^\alpha = (x - (1 - \alpha))^\alpha \leq \frac{\Gamma(x + \alpha)}{\Gamma(x)}. \tag{5.30}$$

Next note that Lemma 5.2.2 proves that for all $x \in (0, \infty)$, $\alpha \in [0, 1]$ it holds that

$$
\begin{aligned}
(\max\{x + \alpha - 1, 0\})^\alpha &= (x + \alpha)^\alpha \left( \frac{\max\{x + \alpha - 1, 0\}}{x + \alpha} \right)^\alpha \\
&= (x + \alpha)^\alpha \left( \max\left\{ 1 - \frac{1}{x + \alpha}, 0 \right\} \right)^\alpha \\
&\leq (x + \alpha)^\alpha \left( 1 - \frac{\alpha}{x + \alpha} \right) = (x + \alpha)^\alpha \left( \frac{x}{x + \alpha} \right) \\
&= \frac{x}{(x + \alpha)^{1-\alpha}}.
\end{aligned}
\tag{5.31}
$$

This and (5.27) establish item (iii). Moreover, observe that induction, item (i) in Lemma 5.2.1, the fact that $\forall\, \alpha \in [0, \infty)\colon \alpha - \lfloor \alpha \rfloor \in [0, 1)$, and item (iii) demonstrate that for all $x \in (0, \infty)$, $\alpha \in [0, \infty)$ it holds that

$$
\begin{aligned}
\frac{\Gamma(x + \alpha)}{\Gamma(x)} &= \left[ \prod_{i=1}^{\lfloor \alpha \rfloor} (x + \alpha - i) \right] \frac{\Gamma(x + \alpha - \lfloor \alpha \rfloor)}{\Gamma(x)} \leq \left[ \prod_{i=1}^{\lfloor \alpha \rfloor} (x + \alpha - i) \right] x^{\alpha - \lfloor \alpha \rfloor} \\
&\leq (x + \alpha - 1)^{\lfloor \alpha \rfloor} x^{\alpha - \lfloor \alpha \rfloor} \\
&\leq (x + \max\{\alpha - 1, 0\})^{\lfloor \alpha \rfloor} (x + \max\{\alpha - 1, 0\})^{\alpha - \lfloor \alpha \rfloor} \\
&= (x + \max\{\alpha - 1, 0\})^\alpha.
\end{aligned}
\tag{5.32}
$$

Furthermore, the fact that $\forall\, \alpha \in [0, \infty)\colon \alpha - \lfloor \alpha \rfloor \in [0, 1)$, item (iii), induction, and item (i) in Lemma 5.2.1 imply that for all $x \in (0, \infty)$, $\alpha \in [0, \infty)$ it holds that

$$
\begin{aligned}
\frac{\Gamma(x + \alpha)}{\Gamma(x)} &= \frac{\Gamma(x + \lfloor \alpha \rfloor + \alpha - \lfloor \alpha \rfloor)}{\Gamma(x)} \\
&\geq (\max\{x + \lfloor \alpha \rfloor + \alpha - \lfloor \alpha \rfloor - 1, 0\})^{\alpha - \lfloor \alpha \rfloor} \left[ \frac{\Gamma(x + \lfloor \alpha \rfloor)}{\Gamma(x)} \right] \\
&= (\max\{x + \alpha - 1, 0\})^{\alpha - \lfloor \alpha \rfloor} \left[ \prod_{i=1}^{\lfloor \alpha \rfloor} (x + \lfloor \alpha \rfloor - i) \right] \frac{\Gamma(x)}{\Gamma(x)} \\
&\geq (\max\{x + \alpha - 1, 0\})^{\alpha - \lfloor \alpha \rfloor} x^{\lfloor \alpha \rfloor} \\
&= (\max\{x + \alpha - 1, 0\})^{\alpha - \lfloor \alpha \rfloor} (\max\{x, 0\})^{\lfloor \alpha \rfloor} \\
&\geq (\max\{x + \min\{\alpha - 1, 0\}, 0\})^{\alpha - \lfloor \alpha \rfloor} (\max\{x + \min\{\alpha - 1, 0\}, 0\})^{\lfloor \alpha \rfloor} \\
&= (\max\{x + \min\{\alpha - 1, 0\}, 0\})^\alpha.
\end{aligned}
\tag{5.33}
$$

Combining this with (5.32) shows item (iv). The proof of Proposition 5.2.3 is thus complete. $\qquad\square$

**Corollary 5.2.4.** *Let* $\mathbb{B}\colon (0, \infty)^2 \to (0, \infty)$ *satisfy for all* $x, y \in (0, \infty)$ *that* $\mathbb{B}(x, y) = \int_0^1 t^{x-1}(1 - t)^{y-1}\, dt$ *and let* $\Gamma\colon (0, \infty) \to (0, \infty)$ *satisfy for all* $x \in (0, \infty)$ *that* $\Gamma(x) = \int_0^\infty t^{x-1} e^{-t}\, dt$. *Then it holds for all* $x, y \in (0, \infty)$ *with* $x + y > 1$ *that*

$$
\frac{\Gamma(x)}{(y + \max\{x - 1, 0\})^x} \leq \mathbb{B}(x, y) \leq \frac{\Gamma(x)}{(y + \min\{x - 1, 0\})^x} \leq \frac{\max\{1, x^x\}}{x(y + \min\{x - 1, 0\})^x}.
\tag{5.34}
$$

*Proof of Corollary 5.2.4.* Note that item (iii) in Lemma 5.2.1 ensures that for all $x, y \in (0, \infty)$ it holds that

$$\mathbb{B}(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(y + x)}. \tag{5.35}$$

In addition, observe that for all $x, y \in (0, \infty)$ with $x + y > 1$ it holds that $y + \min\{x - 1, 0\} > 0$. This and item (iv) in Proposition 5.2.3 demonstrate that for all $x, y \in (0, \infty)$ with $x + y > 1$ it holds that

$$0 < (y + \min\{x - 1, 0\})^x \le \frac{\Gamma(y + x)}{\Gamma(y)} \le (y + \max\{x - 1, 0\})^x. \tag{5.36}$$

Combining this with (5.35) and item (ii) in Proposition 5.2.3 shows that for all $x, y \in (0, \infty)$ with $x + y > 1$ it holds that

$$\frac{\Gamma(x)}{(y + \max\{x - 1, 0\})^x} \le \mathbb{B}(x, y) \le \frac{\Gamma(x)}{(y + \min\{x - 1, 0\})^x} \le \frac{\max\{1, x^x\}}{x(y + \min\{x - 1, 0\})^x}. \tag{5.37}$$

The proof of Corollary 5.2.4 is thus complete. $\qquad\square$

## 5.2.2   Product measurability of continuous random fields

**Lemma 5.2.5** (Projections in metric spaces). *Let $(E, d)$ be a metric space, let $n \in \mathbb{N}$, $e_1, e_2, \ldots, e_n \in E$, and let $P \colon E \to E$ satisfy for all $x \in E$ that*

$$P(x) = e_{\min\{k \in \{1, 2, \ldots, n\} \colon d(x, e_k) = \min\{yd(x, e_1), d(x, e_2), \ldots, d(x, e_n)\}\}}. \tag{5.38}$$

*Then*

*(i) it holds for all $x \in E$ that*

$$d(x, P(x)) = \min_{k \in \{1, 2, \ldots, n\}} d(x, e_k) \tag{5.39}$$

*and*

*(ii) it holds for all $A \subseteq E$ that $P^{-1}(A) \in \mathcal{B}(E)$.*

*Proof of Lemma 5.2.5.* Throughout this proof let $D = (D_1, \ldots, D_n) \colon E \to \mathbb{R}^n$ satisfy for all $x \in E$ that

$$D(x) = (D_1(x), D_2(x), \ldots, D_n(x)) = (d(x, e_1), d(x, e_2), \ldots, d(x, e_n)). \tag{5.40}$$

Note that (5.38) ensures that for all $x \in E$ it holds that

$$d(x, P(x)) = d(x, e_{\min\{k \in \{1, 2, \ldots, n\} \colon d(x, e_k) = \min\{d(x, e_1), d(x, e_2), \ldots, d(x, e_n)\}\}}) = \min_{k \in \{1, 2, \ldots, n\}} d(x, e_k). \tag{5.41}$$

This establishes item (i). It thus remains to prove item (ii). For this observe that the fact that $d \colon E \times E \to [0, \infty)$ is continuous ensures that $D \colon E \to \mathbb{R}^n$ is continuous. Hence, we obtain that $D \colon E \to \mathbb{R}^n$ is $\mathcal{B}(E)/\mathcal{B}(\mathbb{R}^n)$-measurable. Next note that item (i) demonstrates that for all $k \in \{1, 2, \ldots, n\}$, $x \in P^{-1}(\{e_k\})$ it holds that

$$d(x, e_k) = d(x, P(x)) = \min_{l \in \{1, 2, \ldots, n\}} d(x, e_l). \tag{5.42}$$

Hence, we obtain that for all $k \in \{1, 2, \ldots, n\}$, $x \in P^{-1}(\{e_k\})$ it holds that

$$k \geq \min\{l \in \{1, 2, \ldots, n\} \colon d(x, e_l) = \min\{d(x, e_1), d(x, e_2), \ldots, d(x, e_n)\}\}. \tag{5.43}$$

Moreover, note that (5.38) ensures that for all $k \in \{1, 2, \ldots, n\}$, $x \in P^{-1}(\{e_k\})$ it holds that

$$\min\left\{l \in \{1, 2, \ldots, n\} \colon d(x, e_l) = \min_{u \in \{1,2,\ldots,n\}} d(x, e_u)\right\} \\ \in \{l \in \{1, 2, \ldots, n\} \colon e_l = e_k\} \subseteq \{k, k+1, \ldots, n\}. \tag{5.44}$$

Therefore, we obtain that for all $k \in \{1, 2, \ldots, n\}$, $x \in P^{-1}(\{e_k\})$ with $e_k \notin (\cup_{l \in \mathbb{N} \cap [0,k)} \{e_l\})$ it holds that

$$\min\left\{l \in \{1, 2, \ldots, n\} \colon d(x, e_l) = \min_{u \in \{1,2,\ldots,n\}} d(x, e_u)\right\} \geq k. \tag{5.45}$$

Combining this with (5.43) yields that for all $k \in \{1, 2, \ldots, n\}$, $x \in P^{-1}(\{e_k\})$ with $e_k \notin (\cup_{l \in \mathbb{N} \cap [0,k)} \{e_l\})$ it holds that

$$\min\left\{l \in \{1, 2, \ldots, n\} \colon d(x, e_l) = \min_{u \in \{1,2,\ldots,n\}} d(x, e_u)\right\} = k. \tag{5.46}$$

Hence, we obtain that for all $k \in \{1, 2, \ldots, n\}$ with $e_k \notin (\cup_{l \in \mathbb{N} \cap [0,k)} \{e_l\})$ it holds that

$$P^{-1}(\{e_k\}) \subseteq \left\{x \in E \colon \min\left\{l \in \{1, 2, \ldots, n\} \colon d(x, e_l) = \min_{u \in \{1,2,\ldots,n\}} d(x, e_u)\right\} = k\right\}. \tag{5.47}$$

This and (5.38) show that for all $k \in \{1, 2, \ldots, n\}$ with $e_k \notin (\cup_{l \in \mathbb{N} \cap [0,k)} \{e_l\})$ it holds that

$$P^{-1}(\{e_k\}) = \left\{x \in E \colon \min\left\{l \in \{1, 2, \ldots, n\} \colon d(x, e_l) = \min_{u \in \{1,2,\ldots,n\}} d(x, e_u)\right\} = k\right\}. \tag{5.48}$$

Combining (5.40) with the fact that $D \colon E \to \mathbb{R}^n$ is $\mathcal{B}(E)/\mathcal{B}(\mathbb{R}^n)$-measurable therefore demonstrates that for all $k \in \{1, 2, \ldots, n\}$ with $e_k \notin (\cup_{l \in \mathbb{N} \cap [0,k)} \{e_l\})$ it holds that

$$\begin{aligned}
&P^{-1}(\{e_k\}) \\
&= \left\{x \in E \colon \min\left\{l \in \{1, 2, \ldots, n\} \colon d(x, e_l) = \min_{u \in \{1,2,\ldots,n\}} d(x, e_u)\right\} = k\right\} \\
&= \left\{x \in E \colon \min\left\{l \in \{1, 2, \ldots, n\} \colon D_l(x) = \min_{u \in \{1,2,\ldots,n\}} D_u(x)\right\} = k\right\} \\
&= \left\{x \in E \colon \left(\begin{array}{c} \forall\, l \in \mathbb{N} \cap [0, k) \colon D_k(x) < D_l(x) \text{ and} \\ \forall\, l \in \{1, 2, \ldots, n\} \colon D_k(x) \leq D_l(x) \end{array}\right)\right\} \\
&= \left[\bigcap_{l=1}^{k-1} \underbrace{\{x \in E \colon D_k(x) < D_l(x)\}}_{\in \mathcal{B}(E)}\right] \cap \left[\bigcap_{l=1}^{n} \underbrace{\{x \in E \colon D_k(x) \leq D_l(x)\}}_{\in \mathcal{B}(E)}\right] \in \mathcal{B}(E).
\end{aligned} \tag{5.49}$$

Hence, we obtain that for all $f \in \{e_1, e_2, \ldots, e_n\}$ it holds that

$$P^{-1}(\{f\}) \in \mathcal{B}(E). \tag{5.50}$$

Therefore, we obtain that for all $A \subseteq E$ it holds that

$$
\begin{aligned}
P^{-1}(A) &= P^{-1}(A \cap \{e_1, e_2, \ldots, e_n\}) \\
&= \cup_{f \in A \cap \{e_1, e_2, \ldots, e_n\}} \underbrace{P^{-1}(\{f\})}_{\in \mathcal{B}(E)} \in \mathcal{B}(E).
\end{aligned}
\tag{5.51}
$$

This establishes item (ii). The proof of Lemma 5.2.5 is thus complete. $\qquad\square$

**Lemma 5.2.6.** *Let $(E, d)$ be a separable metric space, let $(\mathcal{E}, \delta)$ be a metric space, let $(\Omega, \mathbb{F})$ be a measurable space, let $X \colon E \times \Omega \to \mathcal{E}$, assume for all $e \in E$ that $\Omega \ni \omega \mapsto X(e, \omega) \in \mathcal{E}$ is $\mathbb{F}/\mathcal{B}(\mathcal{E})$-measurable, and assume for all $\omega \in \Omega$ that $E \ni e \mapsto X(e, \omega) \in \mathcal{E}$ is continuous. Then $X \colon E \times \Omega \to \mathcal{E}$ is $(\mathcal{B}(E) \otimes \mathbb{F})/\mathcal{B}(\mathcal{E})$-measurable.*

*Proof of Lemma 5.2.6.* Throughout this proof let $e = (e_m)_{m \in \mathbb{N}} \colon \mathbb{N} \to E$ satisfy $\overline{\{e_m \colon m \in \mathbb{N}\}} = E$, let $P_n \colon E \to E$, $n \in \mathbb{N}$, satisfy for all $n \in \mathbb{N}$, $x \in E$ that

$$
P_n(x) = e_{\min\{k \in \{1,2,\ldots,n\} \colon d(x,e_k) = \min\{d(x,e_1), d(x,e_2), \ldots, d(x,e_n)\}\}},
\tag{5.52}
$$

and let $\mathcal{X}_n \colon E \times \Omega \to \mathcal{E}$, $n \in \mathbb{N}$, satisfy for all $n \in \mathbb{N}$, $x \in E$, $\omega \in \Omega$ that

$$
\mathcal{X}_n(x, \omega) = X(P_n(x), \omega).
\tag{5.53}
$$

Note that (5.53) shows that for all $n \in \mathbb{N}$, $B \in \mathcal{B}(\mathcal{E})$ it holds that

$$
\begin{aligned}
(\mathcal{X}_n)^{-1}(B) &= \{(x, \omega) \in E \times \Omega \colon \mathcal{X}_n(x, \omega) \in B\} \\
&= \bigcup_{y \in \mathrm{Im}(P_n)} \left( \left[ (\mathcal{X}_n)^{-1}(B) \right] \cap \left[ (P_n)^{-1}(\{y\}) \times \Omega \right] \right) \\
&= \bigcup_{y \in \mathrm{Im}(P_n)} \left\{ (x, \omega) \in E \times \Omega \colon \left[ \mathcal{X}_n(x, \omega) \in B \text{ and } x \in (P_n)^{-1}(\{y\}) \right] \right\} \\
&= \bigcup_{y \in \mathrm{Im}(P_n)} \left\{ (x, \omega) \in E \times \Omega \colon \left[ X(P_n(x), \omega) \in B \text{ and } x \in (P_n)^{-1}(\{y\}) \right] \right\}.
\end{aligned}
\tag{5.54}
$$

Item (ii) in Lemma 5.2.5 hence implies that for all $n \in \mathbb{N}$, $B \in \mathcal{B}(\mathcal{E})$ it holds that

$$
\begin{aligned}
(\mathcal{X}_n)^{-1}(B) &= \bigcup_{y \in \mathrm{Im}(P_n)} \left\{ (x, \omega) \in E \times \Omega \colon \left[ X(y, \omega) \in B \text{ and } x \in (P_n)^{-1}(\{y\}) \right] \right\} \\
&= \bigcup_{y \in \mathrm{Im}(P_n)} \left( \{(x, \omega) \in E \times \Omega \colon X(y, \omega) \in B\} \cap \left[ (P_n)^{-1}(\{y\}) \times \Omega \right] \right) \\
&= \bigcup_{y \in \mathrm{Im}(P_n)} \left( \underbrace{\left[ E \times \left( (X(y, \cdot))^{-1}(B) \right) \right]}_{\in (\mathcal{B}(E) \otimes \mathbb{F})} \cap \underbrace{\left[ (P_n)^{-1}(\{y\}) \times \Omega \right]}_{\in (\mathcal{B}(E) \otimes \mathbb{F})} \right) \in (\mathcal{B}(E) \otimes \mathbb{F}).
\end{aligned}
\tag{5.55}
$$

This proves that for all $n \in \mathbb{N}$ it holds that $\mathcal{X}_n$ is $(\mathcal{B}(E) \otimes \mathbb{F})/\mathcal{B}(\mathcal{E})$-measurable. In addition, note that item (i) in Lemma 5.2.5 and the assumption that for all $\omega \in \Omega$ it holds that $E \ni x \mapsto X(x, \omega) \in \mathcal{E}$ is continuous imply that for all $x \in E$, $\omega \in \Omega$ it holds that

$$
\lim_{n \to \infty} \mathcal{X}_n(x, \omega) = \lim_{n \to \infty} X(P_n(x), \omega) = X(x, \omega).
\tag{5.56}
$$

Combining this with the fact that for all $n \in \mathbb{N}$ it holds that $X_n \colon E \times \Omega \to \mathcal{E}$ is $(\mathcal{B}(E) \otimes \mathbb{F})/\mathcal{B}(\mathcal{E})$-measurable shows that $X \colon E \times \Omega \to \mathcal{E}$ is $(\mathcal{B}(E) \otimes \mathbb{F})/\mathcal{B}(\mathcal{E})$-measurable. The proof of Lemma 5.2.6 is thus completed. $\qquad\square$

## 5.2.3 Strong convergences rates for the optimization error

**Proposition 5.2.7.** *Let* $\mathbf{d}, K \in \mathbb{N}$, $L, \alpha \in \mathbb{R}$, $\beta \in (\alpha, \infty)$, *let* $(\Omega, \mathcal{F}, \mathbb{P})$ *be a probability space, let* $\mathscr{R} \colon [\alpha, \beta]^{\mathbf{d}} \times \Omega \to \mathbb{R}$ *be a random field, assume for all* $\theta, \vartheta \in [\alpha, \beta]^{\mathbf{d}}$, $\omega \in \Omega$ *that* $|\mathscr{R}(\theta, \omega) - \mathscr{R}(\vartheta, \omega)| \leq L\|\theta - \vartheta\|_{\infty}$, *let* $\Theta_k \colon \Omega \to [\alpha, \beta]^{\mathbf{d}}$, $k \in \{1, 2, \dots, K\}$, *be i.i.d. random variables, and assume that* $\Theta_1$ *is continuous uniformly distributed on* $[\alpha, \beta]^{\mathbf{d}}$ *(cf. Definition 3.1.16). Then*

(i) *it holds that* $\mathscr{R}$ *is a* $(\mathcal{B}([\alpha, \beta]^{\mathbf{d}}) \otimes \mathcal{F})/\mathcal{B}(\mathbb{R})$-*measurable function and*

(ii) *it holds for all* $\theta \in [\alpha, \beta]^{\mathbf{d}}$, $p \in (0, \infty)$ *that*

$$
\begin{aligned}
&\left(\mathbb{E}\big[\min_{k \in \{1,2,\dots,K\}} |\mathscr{R}(\Theta_k) - \mathscr{R}(\theta)|^p\big]\right)^{1/p} \\
&\leq \frac{L(\beta - \alpha)\max\{1, (p/\mathbf{d})^{1/\mathbf{d}}\}}{K^{1/\mathbf{d}}} \leq \frac{L(\beta - \alpha)\max\{1, p\}}{K^{1/\mathbf{d}}}.
\end{aligned} \tag{5.57}
$$

*Proof of Proposition 5.2.7.* Throughout this proof assume w.l.o.g. that $L > 0$, let $\delta \colon ([\alpha, \beta]^{\mathbf{d}}) \times ([\alpha, \beta]^{\mathbf{d}}) \to [0, \infty)$ satisfy for all $\theta, \vartheta \in [\alpha, \beta]^{\mathbf{d}}$ that $\delta(\theta, \vartheta) = \|\theta - \vartheta\|_{\infty}$, let $\mathbb{B} \colon (0, \infty)^2 \to (0, \infty)$ satisfy for all $x, y \in (0, \infty)$ that $\mathbb{B}(x, y) = \int_0^1 t^{x-1}(1-t)^{y-1} \, dt$, and let $\Theta_{1,1}, \Theta_{1,2}, \dots, \Theta_{1,\mathbf{d}} \colon \Omega \to [\alpha, \beta]$ satisfy $\Theta_1 = (\Theta_{1,1}, \Theta_{1,2}, \dots, \Theta_{1,\mathbf{d}})$. First of all, note that the assumption that $\forall \theta, \vartheta \in [\alpha, \beta]^{\mathbf{d}}$, $\omega \in \Omega \colon |\mathscr{R}(\theta, \omega) - \mathscr{R}(\vartheta, \omega)| \leq L\|\theta - \vartheta\|_{\infty}$ ensures that for all $\omega \in \Omega$ it holds that $[\alpha, \beta]^{\mathbf{d}} \ni \theta \mapsto \mathscr{R}(\theta, \omega) \in \mathbb{R}$ is continuous. Combining this with the fact that $([\alpha, \beta]^{\mathbf{d}}, \delta)$ is a separable metric space, the fact that for all $\theta \in [\alpha, \beta]^{\mathbf{d}}$ it holds that $\Omega \ni \omega \mapsto \mathscr{R}(\theta, \omega) \in \mathbb{R}$ is $\mathcal{F}/\mathcal{B}(\mathbb{R})$-measurable, and Lemma 5.2.6 proves item (i). Next observe that for all $\theta \in [\alpha, \beta]$, $\varepsilon \in [0, \infty)$ it holds that

$$
\begin{aligned}
\min\{\theta + \varepsilon, \beta\} - \max\{\theta - \varepsilon, \alpha\} &= \min\{\theta + \varepsilon, \beta\} + \min\{\varepsilon - \theta, -\alpha\} \\
&= \min\{\theta + \varepsilon + \min\{\varepsilon - \theta, -\alpha\}, \beta + \min\{\varepsilon - \theta, -\alpha\}\} \\
&= \min\{\min\{2\varepsilon, \theta - \alpha + \varepsilon\}, \min\{\beta - \theta + \varepsilon, \beta - \alpha\}\} \\
&\geq \min\{\min\{2\varepsilon, \alpha - \alpha + \varepsilon\}, \min\{\beta - \beta + \varepsilon, \beta - \alpha\}\} \\
&= \min\{2\varepsilon, \varepsilon, \varepsilon, \beta - \alpha\} = \min\{\varepsilon, \beta - \alpha\}.
\end{aligned} \tag{5.58}
$$

The assumption that $\Theta_1$ is continuous uniformly distributed on $[\alpha, \beta]^{\mathbf{d}}$ hence shows that for all $\theta = (\theta_1, \theta_2, \dots, \theta_{\mathbf{d}}) \in [\alpha, \beta]^{\mathbf{d}}$, $\varepsilon \in [0, \infty)$ it holds that

$$
\begin{aligned}
\mathbb{P}(\|\Theta_1 - \theta\|_{\infty} \leq \varepsilon) &= \mathbb{P}\big(\max_{i \in \{1,2,\dots,\mathbf{d}\}} |\Theta_{1,i} - \theta_i| \leq \varepsilon\big) \\
&= \mathbb{P}\big(\forall i \in \{1, 2, \dots, \mathbf{d}\} \colon -\varepsilon \leq \Theta_{1,i} - \theta_i \leq \varepsilon\big) \\
&= \mathbb{P}\big(\forall i \in \{1, 2, \dots, \mathbf{d}\} \colon \theta_i - \varepsilon \leq \Theta_{1,i} \leq \theta_i + \varepsilon\big) \\
&= \mathbb{P}\big(\forall i \in \{1, 2, \dots, \mathbf{d}\} \colon \max\{\theta_i - \varepsilon, \alpha\} \leq \Theta_{1,i} \leq \min\{\theta_i + \varepsilon, \beta\}\big) \\
&= \mathbb{P}\big(\Theta_1 \in \big[\times_{i=1}^{\mathbf{d}} [\max\{\theta_i - \varepsilon, \alpha\}, \min\{\theta_i + \varepsilon, \beta\}]\big]\big) \\
&= \tfrac{1}{(\beta - \alpha)^{\mathbf{d}}} \prod_{i=1}^{\mathbf{d}} (\min\{\theta_i + \varepsilon, \beta\} - \max\{\theta_i - \varepsilon, \alpha\}) \\
&\geq \tfrac{1}{(\beta - \alpha)^{\mathbf{d}}} [\min\{\varepsilon, \beta - \alpha\}]^{\mathbf{d}} = \min\Big\{1, \tfrac{\varepsilon^{\mathbf{d}}}{(\beta - \alpha)^{\mathbf{d}}}\Big\}.
\end{aligned} \tag{5.59}
$$

Therefore, we obtain for all $\theta \in [\alpha, \beta]^{\mathbf{d}}$, $p \in (0, \infty)$, $\varepsilon \in [0, \infty)$ that

$$
\begin{aligned}
\mathbb{P}(\|\Theta_1 - \theta\|_{\infty} > \varepsilon^{1/p}) &= 1 - \mathbb{P}(\|\Theta_1 - \theta\|_{\infty} \leq \varepsilon^{1/p}) \\
&\leq 1 - \min\Big\{1, \tfrac{\varepsilon^{\mathbf{d}/p}}{(\beta - \alpha)^{\mathbf{d}}}\Big\} = \max\Big\{0, 1 - \tfrac{\varepsilon^{\mathbf{d}/p}}{(\beta - \alpha)^{\mathbf{d}}}\Big\}.
\end{aligned} \tag{5.60}
$$

This, item (i), the assumption that $\forall\, \theta, \vartheta \in [\alpha, \beta]^{\mathbf{d}}$, $\omega \in \Omega\colon |\mathscr{R}(\theta, \omega) - \mathscr{R}(\vartheta, \omega)| \leq L\|\theta - \vartheta\|_\infty$, the assumption that $\Theta_k$, $k \in \{1, 2, \ldots, K\}$, are i.i.d. random variables, and Lemma 5.1.2 (applied with $(E, \delta) \curvearrowleft ([\alpha, \beta]^{\mathbf{d}}, \delta)$, $(X_k)_{k\in\{1,2,\ldots,K\}} \curvearrowleft (\Theta_k)_{k\in\{1,2,\ldots,K\}}$ in the notation of Lemma 5.1.2) establish that for all $\theta \in [\alpha, \beta]^{\mathbf{d}}$, $p \in (0, \infty)$ it holds that

$$\mathbb{E}\big[\min_{k\in\{1,2,\ldots,K\}}|\mathscr{R}(\Theta_k) - \mathscr{R}(\theta)|^p\big] \leq L^p \int_0^\infty [\mathbb{P}(\|\Theta_1 - \theta\|_\infty > \varepsilon^{1/p})]^K \, \mathrm{d}\varepsilon$$

$$\leq L^p \int_0^\infty \Big[\max\Big\{0, 1 - \tfrac{\varepsilon^{\mathbf{d}/p}}{(\beta-\alpha)^{\mathbf{d}}}\Big\}\Big]^K \mathrm{d}\varepsilon = L^p \int_0^{(\beta-\alpha)^p} \Big(1 - \tfrac{\varepsilon^{\mathbf{d}/p}}{(\beta-\alpha)^{\mathbf{d}}}\Big)^K \mathrm{d}\varepsilon \qquad (5.61)$$

$$= \tfrac{p}{\mathbf{d}} L^p (\beta - \alpha)^p \int_0^1 t^{p/\mathbf{d}-1}(1-t)^K \, \mathrm{d}t = \tfrac{p}{\mathbf{d}} L^p (\beta - \alpha)^p \int_0^1 t^{p/\mathbf{d}-1}(1-t)^{K+1-1} \, \mathrm{d}t$$

$$= \tfrac{p}{\mathbf{d}} L^p (\beta - \alpha)^p \, \mathbb{B}(p/\mathbf{d}, K + 1).$$

Corollary 5.2.4 (applied with $x \curvearrowleft p/\mathbf{d}$, $y \curvearrowleft K + 1$ for $p \in (0, \infty)$ in the notation of (5.34) in Corollary 5.2.4) hence demonstrates that for all $\theta \in [\alpha, \beta]^{\mathbf{d}}$, $p \in (0, \infty)$ it holds that

$$\mathbb{E}\big[\min_{k\in\{1,2,\ldots,K\}}|\mathscr{R}(\Theta_k) - \mathscr{R}(\theta)|^p\big]$$

$$\leq \frac{\tfrac{p}{\mathbf{d}} L^p (\beta - \alpha)^p \max\{1, (p/\mathbf{d})^{p/\mathbf{d}}\}}{\tfrac{p}{\mathbf{d}}(K + 1 + \min\{p/\mathbf{d} - 1, 0\})^{p/\mathbf{d}}} \leq \frac{L^p (\beta - \alpha)^p \max\{1, (p/\mathbf{d})^{p/\mathbf{d}}\}}{K^{p/\mathbf{d}}}. \qquad (5.62)$$

This implies for all $\theta \in [\alpha, \beta]^{\mathbf{d}}$, $p \in (0, \infty)$ that

$$\big(\mathbb{E}\big[\min_{k\in\{1,2,\ldots,K\}}|\mathscr{R}(\Theta_k) - \mathscr{R}(\theta)|^p\big]\big)^{1/p}$$

$$\leq \frac{L(\beta - \alpha)\max\{1, (p/\mathbf{d})^{1/\mathbf{d}}\}}{K^{1/\mathbf{d}}} \leq \frac{L(\beta - \alpha)\max\{1, p\}}{K^{1/\mathbf{d}}}. \qquad (5.63)$$

This shows item (ii) and thus completes the proof of Proposition 5.2.7. □

# 5.3   Strong convergences rates for the optimization error involving ANNs

## 5.3.1   Local Lipschitz continuity estimates for the parametrization functions associated to ANNs

**Lemma 5.3.1.** *Let* $a, x, y \in \mathbb{R}$. *Then*

$$|\max\{x, a\} - \max\{y, a\}| \leq \max\{x, y\} - \min\{x, y\} = |x - y|. \qquad (5.64)$$

*Proof of Lemma 5.3.1.* Observe that

$$|\max\{x, a\} - \max\{y, a\}| = |\max\{\max\{x, y\}, a\} - \max\{\min\{x, y\}, a\}|$$

$$= \max\{\max\{x, y\}, a\} - \max\{\min\{x, y\}, a\}$$

$$= \max\Big\{\max\{x, y\} - \max\{\min\{x, y\}, a\}, a - \max\{\min\{x, y\}, a\}\Big\}$$

$$\leq \max\Big\{\max\{x, y\} - \max\{\min\{x, y\}, a\}, a - a\Big\} \qquad (5.65)$$

$$= \max\Big\{\max\{x, y\} - \max\{\min\{x, y\}, a\}, 0\Big\} \leq \max\Big\{\max\{x, y\} - \min\{x, y\}, 0\Big\}$$

$$= \max\{x, y\} - \min\{x, y\} = |\max\{x, y\} - \min\{x, y\}| = |x - y|.$$

The proof of Lemma 5.3.1 is thus complete. □

**Corollary 5.3.2.** *Let $a, x, y \in \mathbb{R}$. Then*

$$|\min\{x, a\} - \min\{y, a\}| \le \max\{x, y\} - \min\{x, y\} = |x - y|. \tag{5.66}$$

*Proof of Corollary 5.3.2.* Note that Lemma 5.3.1 ensures that

$$|\min\{x, a\} - \min\{y, a\}| = |-(\min\{x, a\} - \min\{y, a\})| = |\max\{-x, -a\} - \max\{-y, -a\}|$$
$$\le |(-x) - (-y)| = |x - y|. \tag{5.67}$$

The proof of Corollary 5.3.2 is thus complete. $\qquad\square$

**Lemma 5.3.3.** *Let $d \in \mathbb{N}$. Then it holds for all $x, y \in \mathbb{R}^d$ that*

$$\|\mathfrak{R}_d(x) - \mathfrak{R}_d(y)\|_\infty \le \|x - y\|_\infty \tag{5.68}$$

*(cf. Definitions 2.1.7 and 3.1.16).*

*Proof of Lemma 5.3.3.* Note that Lemma 5.3.1 establishes (5.68). The proof of Lemma 5.3.3 is thus complete. $\qquad\square$

**Lemma 5.3.4.** *Let $d \in \mathbb{N}$, $u \in [-\infty, \infty)$, $v \in (u, \infty]$. Then it holds for all $x, y \in \mathbb{R}^d$ that*

$$\|\mathfrak{C}_{u,v,d}(x) - \mathfrak{C}_{u,v,d}(y)\|_\infty \le \|x - y\|_\infty \tag{5.69}$$

*(cf. Definitions 2.1.12 and 3.1.16).*

*Proof of Lemma 5.3.4.* Note that Lemma 5.3.1, Corollary 5.3.2, and the fact that for all $x \in \mathbb{R}$ it holds that $\max\{-\infty, x\} = x = \min\{x, \infty\}$ show that for all $x, y \in \mathbb{R}$ it holds that

$$|\mathfrak{c}_{u,v}(x) - \mathfrak{c}_{u,v}(y)| = |\max\{u, \min\{x, v\}\} - \max\{u, \min\{y, v\}\}| \le |\min\{x, v\} - \min\{y, v\}| \le |x - y| \tag{5.70}$$

(cf. Definition 2.1.11). Hence, we obtain that for all $x = (x_1, x_2, \ldots, x_d), y = (y_1, y_2, \ldots, y_d) \in \mathbb{R}^d$ it holds that

$$\|\mathfrak{C}_{u,v,d}(x) - \mathfrak{C}_{u,v,d}(y)\|_\infty = \max_{i \in \{1,2,\ldots,d\}} |\mathfrak{c}_{u,v}(x_i) - \mathfrak{c}_{u,v}(y_i)| \le \max_{i \in \{1,2,\ldots,d\}} |x_i - y_i| = \|x - y\|_\infty \tag{5.71}$$

(cf. Definitions 2.1.12 and 3.1.16). The proof of Lemma 5.3.4 is thus complete. $\qquad\square$

**Lemma 5.3.5** (Row sum norm, operator norm induced by the maximum norm)**.** *Let $a, b \in \mathbb{N}$, $M = (M_{i,j})_{(i,j) \in \{1,2,\ldots,a\} \times \{1,2,\ldots,b\}} \in \mathbb{R}^{a \times b}$. Then*

$$\sup_{v \in \mathbb{R}^b \setminus \{0\}} \left[ \frac{\|Mv\|_\infty}{\|v\|_\infty} \right] = \max_{i \in \{1,2,\ldots,a\}} \left[ \sum_{j=1}^b |M_{i,j}| \right] \le b \left[ \max_{i \in \{1,2,\ldots,a\}} \max_{j \in \{1,2,\ldots,b\}} |M_{i,j}| \right] \tag{5.72}$$

*(cf. Definition 3.1.16).*

*Proof of Lemma 5.3.5.* Observe that

$$
\sup_{v \in \mathbb{R}^b} \left[ \frac{\|Mv\|_\infty}{\|v\|_\infty} \right] = \sup_{v \in \mathbb{R}^b, \|v\|_\infty \le 1} \|Mv\|_\infty = \sup_{v=(v_1,v_2,\ldots,v_b) \in [-1,1]^b} \|Mv\|_\infty
$$

$$
= \sup_{v=(v_1,v_2,\ldots,v_b) \in [-1,1]^b} \left( \max_{i \in \{1,2,\ldots,a\}} \left| \sum_{j=1}^b M_{i,j} v_j \right| \right)
$$

$$
= \max_{i \in \{1,2,\ldots,a\}} \left( \sup_{v=(v_1,v_2,\ldots,v_b) \in [-1,1]^b} \left| \sum_{j=1}^b M_{i,j} v_j \right| \right) = \max_{i \in \{1,2,\ldots,a\}} \left( \sum_{j=1}^b |M_{i,j}| \right)
\tag{5.73}
$$

(cf. Definition 3.1.16). The proof of Lemma 5.3.5 is thus complete. $\square$

**Theorem 5.3.6.** *Let $a \in \mathbb{R}$, $b \in [a, \infty)$, $d, L \in \mathbb{N}$, $l = (l_0, l_1, \ldots, l_L) \in \mathbb{N}^{L+1}$ satisfy*

$$
d \ge \sum_{k=1}^L l_k(l_{k-1} + 1).
\tag{5.74}
$$

*Then it holds for all $\theta, \vartheta \in \mathbb{R}^d$ that*

$$
\sup_{x \in [a,b]^{l_0}} \|\mathcal{N}_{-\infty,\infty}^{\theta,l}(x) - \mathcal{N}_{-\infty,\infty}^{\vartheta,l}(x)\|_\infty
$$

$$
\le \max\{1, |a|, |b|\} \|\theta - \vartheta\|_\infty \left[ \prod_{m=0}^{L-1} (l_m + 1) \right] \left[ \sum_{n=0}^{L-1} \left( \max\{1, \|\theta\|_\infty^n\} \|\vartheta\|_\infty^{L-1-n} \right) \right]
\tag{5.75}
$$

$$
\le L \max\{1, |a|, |b|\} (\max\{1, \|\theta\|_\infty, \|\vartheta\|_\infty\})^{L-1} \left[ \prod_{m=0}^{L-1} (l_m + 1) \right] \|\theta - \vartheta\|_\infty
$$

$$
\le L \max\{1, |a|, |b|\} \left( \|l\|_\infty + 1 \right)^L (\max\{1, \|\theta\|_\infty, \|\vartheta\|_\infty\})^{L-1} \|\theta - \vartheta\|_\infty
$$

*(cf. Definition 2.1.27 and Definition 3.1.16).*

*Proof of Theorem 5.3.6.* Throughout this proof let $\theta_j = (\theta_{j,1}, \theta_{j,2}, \ldots, \theta_{j,d}) \in \mathbb{R}^d$, $j \in \{1,2\}$, let $\mathfrak{d} \in \mathbb{N}$ satisfy

$$
\mathfrak{d} = \sum_{k=1}^L l_k(l_{k-1} + 1),
\tag{5.76}
$$

let $W_{j,k} \in \mathbb{R}^{l_k \times l_{k-1}}$, $k \in \{1,2,\ldots,L\}$, $j \in \{1,2\}$, and $B_{j,k} \in \mathbb{R}^{l_k}$, $k \in \{1,2,\ldots,L\}$, $j \in \{1,2\}$, satisfy for all $j \in \{1,2\}$, $k \in \{1,2,\ldots,L\}$ that

$$
\mathcal{T}\big(((W_{j,1}, B_{j,1}), (W_{j,2}, B_{j,2}), \ldots, (W_{j,L}, B_{j,L}))\big) = (\theta_{j,1}, \theta_{j,2}, \ldots, \theta_{j,\mathfrak{d}}),
\tag{5.77}
$$

let $\phi_{j,k} \in \mathbf{N}$, $k \in \{1,2,\ldots,L\}$, $j \in \{1,2\}$, satisfy for all $j \in \{1,2\}$, $k \in \{1,2,\ldots,L\}$ that

$$
\phi_{j,k} = \big((W_{j,1}, B_{j,1}), (W_{j,2}, B_{j,2}), \ldots, (W_{j,k}, B_{j,k})\big) \in \left[ \bigtimes_{i=1}^k \big( \mathbb{R}^{l_i \times l_{i-1}} \times \mathbb{R}^{l_i} \big) \right],
\tag{5.78}
$$

let $D = [a,b]^{l_0}$, let $\mathfrak{m}_{j,k} \in [0, \infty)$, $j \in \{1,2\}$, $k \in \{0,1,\ldots,L\}$, satisfy for all $j \in \{1,2\}$, $k \in \{0,1,\ldots,L\}$ that

$$
\mathfrak{m}_{j,k} = \begin{cases} \max\{1, |a|, |b|\} & : k = 0 \\ \max\{1, \sup_{x \in D} \|(\mathcal{R}_{\mathfrak{r}}(\phi_{j,k}))(x)\|_\infty\} & : k > 0, \end{cases}
\tag{5.79}
$$

and let $\mathfrak{e}_k \in [0, \infty)$, $k \in \{0, 1, \ldots, L\}$, satisfy for all $k \in \{0, 1, \ldots, L\}$ that

$$\mathfrak{e}_k = \begin{cases} 0 & : k = 0 \\ \sup_{x \in D} \|(\mathcal{R}_{\mathfrak{r}}(\phi_{1,k}))(x) - (\mathcal{R}_{\mathfrak{r}}(\phi_{2,k}))(x)\|_\infty & : k > 0 \end{cases} \tag{5.80}$$

(cf. Definitions 2.1.6, 2.2.3, 2.2.36, and 3.1.16). Note that Lemma 5.3.5 demonstrates that

$$
\begin{aligned}
\mathfrak{e}_1 &= \sup_{x \in D} \|(\mathcal{R}_{\mathfrak{r}}(\phi_{1,1}))(x) - (\mathcal{R}_{\mathfrak{r}}(\phi_{2,1}))(x)\|_\infty = \sup_{x \in D} \|(W_{1,1}x + B_{1,1}) - (W_{2,1}x + B_{2,1})\|_\infty \\
&\leq \left[ \sup_{x \in D} \|(W_{1,1} - W_{2,1})x\|_\infty \right] + \|B_{1,1} - B_{2,1}\|_\infty \\
&\leq \left[ \sup_{v \in \mathbb{R}^{l_0} \setminus \{0\}} \left( \frac{\|(W_{1,1} - W_{2,1})v\|_\infty}{\|v\|_\infty} \right) \right] \left[ \sup_{x \in D} \|x\|_\infty \right] + \|B_{1,1} - B_{2,1}\|_\infty \\
&\leq l_0 \|\theta_1 - \theta_2\|_\infty \max\{|a|, |b|\} + \|B_{1,1} - B_{2,1}\|_\infty \leq l_0 \|\theta_1 - \theta_2\|_\infty \max\{|a|, |b|\} + \|\theta_1 - \theta_2\|_\infty \\
&= \|\theta_1 - \theta_2\|_\infty (l_0 \max\{|a|, |b|\} + 1) \leq \mathfrak{m}_{1,0} \|\theta_1 - \theta_2\|_\infty (l_0 + 1).
\end{aligned}
\tag{5.81}
$$

Moreover, observe that the triangle inequality assures that for all $k \in \{1, 2, \ldots, L\} \cap (1, \infty)$ it holds that

$$
\begin{aligned}
\mathfrak{e}_k &= \sup_{x \in D} \|(\mathcal{R}_{\mathfrak{r}}(\phi_{1,k}))(x) - (\mathcal{R}_{\mathfrak{r}}(\phi_{2,k}))(x)\|_\infty \\
&= \sup_{x \in D} \left\| \left[ W_{1,k}\Big(\mathfrak{R}_{l_{k-1}}\big((\mathcal{R}_{\mathfrak{r}}(\phi_{1,k-1}))(x)\big)\Big) + B_{1,k} \right] - \left[ W_{2,k}\Big(\mathfrak{R}_{l_{k-1}}\big((\mathcal{R}_{\mathfrak{r}}(\phi_{2,k-1}))(x)\big)\Big) + B_{2,k} \right] \right\|_\infty \\
&\leq \left[ \sup_{x \in D} \left\| W_{1,k}\Big(\mathfrak{R}_{l_{k-1}}\big((\mathcal{R}_{\mathfrak{r}}(\phi_{1,k-1}))(x)\big)\Big) - W_{2,k}\Big(\mathfrak{R}_{l_{k-1}}\big((\mathcal{R}_{\mathfrak{r}}(\phi_{2,k-1}))(x)\big)\Big) \right\|_\infty \right] + \|\theta_1 - \theta_2\|_\infty.
\end{aligned}
\tag{5.82}
$$

The triangle inequality hence implies that for all $j \in \{1, 2\}$, $k \in \{1, 2, \ldots, L\} \cap (1, \infty)$ it holds that

$$
\begin{aligned}
\mathfrak{e}_k &\leq \left[ \sup_{x \in D} \left\| (W_{1,k} - W_{2,k})\big(\mathfrak{R}_{l_{k-1}}\big((\mathcal{R}_{\mathfrak{r}}(\phi_{j,k-1}))(x)\big)\big) \right\|_\infty \right] \\
&\quad + \left[ \sup_{x \in D} \left\| W_{3-j,k}\Big(\mathfrak{R}_{l_{k-1}}\big((\mathcal{R}_{\mathfrak{r}}(\phi_{1,k-1}))(x)\big) - \mathfrak{R}_{l_{k-1}}\big((\mathcal{R}_{\mathfrak{r}}(\phi_{2,k-1}))(x)\big)\Big) \right\|_\infty \right] + \|\theta_1 - \theta_2\|_\infty \\
&\leq \left[ \sup_{v \in \mathbb{R}^{l_{k-1}} \setminus \{0\}} \left( \frac{\|(W_{1,k} - W_{2,k})v\|_\infty}{\|v\|_\infty} \right) \right] \left[ \sup_{x \in D} \left\| \mathfrak{R}_{l_{k-1}}\big((\mathcal{R}_{\mathfrak{r}}(\phi_{j,k-1}))(x)\big) \right\|_\infty \right] + \|\theta_1 - \theta_2\|_\infty \\
&\quad + \left[ \sup_{v \in \mathbb{R}^{l_{k-1}} \setminus \{0\}} \left( \frac{\|W_{3-j,k}v\|_\infty}{\|v\|_\infty} \right) \right] \left[ \sup_{x \in D} \left\| \mathfrak{R}_{l_{k-1}}\big((\mathcal{R}_{\mathfrak{r}}(\phi_{1,k-1}))(x)\big) - \mathfrak{R}_{l_{k-1}}\big((\mathcal{R}_{\mathfrak{r}}(\phi_{2,k-1}))(x)\big) \right\|_\infty \right].
\end{aligned}
\tag{5.83}
$$

Lemma 5.3.5 and Lemma 5.3.3 therefore show that for all $j \in \{1, 2\}$, $k \in \{1, 2, \ldots, L\} \cap$

$(1, \infty)$ it holds that

$$
\begin{aligned}
\mathfrak{e}_k &\leq l_{k-1} \|\theta_1 - \theta_2\|_\infty \left[ \sup_{x \in D} \left\| \mathfrak{R}_{l_{k-1}} \big( (\mathcal{R}_\mathfrak{r}(\phi_{j,k-1}))(x) \big) \right\|_\infty \right] + \|\theta_1 - \theta_2\|_\infty \\
&\quad + l_{k-1} \|\theta_{3-j}\|_\infty \left[ \sup_{x \in D} \left\| \mathfrak{R}_{l_{k-1}} \big( (\mathcal{R}_\mathfrak{r}(\phi_{1,k-1}))(x) \big) - \mathfrak{R}_{l_{k-1}} \big( (\mathcal{R}_\mathfrak{r}(\phi_{2,k-1}))(x) \big) \right\|_\infty \right] \\
&\leq l_{k-1} \|\theta_1 - \theta_2\|_\infty \left[ \sup_{x \in D} \left\| (\mathcal{R}_\mathfrak{r}(\phi_{j,k-1}))(x) \right\|_\infty \right] + \|\theta_1 - \theta_2\|_\infty \\
&\quad + l_{k-1} \|\theta_{3-j}\|_\infty \left[ \sup_{x \in D} \left\| (\mathcal{R}_\mathfrak{r}(\phi_{1,k-1}))(x) - (\mathcal{R}_\mathfrak{r}(\phi_{2,k-1}))(x) \right\|_\infty \right] \\
&\leq \|\theta_1 - \theta_2\|_\infty (l_{k-1} \mathfrak{m}_{j,k-1} + 1) + l_{k-1} \|\theta_{3-j}\|_\infty \mathfrak{e}_{k-1}.
\end{aligned}
\tag{5.84}
$$

Hence, we obtain that for all $j \in \{1, 2\}$, $k \in \{1, 2, \ldots, L\} \cap (1, \infty)$ it holds that

$$
\mathfrak{e}_k \leq \mathfrak{m}_{j,k-1} \|\theta_1 - \theta_2\|_\infty (l_{k-1} + 1) + l_{k-1} \|\theta_{3-j}\|_\infty \mathfrak{e}_{k-1}.
\tag{5.85}
$$

Combining this with (5.81), the fact that $\mathfrak{e}_0 = 0$, and the fact that $\mathfrak{m}_{1,0} = \mathfrak{m}_{2,0}$ demonstrates that for all $j \in \{1, 2\}$, $k \in \{1, 2, \ldots, L\}$ it holds that

$$
\mathfrak{e}_k \leq \mathfrak{m}_{j,k-1}(l_{k-1} + 1) \|\theta_1 - \theta_2\|_\infty + l_{k-1} \|\theta_{3-j}\|_\infty \mathfrak{e}_{k-1}.
\tag{5.86}
$$

This shows that for all $j = (j_n)_{n \in \{0,1,\ldots,L\}} \colon \{0, 1, \ldots, L\} \to \{1, 2\}$ and all $k \in \{1, 2, \ldots, L\}$ it holds that

$$
\mathfrak{e}_k \leq \mathfrak{m}_{j_{k-1}, k-1}(l_{k-1} + 1) \|\theta_1 - \theta_2\|_\infty + l_{k-1} \|\theta_{3-j_{k-1}}\|_\infty \mathfrak{e}_{k-1}.
\tag{5.87}
$$

Therefore, we obtain that for all $j = (j_n)_{n \in \{0,1,\ldots,L\}} \colon \{0, 1, \ldots, L\} \to \{1, 2\}$ and all $k \in \{1, 2, \ldots, L\}$ it holds that

$$
\begin{aligned}
\mathfrak{e}_k &\leq \sum_{n=0}^{k-1} \left( \left[ \prod_{m=n+1}^{k-1} \big( l_m \|\theta_{3-j_m}\|_\infty \big) \right] \mathfrak{m}_{j_n, n}(l_n + 1) \|\theta_1 - \theta_2\|_\infty \right) \\
&= \|\theta_1 - \theta_2\|_\infty \left[ \sum_{n=0}^{k-1} \left( \left[ \prod_{m=n+1}^{k-1} \big( l_m \|\theta_{3-j_m}\|_\infty \big) \right] \mathfrak{m}_{j_n, n}(l_n + 1) \right) \right].
\end{aligned}
\tag{5.88}
$$

Next observe that Lemma 5.3.5 ensures that for all $j \in \{1, 2\}$, $k \in \{1, 2, \ldots, L\} \cap (1, \infty)$, $x \in D$ it holds that

$$
\begin{aligned}
\|(\mathcal{R}_\mathfrak{r}(\phi_{j,k}))(x)\|_\infty &= \left\| W_{j,k} \Big( \mathfrak{R}_{l_{k-1}} \big( (\mathcal{R}_\mathfrak{r}(\phi_{j,k-1}))(x) \big) \Big) + B_{j,k} \right\|_\infty \\
&\leq \left[ \sup_{v \in \mathbb{R}^{l_{k-1}} \setminus \{0\}} \frac{\|W_{j,k} v\|_\infty}{\|v\|_\infty} \right] \left\| \mathfrak{R}_{l_{k-1}} \big( (\mathcal{R}_\mathfrak{r}(\phi_{j,k-1}))(x) \big) \right\|_\infty + \|B_{j,k}\|_\infty \\
&\leq l_{k-1} \|\theta_j\|_\infty \left\| \mathfrak{R}_{l_{k-1}} \big( (\mathcal{R}_\mathfrak{r}(\phi_{j,k-1}))(x) \big) \right\|_\infty + \|\theta_j\|_\infty \\
&\leq l_{k-1} \|\theta_j\|_\infty \|(\mathcal{R}_\mathfrak{r}(\phi_{j,k-1}))(x)\|_\infty + \|\theta_j\|_\infty \\
&= \big( l_{k-1} \|(\mathcal{R}_\mathfrak{r}(\phi_{j,k-1}))(x)\|_\infty + 1 \big) \|\theta_j\|_\infty \\
&\leq (l_{k-1} \mathfrak{m}_{j,k-1} + 1) \|\theta_j\|_\infty \leq \mathfrak{m}_{j,k-1}(l_{k-1} + 1) \|\theta_j\|_\infty.
\end{aligned}
\tag{5.89}
$$

Hence, we obtain for all $j \in \{1, 2\}$, $k \in \{1, 2, \ldots, L\} \cap (1, \infty)$ that

$$\mathfrak{m}_{j,k} \leq \max\{1, \mathfrak{m}_{j,k-1}(l_{k-1}+1)\|\theta_j\|_\infty\}. \tag{5.90}$$

Furthermore, note that Lemma 5.3.5 assures that for all $j \in \{1, 2\}$, $x \in D$ it holds that

$$
\begin{aligned}
\|(\mathcal{R}_{\mathfrak{r}}(\phi_{j,1}))(x)\|_\infty &= \|W_{j,1}x + B_{j,1}\|_\infty \\
&\leq \left[\sup_{v \in \mathbb{R}^{l_0} \setminus \{0\}} \frac{\|W_{j,1}v\|_\infty}{\|v\|_\infty}\right]\|x\|_\infty + \|B_{j,1}\|_\infty \\
&\leq l_0 \|\theta_j\|_\infty \|x\|_\infty + \|\theta_j\|_\infty \leq l_0 \|\theta_j\|_\infty \max\{|a|, |b|\} + \|\theta_j\|_\infty \\
&= (l_0 \max\{|a|, |b|\} + 1)\|\theta_j\|_\infty \leq \mathfrak{m}_{1,0}(l_0+1)\|\theta_j\|_\infty.
\end{aligned}
\tag{5.91}
$$

Therefore, we obtain that for all $j \in \{1, 2\}$ it holds that

$$\mathfrak{m}_{j,1} \leq \max\{1, \mathfrak{m}_{j,0}(l_0+1)\|\theta_j\|_\infty\}. \tag{5.92}$$

Combining this with (5.90) demonstrates that for all $j \in \{1, 2\}$, $k \in \{1, 2, \ldots, L\}$ it holds that

$$\mathfrak{m}_{j,k} \leq \max\{1, \mathfrak{m}_{j,k-1}(l_{k-1}+1)\|\theta_j\|_\infty\}. \tag{5.93}$$

Hence, we obtain that for all $j \in \{1, 2\}$, $k \in \{0, 1, \ldots, L\}$ it holds that

$$\mathfrak{m}_{j,k} \leq \mathfrak{m}_{j,0}\left[\prod_{n=0}^{k-1}(l_n+1)\right]\left[\max\{1, \|\theta_j\|_\infty\}\right]^k. \tag{5.94}$$

Combining this with (5.88) proves that for all $j = (j_n)_{n \in \{0,1,\ldots,L\}} \colon \{0, 1, \ldots, L\} \to \{1, 2\}$ and all $k \in \{1, 2, \ldots, L\}$ it holds that

$$
\begin{aligned}
\mathfrak{e}_k &\leq \|\theta_1 - \theta_2\|_\infty \left[\sum_{n=0}^{k-1}\left(\left[\prod_{m=n+1}^{k-1}(l_m\|\theta_{3-j_m}\|_\infty)\right]\left(\mathfrak{m}_{j_n,0}\left[\prod_{v=0}^{n-1}(l_v+1)\right]\max\{1, \|\theta_{j_n}\|_\infty^n\}(l_n+1)\right)\right)\right] \\
&= \mathfrak{m}_{1,0}\|\theta_1 - \theta_2\|_\infty\left[\sum_{n=0}^{k-1}\left(\left[\prod_{m=n+1}^{k-1}(l_m\|\theta_{3-j_m}\|_\infty)\right]\left(\left[\prod_{v=0}^{n}(l_v+1)\right]\max\{1, \|\theta_{j_n}\|_\infty^n\}\right)\right)\right] \\
&\leq \mathfrak{m}_{1,0}\|\theta_1 - \theta_2\|_\infty\left[\sum_{n=0}^{k-1}\left(\left[\prod_{m=n+1}^{k-1}\|\theta_{3-j_m}\|_\infty\right]\left[\prod_{v=0}^{k-1}(l_v+1)\right]\max\{1, \|\theta_{j_n}\|_\infty^n\}\right)\right] \\
&= \mathfrak{m}_{1,0}\|\theta_1 - \theta_2\|_\infty\left[\prod_{n=0}^{k-1}(l_n+1)\right]\left[\sum_{n=0}^{k-1}\left(\left[\prod_{m=n+1}^{k-1}\|\theta_{3-j_m}\|_\infty\right]\max\{1, \|\theta_{j_n}\|_\infty^n\}\right)\right].
\end{aligned}
\tag{5.95}
$$

Therefore, we obtain that for all $j \in \{1, 2\}$, $k \in \{1, 2, \ldots, L\}$ it holds that

$$
\begin{aligned}
\mathfrak{e}_k &\leq \mathfrak{m}_{1,0}\|\theta_1 - \theta_2\|_\infty\left[\prod_{n=0}^{k-1}(l_n+1)\right]\left[\sum_{n=0}^{k-1}\left(\left[\prod_{m=n+1}^{k-1}\|\theta_{3-j}\|_\infty\right]\max\{1, \|\theta_j\|_\infty^n\}\right)\right] \\
&= \mathfrak{m}_{1,0}\|\theta_1 - \theta_2\|_\infty\left[\prod_{n=0}^{k-1}(l_n+1)\right]\left[\sum_{n=0}^{k-1}(\max\{1, \|\theta_j\|_\infty^n\}\|\theta_{3-j}\|_\infty^{k-1-n})\right] \\
&\leq k\,\mathfrak{m}_{1,0}\|\theta_1 - \theta_2\|_\infty(\max\{1, \|\theta_1\|_\infty, \|\theta_2\|_\infty\})^{k-1}\left[\prod_{m=0}^{k-1}(l_m+1)\right].
\end{aligned}
\tag{5.96}
$$

The proof of Theorem 5.3.6 is thus complete. $\qquad \square$

**Corollary 5.3.7.** *Let $a \in \mathbb{R}$, $b \in [a, \infty)$, $u \in [-\infty, \infty)$, $v \in (u, \infty]$, $d, L \in \mathbb{N}$, $l = (l_0, l_1, \ldots, l_L) \in \mathbb{N}^{L+1}$ satisfy*

$$d \geq \sum_{k=1}^{L} l_k(l_{k-1} + 1). \tag{5.97}$$

*Then it holds for all $\theta, \vartheta \in \mathbb{R}^d$ that*

$$\sup_{x \in [a,b]^{l_0}} \|\mathcal{N}_{u,v}^{\theta,l}(x) - \mathcal{N}_{u,v}^{\vartheta,l}(x)\|_\infty \leq L \max\{1, |a|, |b|\} (\|l\|_\infty + 1)^L (\max\{1, \|\theta\|_\infty, \|\vartheta\|_\infty\})^{L-1} \|\theta - \vartheta\|_\infty \tag{5.98}$$

*(cf. Definitions 2.1.27 and 3.1.16).*

*Proof of Corollary 5.3.7.* Note that Lemma 5.3.4 and Theorem 5.3.6 demonstrate that for all $\theta, \vartheta \in \mathbb{R}^d$ it holds that

$$
\begin{aligned}
&\sup_{x \in [a,b]^{l_0}} \|\mathcal{N}_{u,v}^{\theta,l}(x) - \mathcal{N}_{u,v}^{\vartheta,l}(x)\|_\infty \\
&= \sup_{x \in [a,b]^{l_0}} \|\mathfrak{C}_{u,v,l_L}(\mathcal{N}_{-\infty,\infty}^{\theta,l}(x)) - \mathfrak{C}_{u,v,l_L}(\mathcal{N}_{-\infty,\infty}^{\vartheta,l}(x))\|_\infty \\
&\leq \sup_{x \in [a,b]^{l_0}} \|\mathcal{N}_{-\infty,\infty}^{\theta,l}(x) - \mathcal{N}_{-\infty,\infty}^{\vartheta,l}(x)\|_\infty \\
&\leq L \max\{1, |a|, |b|\} (\|l\|_\infty + 1)^L (\max\{1, \|\theta\|_\infty, \|\vartheta\|_\infty\})^{L-1} \|\theta - \vartheta\|_\infty
\end{aligned}
\tag{5.99}
$$

(cf. Definitions 2.1.12, 2.1.27, and 3.1.16). This completes the proof of Corollary 5.3.7. □

## 5.3.2 Strong convergences rates for the optimization error involving ANNs

**Lemma 5.3.8.** *Let $d, \mathbf{d}, \mathbf{L}, M \in \mathbb{N}$, $B, b \in [1, \infty)$, $u \in \mathbb{R}$, $v \in (u, \infty)$, $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \ldots, \mathbf{l_L}) \in \mathbb{N}^{\mathbf{L}+1}$, $D \subseteq [-b, b]^d$, assume $\mathbf{l}_0 = d$, $\mathbf{l_L} = 1$, and $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1} + 1)$, let $\Omega$ be a set, let $X_j : \Omega \to D$, $j \in \{1, 2, \ldots, M\}$, and $Y_j : \Omega \to [u, v]$, $j \in \{1, 2, \ldots, M\}$, be functions, and let $\mathscr{R} : [-B, B]^{\mathbf{d}} \times \Omega \to [0, \infty)$ satisfy for all $\theta \in [-B, B]^{\mathbf{d}}$, $\omega \in \Omega$ that*

$$\mathscr{R}(\theta, \omega) = \frac{1}{M} \left[ \sum_{j=1}^{M} |\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)|^2 \right] \tag{5.100}$$

*(cf. Definition 2.1.27). Then it holds for all $\theta, \vartheta \in [-B, B]^{\mathbf{d}}$, $\omega \in \Omega$ that*

$$|\mathscr{R}(\theta, \omega) - \mathscr{R}(\vartheta, \omega)| \leq 2(v - u)b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}} B^{\mathbf{L}-1} \|\theta - \vartheta\|_\infty \tag{5.101}$$

*(cf. Definition 3.1.16).*

*Proof of Lemma 5.3.8.* Observe that the fact that $\forall x_1, x_2, y \in \mathbb{R}: (x_1 - y)^2 - (x_2 - y)^2 = (x_1 - x_2)((x_1 - y) + (x_2 - y))$, the fact that $\forall \theta \in \mathbb{R}^{\mathbf{d}}, x \in \mathbb{R}^d: \mathcal{N}_{u,v}^{\theta,\mathbf{l}}(x) \in [u, v]$, and the assumption that $\forall j \in \{1, 2, \ldots, M\}, \omega \in \Omega: Y_j(\omega) \in [u, v]$ prove that for all

$\theta, \vartheta \in [-B, B]^{\mathbf{d}}$, $\omega \in \Omega$ it holds that

$$|\mathscr{R}(\theta, \omega) - \mathscr{R}(\vartheta, \omega)|$$

$$= \frac{1}{M} \left| \left[ \sum_{j=1}^{M} |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)|^2 \right] - \left[ \sum_{j=1}^{M} |\mathscr{N}_{u,v}^{\vartheta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)|^2 \right] \right|$$

$$\leq \frac{1}{M} \left[ \sum_{j=1}^{M} \left| [\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)]^2 - [\mathscr{N}_{u,v}^{\vartheta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)]^2 \right| \right]$$

$$= \frac{1}{M} \left[ \sum_{j=1}^{M} \left( |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j(\omega)) - \mathscr{N}_{u,v}^{\vartheta,\mathbf{l}}(X_j(\omega))| \right. \right. \tag{5.102}$$

$$\left. \left. \cdot \left| [\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)] + [\mathscr{N}_{u,v}^{\vartheta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)] \right| \right) \right]$$

$$\leq \frac{2}{M} \left[ \sum_{j=1}^{M} \left( \left[ \sup_{x \in D} |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(x) - \mathscr{N}_{u,v}^{\vartheta,\mathbf{l}}(x)| \right] \left[ \sup_{y_1, y_2 \in [u,v]} |y_1 - y_2| \right] \right) \right]$$

$$= 2(v - u) \left[ \sup_{x \in D} |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(x) - \mathscr{N}_{u,v}^{\vartheta,\mathbf{l}}(x)| \right].$$

In addition, combining the assumptions that $D \subseteq [-b, b]^d$, $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1} + 1)$, $\mathbf{l}_0 = d$, $\mathbf{l}_\mathbf{L} = 1$, $b \geq 1$, and $B \geq 1$ with Corollary 5.3.7 (applied with $a \curvearrowright -b$, $b \curvearrowright b$, $u \curvearrowright u$, $v \curvearrowright v$, $d \curvearrowright \mathbf{d}$, $L \curvearrowright \mathbf{L}$, $l \curvearrowright \mathbf{l}$ in the notation of Corollary 5.3.7) shows that for all $\theta, \vartheta \in [-B, B]^{\mathbf{d}}$ it holds that

$$\sup_{x \in D} |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(x) - \mathscr{N}_{u,v}^{\vartheta,\mathbf{l}}(x)| \leq \sup_{x \in [-b,b]^d} |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(x) - \mathscr{N}_{u,v}^{\vartheta,\mathbf{l}}(x)|$$

$$\leq \mathbf{L} \max\{1, b\}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}} (\max\{1, \|\theta\|_\infty, \|\vartheta\|_\infty\})^{\mathbf{L}-1} \|\theta - \vartheta\|_\infty \tag{5.103}$$

$$\leq b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}} B^{\mathbf{L}-1} \|\theta - \vartheta\|_\infty.$$

This and (5.102) imply that for all $\theta, \vartheta \in [-B, B]^{\mathbf{d}}$, $\omega \in \Omega$ it holds that

$$|\mathscr{R}(\theta, \omega) - \mathscr{R}(\vartheta, \omega)| \leq 2(v-u)b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}} B^{\mathbf{L}-1} \|\theta - \vartheta\|_\infty. \tag{5.104}$$

The proof of Lemma 5.3.8 is thus complete. $\qquad\square$

**Corollary 5.3.9.** *Let* $d, \mathbf{d}, \mathfrak{d}, \mathbf{L}, M, K \in \mathbb{N}$, $B, b \in [1, \infty)$, $u \in \mathbb{R}$, $v \in (u, \infty)$, $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \ldots, \mathbf{l}_\mathbf{L}) \in \mathbb{N}^{\mathbf{L}+1}$, $D \subseteq [-b, b]^d$, *assume* $\mathbf{l}_0 = d$, $\mathbf{l}_\mathbf{L} = 1$, *and* $\mathbf{d} \geq \mathfrak{d} = \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1}+1)$, *let* $(\Omega, \mathcal{F}, \mathbb{P})$ *be a probability space, let* $\Theta_k \colon \Omega \to [-B, B]^{\mathbf{d}}$, $k \in \{1, 2, \ldots, K\}$, *be i.i.d. random variables, assume that* $\Theta_1$ *is continuous uniformly distributed on* $[-B, B]^{\mathbf{d}}$, *let* $X_j \colon \Omega \to D$, $j \in \{1, 2, \ldots, M\}$, *and* $Y_j \colon \Omega \to [u, v]$, $j \in \{1, 2, \ldots, M\}$, *be random variables, and let* $\mathscr{R} \colon [-B, B]^{\mathbf{d}} \times \Omega \to [0, \infty)$ *satisfy for all* $\theta \in [-B, B]^{\mathbf{d}}$, $\omega \in \Omega$ *that*

$$\mathscr{R}(\theta, \omega) = \frac{1}{M} \left[ \sum_{j=1}^{M} |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)|^2 \right] \tag{5.105}$$

*(cf. Definition 2.1.27). Then*

*(i) it holds that* $\mathscr{R}$ *is a* $(\mathcal{B}([-B, B]^{\mathbf{d}}) \otimes \mathcal{F})/\mathcal{B}([0, \infty))$-*measurable function and*

*(ii) it holds for all* $\theta \in [-B, B]^{\mathbf{d}}$, $p \in (0, \infty)$ *that*

$$\left( \mathbb{E} \left[ \min_{k \in \{1, 2, \ldots, K\}} |\mathscr{R}(\Theta_k) - \mathscr{R}(\theta)|^p \right] \right)^{1/p} \tag{5.106}$$

$$\leq \frac{4(v-u)b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}} B^{\mathbf{L}} \sqrt{\max\{1, p/\mathfrak{d}\}}}{K^{1/\mathfrak{d}}} \leq \frac{4(v-u)b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}} B^{\mathbf{L}} \max\{1, p\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}}$$

*(cf. Definition 3.1.16).*

*Proof of Corollary 5.3.9.* Throughout this proof let $L \in \mathbb{R}$ be given by $L = 2(v - u)b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^\mathbf{L}B^{\mathbf{L}-1}$, let $P\colon [-B,B]^\mathbf{d} \to [-B,B]^\mathfrak{d}$ satisfy for all $\theta = (\theta_1, \theta_2, \ldots, \theta_\mathbf{d}) \in [-B,B]^\mathbf{d}$ that $P(\theta) = (\theta_1, \theta_2, \ldots, \theta_\mathfrak{d})$, and let $R\colon [-B,B]^\mathfrak{d} \times \Omega \to \mathbb{R}$ satisfy for all $\theta \in [-B,B]^\mathfrak{d}$, $\omega \in \Omega$ that

$$R(\theta,\omega) = \frac{1}{M}\left[\sum_{j=1}^M |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)|^2\right]. \tag{5.107}$$

Note that the fact that $\forall\, \theta \in [-B,B]^\mathbf{d}\colon \mathscr{N}_{u,v}^{\theta,\mathbf{l}} = \mathscr{N}_{u,v}^{P(\theta),\mathbf{l}}$ implies that for all $\theta \in [-B,B]^\mathbf{d}$, $\omega \in \Omega$ it holds that

$$\begin{aligned}\mathscr{R}(\theta,\omega) &= \frac{1}{M}\left[\sum_{j=1}^M |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)|^2\right]\\ &= \frac{1}{M}\left[\sum_{j=1}^M |\mathscr{N}_{u,v}^{P(\theta),\mathbf{l}}(X_j(\omega)) - Y_j(\omega)|^2\right] = R(P(\theta),\omega).\end{aligned} \tag{5.108}$$

Furthermore, Lemma 5.3.8 (applied with $\mathbf{d} \curvearrowright \mathfrak{d}$, $\mathscr{R} \curvearrowright ([-B,B]^\mathfrak{d}\times\Omega \ni (\theta,\omega) \mapsto R(\theta,\omega) \in [0,\infty))$ in the notation of Lemma 5.3.8) demonstrates that for all $\theta, \vartheta \in [-B,B]^\mathfrak{d}$, $\omega \in \Omega$ it holds that

$$|R(\theta,\omega) - R(\vartheta,\omega)| \le 2(v-u)b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^\mathbf{L}B^{\mathbf{L}-1}\|\theta - \vartheta\|_\infty = L\|\theta - \vartheta\|_\infty. \tag{5.109}$$

Moreover, observe that the assumption that $X_j$, $j \in \{1,2,\ldots,M\}$, and $Y_j$, $j \in \{1,2,\ldots,M\}$, are random variables ensures that $R\colon [-B,B]^\mathfrak{d}\times\Omega \to \mathbb{R}$ is a random field. This, (5.109), the fact that $P{\circ}\Theta_k\colon \Omega \to [-B,B]^\mathfrak{d}$, $k \in \{1,2,\ldots,K\}$, are i.i.d. random variables, the fact that $P{\circ}\Theta_1$ is continuous uniformly distributed on $[-B,B]^\mathfrak{d}$, and Proposition 5.2.7 (applied with $\mathbf{d} \curvearrowright \mathfrak{d}$, $\alpha \curvearrowright -B$, $\beta \curvearrowright B$, $\mathscr{R} \curvearrowright R$, $(\Theta_k)_{k\in\{1,2,\ldots,K\}} \curvearrowright (P \circ \Theta_k)_{k\in\{1,2,\ldots,K\}}$ in the notation of Proposition 5.2.7) prove that for all $\theta \in [-B,B]^\mathbf{d}$, $p \in (0,\infty)$ it holds that $R$ is $(\mathcal{B}([-B,B]^\mathfrak{d}) \otimes \mathcal{F})/\mathcal{B}(\mathbb{R})$-measurable and

$$\begin{aligned}&\left(\mathbb{E}\big[\min_{k\in\{1,2,\ldots,K\}}|R(P(\Theta_k)) - R(P(\theta))|^p\big]\right)^{1/p}\\ &\le \frac{L(2B)\max\{1,(p/\mathfrak{d})^{1/\mathfrak{d}}\}}{K^{1/\mathfrak{d}}} = \frac{4(v-u)b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^\mathbf{L}B^\mathbf{L}\max\{1,(p/\mathfrak{d})^{1/\mathfrak{d}}\}}{K^{1/\mathfrak{d}}}.\end{aligned} \tag{5.110}$$

The fact that $P$ is $\mathcal{B}([-B,B]^\mathbf{d})/\mathcal{B}([-B,B]^\mathfrak{d})$-measurable and (5.108) hence show item (i). In addition, (5.108), (5.110), and the fact that $2 \le \mathfrak{d} = \sum_{i=1}^\mathbf{L} \mathbf{l}_i(\mathbf{l}_{i-1} + 1) \le \mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2$ yield that for all $\theta \in [-B,B]^\mathbf{d}$, $p \in (0,\infty)$ it holds that

$$\begin{aligned}&\left(\mathbb{E}\big[\min_{k\in\{1,2,\ldots,K\}}|\mathscr{R}(\Theta_k) - \mathscr{R}(\theta)|^p\big]\right)^{1/p}\\ &= \left(\mathbb{E}\big[\min_{k\in\{1,2,\ldots,K\}}|R(P(\Theta_k)) - R(P(\theta))|^p\big]\right)^{1/p}\\ &\le \frac{4(v-u)b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^\mathbf{L}B^\mathbf{L}\sqrt{\max\{1,p/\mathfrak{d}\}}}{K^{1/\mathfrak{d}}} \le \frac{4(v-u)b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^\mathbf{L}B^\mathbf{L}\max\{1,p\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}}.\end{aligned} \tag{5.111}$$

This establishes item (ii). The proof of Corollary 5.3.9 is thus complete. $\square$

# Chapter 6

# Analysis of the generalisation error

## 6.1 Monte Carlo estimates

**Lemma 6.1.1.** *Let $d, M \in \mathbb{N}$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $X_j \colon \Omega \to \mathbb{R}^d$, $j \in \{1, 2, \dots, M\}$, be independent random variables, and assume $\max_{j \in \{1, 2, \dots, M\}} \mathbb{E}[\|X_j\|_2] < \infty$ (cf. Definition 3.1.16). Then*

$$\left( \mathbb{E}\left[ \left\| \frac{1}{M} \left[ \sum_{j=1}^{M} X_j \right] - \mathbb{E}\left[ \frac{1}{M} \sum_{j=1}^{M} X_j \right] \right\|_2^2 \right] \right)^{1/2} \leq \frac{1}{\sqrt{M}} \left[ \max_{j \in \{1, 2, \dots, M\}} \left( \mathbb{E}\left[ \|X_j - \mathbb{E}[X_j]\|_2^2 \right] \right)^{1/2} \right]. \quad (6.1)$$

*Proof of Lemma 6.1.1.* Throughout this proof let $\langle \cdot, \cdot \rangle \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ satisfy for all $x = (x_1, x_2, \dots, x_d)$, $y = (y_1, y_2, \dots, y_d) \in \mathbb{R}^d$ that $\langle x, y \rangle = \sum_{i=1}^{d} x_i y_i$. Note that the fact that for all $x \in \mathbb{R}^d$ it holds that $\langle x, x \rangle = \|x\|_2^2$ shows that

$$\left\| \frac{1}{M} \left[ \sum_{j=1}^{M} X_j \right] - \mathbb{E}\left[ \frac{1}{M} \sum_{j=1}^{M} X_j \right] \right\|_2^2$$

$$= \frac{1}{M^2} \left\| \left[ \sum_{j=1}^{M} X_j \right] - \mathbb{E}\left[ \sum_{j=1}^{M} X_j \right] \right\|_2^2$$

$$= \frac{1}{M^2} \left\| \sum_{j=1}^{M} \left( X_j - \mathbb{E}[X_j] \right) \right\|_2^2 \qquad (6.2)$$

$$= \frac{1}{M^2} \left[ \sum_{i,j=1}^{M} \left\langle X_i - \mathbb{E}[X_i], X_j - \mathbb{E}[X_j] \right\rangle \right]$$

$$= \frac{1}{M^2} \left[ \sum_{j=1}^{M} \|X_j - \mathbb{E}[X_j]\|_2^2 \right] + \frac{1}{M^2} \left[ \sum_{(i,j) \in \{1,2,\dots,M\}^2, \, i \neq j} \left\langle X_i - \mathbb{E}[X_i], X_j - \mathbb{E}[X_j] \right\rangle \right].$$

This, the fact that for all independent random variables $Y, Z \colon \Omega \to \mathbb{R}^d$ with $\mathbb{E}[\|Y\|_2 + \|Z\|_2] < \infty$ it holds that $\mathbb{E}[|\langle X, Y \rangle|] < \infty$ and $\mathbb{E}[\langle Y, Z \rangle] = \langle \mathbb{E}[Y], \mathbb{E}[Z] \rangle$, and the assumption that $X_j \colon \Omega \to \mathbb{R}^d$, $j \in \{1, 2, \dots, M\}$, are independent random variables imply

that

$$
\mathbb{E}\!\left[\left\|\frac{1}{M}\!\left[\sum_{j=1}^{M} X_j\right] - \mathbb{E}\!\left[\frac{1}{M}\sum_{j=1}^{M} X_j\right]\right\|_2^2\right]
$$

$$
= \frac{1}{M^2}\!\left[\sum_{j=1}^{M} \mathbb{E}\big[\|X_j - \mathbb{E}[X_j]\|_2^2\big]\right] + \frac{1}{M^2}\!\left[\sum_{(i,j)\in\{1,2,\ldots,M\}^2,\, i\neq j} \big\langle \mathbb{E}[X_i - \mathbb{E}[X_i]], \mathbb{E}[X_j - \mathbb{E}[X_j]]\big\rangle\right]
$$

$$
= \frac{1}{M^2}\!\left[\sum_{j=1}^{M} \mathbb{E}\big[\|X_j - \mathbb{E}[X_j]\|_2^2\big]\right]
$$

$$
\leq \frac{1}{M}\!\left[\max_{j\in\{1,2,\ldots,M\}} \mathbb{E}\big[\|X_j - \mathbb{E}[X_j]\|_2^2\big]\right].
$$

$$(6.3)$$

This completes the proof of Lemma 6.1.1. $\qquad\qquad\square$

**Definition 6.1.2** (Rademacher family)**.** *Let $(\Omega,\mathcal{F},\mathbb{P})$ be a probability space and let $J$ be a set. Then we say that $(r_j)_{j\in J}$ is a $\mathbb{P}$-Rademacher family if and only if it holds that $r_j\colon \Omega \to \{-1,1\}$, $j\in J$, are independent random variables with $\forall\, j\in J\colon \mathbb{P}(r_j = 1) = \mathbb{P}(r_j = -1)$.*

**Definition 6.1.3** ($p$-Kahane–Khintchine constant)**.** *Let $p\in(0,\infty)$. Then we denote by $\mathfrak{K}_p \in (0,\infty]$ the extended real number given by*

$$
\mathfrak{K}_p = \sup\left\{c\in[0,\infty)\colon \left[\begin{array}{c} \exists\,\mathbb{R}\text{-}Banach\ space\ (E,\|\cdot\|_E)\colon \\ \exists\,probability\ space\ (\Omega,\mathcal{F},\mathbb{P})\colon \\ \exists\,\mathbb{P}\text{-}Rademacher\ family\ (r_j)_{j\in\mathbb{N}}\colon \\ \exists\,k\in\mathbb{N}\colon \exists\,x_1,x_2,\ldots,x_k\in E\setminus\{0\}\colon \\ \left(\mathbb{E}\big[\|\sum_{j=1}^{k} r_j x_j\|_E^p\big]\right)^{1/p} = c\left(\mathbb{E}\big[\|\sum_{j=1}^{k} r_j x_j\|_E^2\big]\right)^{1/2} \end{array}\right]\right\}
$$

$$(6.4)$$

*(cf. Definition 6.1.2).*

**Lemma 6.1.4.** *It holds for all $p\in[2,\infty)$ that $\mathfrak{K}_p \leq \sqrt{p-1} < \infty$ (cf. Definition 6.1.3).*

**Proposition 6.1.5.** *Let $d,M\in\mathbb{N}$, $p\in[2,\infty)$, let $(\Omega,\mathcal{F},\mathbb{P})$ be a probability space, let $X_j\colon \Omega \to \mathbb{R}^d$, $j\in\{1,2,\ldots,M\}$, be independent random variables, and assume $\max_{j\in\{1,2,\ldots,M\}} \mathbb{E}[\|X_j\|_2] < \infty$ (cf. Definition 3.1.16). Then*

$$
\left(\mathbb{E}\!\left[\left\|\left[\sum_{j=1}^{M} X_j\right] - \mathbb{E}\!\left[\sum_{j=1}^{M} X_j\right]\right\|_2^p\right]\right)^{1/p} \leq 2\mathfrak{K}_p\left[\sum_{j=1}^{M}\big(\mathbb{E}\big[\|X_j - \mathbb{E}[X_j]\|_2^p\big]\big)^{2/p}\right]^{1/2}
$$

$$(6.5)$$

*(cf. Definition 6.1.3 and Lemma 6.1.4).*

**Corollary 6.1.6.** *Let $d,M\in\mathbb{N}$, $p\in[2,\infty)$, let $(\Omega,\mathcal{F},\mathbb{P})$ be a probability space, let $X_j\colon \Omega \to \mathbb{R}^d$, $j\in\{1,2,\ldots,M\}$, be independent random variables, and assume $\max_{j\in\{1,2,\ldots,M\}} \mathbb{E}[\|X_j\|_2] < \infty$ (cf. Definition 3.1.16). Then*

$$
\left(\mathbb{E}\!\left[\left\|\frac{1}{M}\!\left[\sum_{j=1}^{M} X_j\right] - \mathbb{E}\!\left[\frac{1}{M}\sum_{j=1}^{M} X_j\right]\right\|_2^p\right]\right)^{1/p} \leq \frac{2\sqrt{p-1}}{\sqrt{M}}\!\left[\max_{j\in\{1,2,\ldots,M\}}\big(\mathbb{E}\big[\|X_j - \mathbb{E}[X_j]\|_2^p\big]\big)^{1/p}\right].
$$

$$(6.6)$$

*Proof of Corollary 6.1.6.* Observe that Proposition 6.1.5 and Lemma 6.1.4 imply that

$$
\begin{aligned}
&\left(\mathbb{E}\left[\left\|\frac{1}{M}\left[\sum_{j=1}^{M} X_j\right] - \mathbb{E}\left[\frac{1}{M}\sum_{j=1}^{M} X_j\right]\right\|_2^p\right]\right)^{1/p} \\
&= \frac{1}{M}\left(\mathbb{E}\left[\left\|\left[\sum_{j=1}^{M} X_j\right] - \mathbb{E}\left[\sum_{j=1}^{M} X_j\right]\right\|_2^p\right]\right)^{1/p} \\
&\leq \frac{2\mathfrak{K}_p}{M}\left[\sum_{j=1}^{M}\left(\mathbb{E}\left[\|X_j - \mathbb{E}[X_j]\|_2^p\right]\right)^{2/p}\right]^{1/2} \\
&\leq \frac{2\mathfrak{K}_p}{M}\left[M\left(\max_{j\in\{1,2,\ldots,M\}}\left(\mathbb{E}\left[\|X_j - \mathbb{E}[X_j]\|_2^p\right]\right)^{2/p}\right)\right]^{1/2} \\
&= \frac{2\mathfrak{K}_p}{\sqrt{M}}\left[\max_{j\in\{1,2,\ldots,M\}}\left(\mathbb{E}\left[\|X_j - \mathbb{E}[X_j]\|_2^p\right]\right)^{1/p}\right] \\
&\leq \frac{2\sqrt{p-1}}{\sqrt{M}}\left[\max_{j\in\{1,2,\ldots,M\}}\left(\mathbb{E}\left[\|X_j - \mathbb{E}[X_j]\|_2^p\right]\right)^{1/p}\right]
\end{aligned}
\tag{6.7}
$$

(cf. Definition 6.1.3). The proof of Corollary 6.1.6 is thus complete. □

## 6.2 Uniform strong error estimates for random fields

**Lemma 6.2.1.** *Let $(E, \mathscr{E})$ be a separable topological space, let $(\Omega, \mathcal{F})$ be a measurable space, let $f_x\colon \Omega \to \mathbb{R}$, $x \in E$, be $\mathcal{F}/\mathcal{B}(\mathbb{R})$-measurable, and assume for all $\omega \in \Omega$ that $E \ni x \mapsto f_x(\omega) \in \mathbb{R}$ is continuous. Then it holds that*

$$
\Omega \ni \omega \mapsto \sup\left(\{f_x(\omega)\colon x \in E\} \cup \{0\}\right) \in \mathbb{R} \cup \{\infty\}
\tag{6.8}
$$

*is $\mathcal{F}/\mathcal{B}(\mathbb{R} \cup \{\infty\})$-measurable.*

*Proof of Lemma 6.2.1.* Throughout this proof assume w.l.o.g. that $E \neq \emptyset$, let $F\colon \Omega \to \mathbb{R} \cup \{\infty\}$ satisfy for all $\omega \in \Omega$ that $F(\omega) = \sup_{x\in E} f_x(\omega)$, and let $\mathbf{E} \subseteq E$ be an at most countable and dense subset of $E$. Note that the fact that $\mathbf{E}$ is dense in $E$ implies that for all $g \in C(E, \mathbb{R})$ it holds that

$$
\sup_{x\in E} g(x) = \sup_{x\in\mathbf{E}} g(x).
\tag{6.9}
$$

This and the assumption that for all $\omega \in \Omega$ it holds that $E \ni x \mapsto f_x(\omega) \in \mathbb{R}$ is a continuous function show that for all $\omega \in \Omega$ it holds that

$$
F(\omega) = \sup_{x\in E} f_x(\omega) = \sup_{x\in\mathbf{E}} f_x(\omega).
\tag{6.10}
$$

The assumption that for all $x \in E$ it holds that $f_x$ is $\mathcal{F}/\mathcal{B}(\mathbb{R})$-measurable hence demonstrates that $F$ is $\mathcal{F}/\mathcal{B}(\mathbb{R})$-measurable. The proof of Lemma 6.2.1 is thus complete. □

**Lemma 6.2.2.** *Let $(E, \delta)$ be a separable metric space, assume $E \neq \emptyset$, let $L \in \mathbb{R}$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $Z_x\colon \Omega \to \mathbb{R}$, $x \in E$, be random variables, and assume for all $x, y \in E$ that $\mathbb{E}[|Z_x|] < \infty$ and $|Z_x - Z_y| \leq L\delta(x, y)$. Then it holds that*

$$
\Omega \ni \omega \mapsto \sup_{x\in E}|Z_x(\omega) - \mathbb{E}[Z_x]| \in [0, \infty]
\tag{6.11}
$$

*is $\mathcal{F}/\mathcal{B}([0, \infty])$-measurable.*

*Proof of Lemma 6.2.2.* Note that the assumption that for all $x, y \in E$ it holds that $|Z_x - Z_y| \leq L\delta(x,y)$ shows that for all $x, y \in E$ it holds that

$$
\begin{aligned}
|(Z_x - \mathbb{E}[Z_x]) - (Z_y - \mathbb{E}[Z_y])| &= |(Z_x - Z_y) + (\mathbb{E}[Z_y] - \mathbb{E}[Z_x])| \leq |Z_x - Z_y| + |\mathbb{E}[Z_x] - \mathbb{E}[Z_y]| \\
&\leq L\delta(x,y) + |\mathbb{E}[Z_x] - \mathbb{E}[Z_y]| = L\delta(x,y) + |\mathbb{E}[Z_x - Z_y]| \\
&\leq L\delta(x,y) + \mathbb{E}[|Z_x - Z_y|] \leq L\delta(x,y) + L\delta(x,y) = 2L\delta(x,y).
\end{aligned}
\tag{6.12}
$$

This implies that for all $\omega \in \Omega$ it holds that $E \ni x \mapsto |Z_x(\omega) - \mathbb{E}[Z_x]| \in \mathbb{R}$ is a continuous function. Combining this and the assumption that $E$ is separable with Lemma 6.2.1 completes the proof of Lemma 6.2.2. $\qquad\square$

**Lemma 6.2.3.** *Let $(E, \delta)$ be a separable metric space, let $N \in \mathbb{N}$, $p, L, r_1, r_2, \ldots, r_N \in [0, \infty)$, $z_1, z_2, \ldots, z_N \in E$ satisfy $E \subseteq \left( \bigcup_{i=1}^{N} \{ x \in E \colon \delta(x, z_i) \leq r_i \} \right)$, let $(\Omega, \mathcal{F}, \mu)$ be a measure space, let $Z_x \colon \Omega \to \mathbb{R}$, $x \in E$, be $\mathcal{F}/\mathcal{B}(\mathbb{R})$-measurable, and assume for all $x, y \in E$ that $|Z_x - Z_y| \leq L\delta(x, y)$. Then*

$$
\int_\Omega \sup_{x \in E} |Z_x(\omega)|^p \, \mu(\mathrm{d}\omega) \leq \sum_{i=1}^{N} \int_\Omega (Lr_i + |Z_{z_i}(\omega)|)^p \, \mu(\mathrm{d}\omega)
\tag{6.13}
$$

*(cf. Lemma 6.2.1).*

*Proof of Lemma 6.2.3.* Throughout this proof let $B_1, B_2, \ldots, B_N \subseteq E$ satisfy for all $i \in \{1, 2, \ldots, N\}$ that $B_i = \{ x \in E \colon \delta(x, z_i) \leq r_i \}$. Note that the fact that $E = \bigcup_{i=1}^{N} B_i$ shows that

$$
\sup_{x \in E} |Z_x| = \sup_{x \in \left( \bigcup_{i=1}^{N} B_i \right)} |Z_x| = \max_{i \in \{1, 2, \ldots, N\}} \sup_{x \in B_i} |Z_x|.
\tag{6.14}
$$

This establishes that

$$
\begin{aligned}
\int_\Omega \sup_{x \in E} |Z_x(\omega)|^p \, \mu(\mathrm{d}\omega) &= \int_\Omega \max_{i \in \{1, 2, \ldots, N\}} \sup_{x \in B_i} |Z_x(\omega)|^p \, \mu(\mathrm{d}\omega) \\
&\leq \int_\Omega \sum_{i=1}^{N} \sup_{x \in B_i} |Z_x(\omega)|^p \, \mu(\mathrm{d}\omega) = \sum_{i=1}^{N} \int_\Omega \sup_{x \in B_i} |Z_x(\omega)|^p \, \mu(\mathrm{d}\omega).
\end{aligned}
\tag{6.15}
$$

Furthermore, the assumption that $\forall\, x, y \in E \colon |Z_x - Z_y| \leq L\delta(x, y)$ implies that for all $i \in \{1, 2, \ldots, N\}$, $x \in B_i$ it holds that

$$
|Z_x| = |Z_x - Z_{z_i} + Z_{z_i}| \leq |Z_x - Z_{z_i}| + |Z_{z_i}| \leq L\delta(x, z_i) + |Z_{z_i}| \leq Lr_i + |Z_{z_i}|.
\tag{6.16}
$$

Combining this with (6.15) proves that

$$
\int_\Omega \sup_{x \in E} |Z_x(\omega)|^p \, \mu(\mathrm{d}\omega) \leq \sum_{i=1}^{N} \int_\Omega (Lr_i + |Z_{z_i}(\omega)|)^p \, \mu(\mathrm{d}\omega).
\tag{6.17}
$$

The proof of Lemma 6.2.3 is thus complete. $\qquad\square$

**Lemma 6.2.4.** *Let $p, L, r \in (0, \infty)$, let $(E, \delta)$ be a separable metric space, let $(\Omega, \mathcal{F}, \mu)$ be a measure space, assume $E \neq \emptyset$ and $\mu(\Omega) \neq 0$, let $Z_x \colon \Omega \to \mathbb{R}$, $x \in E$, be $\mathcal{F}/\mathcal{B}(\mathbb{R})$-measurable, and assume for all $x, y \in E$ that $|Z_x - Z_y| \leq L\delta(x, y)$. Then*

$$
\int_\Omega \sup_{x \in E} |Z_x(\omega)|^p \, \mu(\mathrm{d}\omega) \leq \mathcal{C}^{(E, \delta), r} \left[ \sup_{x \in E} \int_\Omega (Lr + |Z_x(\omega)|)^p \, \mu(\mathrm{d}\omega) \right]
\tag{6.18}
$$

*(cf. Definition 3.2.13 and Lemma 6.2.1).*

*Proof of Lemma 6.2.4.* Throughout this proof assume w.l.o.g. that $\mathcal{C}^{(E,\delta),r} < \infty$, let $N \in \mathbb{N}$ be given by $N = \mathcal{C}^{(E,\delta),r}$, and let $z_1, z_2, \ldots, z_N \in E$ satisfy $E \subseteq \bigcup_{i=1}^{N}\{x \in E \colon \delta(x, z_i) \leq r\}$. Note that Lemma 6.2.3 (applied with $r_1 \curvearrowleft r, r_2 \curvearrowleft r, \ldots, r_N \curvearrowleft r$ in the notation of Lemma 6.2.3) establishes that

$$
\begin{aligned}
\int_{\Omega} \sup_{x \in E} |Z_x(\omega)|^p \, \mu(\mathrm{d}\omega) &\leq \sum_{i=1}^{N} \int_{\Omega} (Lr + |Z_{z_i}(\omega)|)^p \, \mu(\mathrm{d}\omega) \\
&\leq \sum_{i=1}^{N} \left[ \sup_{x \in E} \int_{\Omega} (Lr + |Z_x(\omega)|)^p \, \mu(\mathrm{d}\omega) \right] = N \left[ \sup_{x \in E} \int_{\Omega} (Lr + |Z_x(\omega)|)^p \, \mu(\mathrm{d}\omega) \right].
\end{aligned}
\tag{6.19}
$$

The proof of Lemma 6.2.4 is thus complete. $\qquad\square$

**Lemma 6.2.5.** *Let $p \in [1, \infty)$, $L, r \in (0, \infty)$, let $(E, \delta)$ be a separable metric space, assume $E \neq \emptyset$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $Z_x \colon \Omega \to \mathbb{R}$, $x \in E$, be random variables, and assume for all $x, y \in E$ that $\mathbb{E}[|Z_x|] < \infty$ and $|Z_x - Z_y| \leq L\delta(x, y)$. Then*

$$
\left(\mathbb{E}\big[\sup_{x \in E} |Z_x - \mathbb{E}[Z_x]|^p\big]\right)^{1/p} \leq (\mathcal{C}^{(E,\delta),r})^{1/p}\left[2Lr + \sup_{x \in E}\big(\mathbb{E}\big[|Z_x - \mathbb{E}[Z_x]|^p\big]\big)^{1/p}\right] \tag{6.20}
$$

*(cf. Definition 3.2.13 and Lemma 6.2.2).*

*Proof of Lemma 6.2.5.* Throughout this proof let $Y_x \colon \Omega \to \mathbb{R}$, $x \in E$, satisfy for all $x \in E$, $\omega \in \Omega$ that $Y_x(\omega) = Z_x(\omega) - \mathbb{E}[Z_x]$. Note that it holds for all $x, y \in E$ that

$$
\begin{aligned}
|Y_x - Y_y| = |(Z_x - \mathbb{E}[Z_x]) - (Z_y - \mathbb{E}[Z_y])| &\leq |Z_x - Z_y| + |\mathbb{E}[Z_x] - \mathbb{E}[Z_y]| \\
&\leq L\delta(x, y) + \mathbb{E}[|Z_x - Z_y|] \leq 2L\delta(x, y).
\end{aligned}
\tag{6.21}
$$

Combining this with Lemma 6.2.4 (applied with $L \curvearrowleft 2L$, $(\Omega, \mathcal{F}, \mu) \curvearrowleft (\Omega, \mathcal{F}, \mathbb{P})$, $(Z_x)_{x \in E} \curvearrowleft (Y_x)_{x \in E}$ in the notation of Lemma 6.2.4) implies that

$$
\begin{aligned}
\left(\mathbb{E}\big[\sup_{x \in E} |Z_x - \mathbb{E}[Z_x]|^p\big]\right)^{1/p} &= \left(\mathbb{E}\big[\sup_{x \in E} |Y_x|^p\big]\right)^{1/p} \\
&\leq (\mathcal{C}^{(E,\delta),r})^{1/p}\left[\sup_{x \in E}\big(\mathbb{E}\big[(2Lr + |Y_x|)^p\big]\big)^{1/p}\right] \\
&\leq (\mathcal{C}^{(E,\delta),r})^{1/p}\left[2Lr + \sup_{x \in E}\big(\mathbb{E}\big[|Y_x|^p\big]\big)^{1/p}\right] \\
&= (\mathcal{C}^{(E,\delta),r})^{1/p}\left[2Lr + \sup_{x \in E}\big(\mathbb{E}\big[|Z_x - \mathbb{E}[Z_x]|^p\big]\big)^{1/p}\right].
\end{aligned}
\tag{6.22}
$$

The proof of Lemma 6.2.5 is thus complete. $\qquad\square$

**Lemma 6.2.6.** *Let $M \in \mathbb{N}$, $p \in [2, \infty)$, $L, r \in (0, \infty)$, let $(E, \delta)$ be a separable metric space, assume $E \neq \emptyset$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, for every $x \in E$ let $Y_{x,j} \colon \Omega \to \mathbb{R}$, $j \in \{1, 2, \ldots, M\}$, be independent random variables, assume for all $x, y \in E$, $j \in \{1, 2, \ldots, M\}$ that $\mathbb{E}[|Y_{x,j}|] < \infty$ and $|Y_{x,j} - Y_{y,j}| \leq L\delta(x, y)$, and let $Z_x \colon \Omega \to \mathbb{R}$, $x \in E$, satisfy for all $x \in E$ that*

$$
Z_x = \frac{1}{M}\left[\sum_{j=1}^{M} Y_{x,j}\right]. \tag{6.23}
$$

*Then*

*(i) it holds for all $x \in E$ that $\mathbb{E}[|Z_x|] < \infty$,*

*(ii) it holds that $\Omega \ni \omega \mapsto \sup_{x \in E}|Z_x(\omega) - \mathbb{E}[Z_x]| \in [0, \infty]$ is $\mathcal{F}/\mathcal{B}([0, \infty])$-measurable, and*

*(iii) it holds that*

$$\left(\mathbb{E}\big[\sup_{x\in E}|Z_x - \mathbb{E}[Z_x]|^p\big]\right)^{1/p}$$
$$\leq 2(\mathcal{C}^{(E,\delta),r})^{1/p}\Big[Lr + \tfrac{\sqrt{p-1}}{\sqrt{M}}\Big(\sup_{x\in E}\max_{j\in\{1,2,\dots,M\}}\big(\mathbb{E}\big[|Y_{x,j} - \mathbb{E}[Y_{x,j}]|^p\big]\big)^{1/p}\Big)\Big] \tag{6.24}$$

*(cf. Definition 3.2.13).*

*Proof of Lemma 6.2.6.* Note that the assumption that $\forall\, x \in E,\ j \in \{1,2,\dots,M\}\colon \mathbb{E}[|Y_{x,j}|] < \infty$ implies that for all $x \in E$ it holds that

$$\mathbb{E}[|Z_x|] = \mathbb{E}\left[\frac{1}{M}\left|\sum_{j=1}^M Y_{x,j}\right|\right] \leq \frac{1}{M}\left[\sum_{j=1}^M \mathbb{E}[|Y_{x,j}|]\right] \leq \max_{j\in\{1,2,\dots,M\}}\mathbb{E}[|Y_{x,j}|] < \infty. \tag{6.25}$$

This proves item (i). Next observe that the assumption that $\forall\, x,y \in E,\ j \in \{1,2,\dots,M\}\colon |Y_{x,j} - Y_{y,j}| \leq L\delta(x,y)$ demonstrates that for all $x,y \in E$ it holds that

$$|Z_x - Z_y| = \frac{1}{M}\left|\left[\sum_{j=1}^M Y_{x,j}\right] - \left[\sum_{j=1}^M Y_{y,j}\right]\right| \leq \frac{1}{M}\left[\sum_{j=1}^M |Y_{x,j} - Y_{y,j}|\right] \leq L\delta(x,y). \tag{6.26}$$

Combining this with item (i) and Lemma 6.2.2 establishes item (ii). It thus remains to show item (iii). For this note that item (i), (6.26), and Lemma 6.2.5 yield that

$$\left(\mathbb{E}\big[\sup_{x\in E}|Z_x - \mathbb{E}[Z_x]|^p\big]\right)^{1/p} \leq (\mathcal{C}^{(E,\delta),r})^{1/p}\Big[2Lr + \sup_{x\in E}\big(\mathbb{E}\big[|Z_x - \mathbb{E}[Z_x]|^p\big]\big)^{1/p}\Big]. \tag{6.27}$$

Moreover, (6.25) and Corollary 6.1.6 (applied with $d \curvearrowleft 1$, $(X_j)_{j\in\{1,2,\dots,M\}} \curvearrowleft (Y_{x,j})_{j\in\{1,2,\dots,M\}}$ for $x \in E$ in the notation of Corollary 6.1.6) prove that for all $x \in E$ it holds that

$$\left(\mathbb{E}\big[|Z_x - \mathbb{E}[Z_x]|^p\big]\right)^{1/p} = \left(\mathbb{E}\left[\left|\frac{1}{M}\left[\sum_{j=1}^M Y_{x,j}\right] - \mathbb{E}\left[\frac{1}{M}\sum_{j=1}^M Y_{x,j}\right]\right|^p\right]\right)^{1/p}$$
$$\leq \frac{2\sqrt{p-1}}{\sqrt{M}}\left[\max_{j\in\{1,2,\dots,M\}}\big(\mathbb{E}\big[|Y_{x,j} - \mathbb{E}[Y_{x,j}]|^p\big]\big)^{1/p}\right]. \tag{6.28}$$

This and (6.27) imply that

$$\left(\mathbb{E}\big[\sup_{x\in E}|Z_x - \mathbb{E}[Z_x]|^p\big]\right)^{1/p}$$
$$\leq (\mathcal{C}^{(E,\delta),r})^{1/p}\Big[2Lr + \tfrac{2\sqrt{p-1}}{\sqrt{M}}\Big(\sup_{x\in E}\max_{j\in\{1,2,\dots,M\}}\big(\mathbb{E}\big[|Y_{x,j} - \mathbb{E}[Y_{x,j}]|^p\big]\big)^{1/p}\Big)\Big] \tag{6.29}$$
$$= 2(\mathcal{C}^{(E,\delta),r})^{1/p}\Big[Lr + \tfrac{\sqrt{p-1}}{\sqrt{M}}\Big(\sup_{x\in E}\max_{j\in\{1,2,\dots,M\}}\big(\mathbb{E}\big[|Y_{x,j} - \mathbb{E}[Y_{x,j}]|^p\big]\big)^{1/p}\Big)\Big].$$

The proof of Lemma 6.2.6 is thus complete. $\qquad\square$

**Corollary 6.2.7.** *Let $M \in \mathbb{N}$, $p \in [2,\infty)$, $L,C \in (0,\infty)$, let $(E,\delta)$ be a separable metric space, assume $E \neq \emptyset$, let $(\Omega,\mathcal{F},\mathbb{P})$ be a probability space, for every $x \in E$ let $Y_{x,j}\colon \Omega \to \mathbb{R}$, $j \in \{1,2,\dots,M\}$, be independent random variables, assume for all $x,y \in E$, $j \in \{1,2,\dots,M\}$ that $\mathbb{E}[|Y_{x,j}|] < \infty$ and $|Y_{x,j} - Y_{y,j}| \leq L\delta(x,y)$, and let $Z_x\colon \Omega \to \mathbb{R}$, $x \in E$, satisfy for all $x \in E$ that*

$$Z_x = \frac{1}{M}\left[\sum_{j=1}^M Y_{x,j}\right]. \tag{6.30}$$

*Then*

(i) it holds for all $x \in E$ that $\mathbb{E}[|Z_x|] < \infty$,

(ii) it holds that $\Omega \ni \omega \mapsto \sup_{x \in E}|Z_x(\omega) - \mathbb{E}[Z_x]| \in [0,\infty]$ *is $\mathcal{F}/\mathcal{B}([0,\infty])$-measurable, and*

(iii) it holds that

$$
\left(\mathbb{E}\big[\sup_{x \in E}|Z_x - \mathbb{E}[Z_x]|^p\big]\right)^{1/p}
$$
$$
\leq \tfrac{2\sqrt{p-1}}{\sqrt{M}}\left(\mathcal{C}^{(E,\delta),\frac{C\sqrt{p-1}}{L\sqrt{M}}}\right)^{1/p}\left[C + \sup_{x \in E}\max_{j \in \{1,2,\ldots,M\}}\left(\mathbb{E}\big[|Y_{x,j} - \mathbb{E}[Y_{x,j}]|^p\big]\right)^{1/p}\right]
$$
(6.31)

*(cf. Definition 3.2.13).*

*Proof of Corollary 6.2.7.* Note that Lemma 6.2.6 shows items (i) and (ii). In addition, Lemma 6.2.6 (applied with $r \curvearrowleft {}^{C\sqrt{p-1}}\!/\!_{(L\sqrt{M})}$ in the notation of Lemma 6.2.6) ensures that

$$
\left(\mathbb{E}\big[\sup_{x \in E}|Z_x - \mathbb{E}[Z_x]|^p\big]\right)^{1/p}
$$
$$
\leq 2\left(\mathcal{C}^{(E,\delta),\frac{C\sqrt{p-1}}{L\sqrt{M}}}\right)^{1/p}\left[L\frac{C\sqrt{p-1}}{L\sqrt{M}} + \frac{\sqrt{p-1}}{\sqrt{M}}\left(\sup_{x \in E}\max_{j \in \{1,2,\ldots,M\}}\left(\mathbb{E}\big[|Y_{x,j} - \mathbb{E}[Y_{x,j}]|^p\big]\right)^{1/p}\right)\right]
$$
$$
= \tfrac{2\sqrt{p-1}}{\sqrt{M}}\left(\mathcal{C}^{(E,\delta),\frac{C\sqrt{p-1}}{L\sqrt{M}}}\right)^{1/p}\left[C + \sup_{x \in E}\max_{j \in \{1,2,\ldots,M\}}\left(\mathbb{E}\big[|Y_{x,j} - \mathbb{E}[Y_{x,j}]|^p\big]\right)^{1/p}\right].
$$
(6.32)

This establishes item (iii) and thus completes the proof of Corollary 6.2.7. $\qquad\square$

## 6.3 Strong convergence rates for the generalisation error

**Lemma 6.3.1.** *Let $M \in \mathbb{N}$, $p \in [2,\infty)$, $L,C,b \in (0,\infty)$, let $(E,\delta)$ be a separable metric space, assume $E \neq \emptyset$, let $(\Omega,\mathcal{F},\mathbb{P})$ be a probability space, let $X_{x,j}\colon \Omega \to \mathbb{R}$, $j \in \{1,2,\ldots,M\}$, $x \in E$, and $Y_j\colon \Omega \to \mathbb{R}$, $j \in \{1,2,\ldots,M\}$, be functions, assume for all $x \in E$ that $(X_{x,j},Y_j)$, $j \in \{1,2,\ldots,M\}$, are i.i.d. random variables, assume for all $x,y \in E$, $j \in \{1,2,\ldots,M\}$ that $|X_{x,j} - Y_j| \leq b$ and $|X_{x,j} - X_{y,j}| \leq L\delta(x,y)$, let $\mathbf{R}\colon E \to [0,\infty)$ satisfy for all $x \in E$ that $\mathbf{R}(x) = \mathbb{E}[|X_{x,1} - Y_1|^2]$, and let $\mathscr{R}\colon E \times \Omega \to [0,\infty)$ satisfy for all $x \in E$, $\omega \in \Omega$ that*

$$
\mathscr{R}(x,\omega) = \frac{1}{M}\left[\sum_{j=1}^{M}|X_{x,j}(\omega) - Y_j(\omega)|^2\right].
$$
(6.33)

*Then*

(i) *it holds that $\Omega \ni \omega \mapsto \sup_{x \in E}|\mathscr{R}(x,\omega) - \mathbf{R}(x)| \in [0,\infty]$ is $\mathcal{F}/\mathcal{B}([0,\infty])$-measurable and*

(ii) *it holds that*

$$
\left(\mathbb{E}\big[\sup_{x \in E}|\mathscr{R}(x) - \mathbf{R}(x)|^p\big]\right)^{1/p} \leq \left(\mathcal{C}^{(E,\delta),\frac{Cb\sqrt{p-1}}{2L\sqrt{M}}}\right)^{1/p}\left[\frac{2(C+1)b^2\sqrt{p-1}}{\sqrt{M}}\right]
$$
(6.34)

*(cf. Definition 3.2.13).*

*Proof of Lemma 6.3.1.* Throughout this proof let $\mathcal{Y}_{x,j} \colon \Omega \to \mathbb{R}$, $j \in \{1, 2, \ldots, M\}$, $x \in E$, satisfy for all $x \in E$, $j \in \{1, 2, \ldots, M\}$ that $\mathcal{Y}_{x,j} = |X_{x,j} - Y_j|^2$. Note that the assumption that for all $x \in E$ it holds that $(X_{x,j}, Y_j)$, $j \in \{1, 2, \ldots, M\}$, are i.i.d. random variables ensures that for all $x \in E$ it holds that

$$\mathbb{E}[\mathscr{R}(x)] = \frac{1}{M}\left[\sum_{j=1}^{M} \mathbb{E}\big[|X_{x,j} - Y_j|^2\big]\right] = \frac{M\,\mathbb{E}\big[|X_{x,1} - Y_1|^2\big]}{M} = \mathbf{R}(x). \tag{6.35}$$

Furthermore, the assumption that $\forall\, x \in E$, $j \in \{1, 2, \ldots, M\} \colon |X_{x,j} - Y_j| \le b$ shows that for all $x \in E$, $j \in \{1, 2, \ldots, M\}$ it holds that

$$\mathbb{E}[|\mathcal{Y}_{x,j}|] = \mathbb{E}[|X_{x,j} - Y_j|^2] \le b^2 < \infty, \tag{6.36}$$

$$\mathcal{Y}_{x,j} - \mathbb{E}[\mathcal{Y}_{x,j}] = |X_{x,j} - Y_j|^2 - \mathbb{E}\big[|X_{x,j} - Y_j|^2\big] \le |X_{x,j} - Y_j|^2 \le b^2, \tag{6.37}$$

and

$$\mathbb{E}[\mathcal{Y}_{x,j}] - \mathcal{Y}_{x,j} = \mathbb{E}\big[|X_{x,j} - Y_j|^2\big] - |X_{x,j} - Y_j|^2 \le \mathbb{E}\big[|X_{x,j} - Y_j|^2\big] \le b^2. \tag{6.38}$$

Combining (6.36)–(6.38) implies for all $x \in E$, $j \in \{1, 2, \ldots, M\}$ that

$$\big(\mathbb{E}\big[|\mathcal{Y}_{x,j} - \mathbb{E}[\mathcal{Y}_{x,j}]|^p\big]\big)^{1/p} \le \big(\mathbb{E}\big[b^{2p}\big]\big)^{1/p} = b^2. \tag{6.39}$$

Moreover, note that the assumptions that $\forall\, x, y \in E$, $j \in \{1, 2, \ldots, M\} \colon [|X_{x,j} - Y_j| \le b$ and $|X_{x,j} - X_{y,j}| \le L\delta(x,y)]$ and the fact that $\forall\, x_1, x_2, y \in \mathbb{R} \colon (x_1 - y)^2 - (x_2 - y)^2 = (x_1 - x_2)((x_1 - y) + (x_2 - y))$ establish that for all $x, y \in E$, $j \in \{1, 2, \ldots, M\}$ it holds that

$$\begin{aligned}
|\mathcal{Y}_{x,j} - \mathcal{Y}_{y,j}| &= |(X_{x,j} - Y_j)^2 - (X_{y,j} - Y_j)^2| \\
&\le |X_{x,j} - X_{y,j}|(|X_{x,j} - Y_j| + |X_{y,j} - Y_j|) \\
&\le 2b|X_{x,j} - X_{y,j}| \le 2bL\delta(x,y).
\end{aligned} \tag{6.40}$$

Combining this, (6.35), (6.36), and the fact that for all $x \in E$ it holds that $\mathcal{Y}_{x,j}$, $j \in \{1, 2, \ldots, M\}$, are independent random variables with Corollary 6.2.7 (applied with $L \curvearrowleft 2bL$, $C \curvearrowleft Cb^2$, $(Y_{x,j})_{x \in E, j \in \{1,2,\ldots,M\}} \curvearrowleft (\mathcal{Y}_{x,j})_{x \in E, j \in \{1,2,\ldots,M\}}$, $(Z_x)_{x \in E} \curvearrowleft (\Omega \ni \omega \mapsto \mathscr{R}(x,\omega) \in \mathbb{R})_{x \in E}$ in the notation of Corollary 6.2.7) and (6.39) proves item (i) and

$$\begin{aligned}
\big(\mathbb{E}\big[\sup_{x \in E}|\mathscr{R}(x) - \mathbf{R}(x)|^p\big]\big)^{1/p} &= \big(\mathbb{E}\big[\sup_{x \in E}|\mathscr{R}(x) - \mathbb{E}[\mathscr{R}(x)]|^p\big]\big)^{1/p} \\
&\le \tfrac{2\sqrt{p-1}}{\sqrt{M}}\Big(\mathcal{C}^{(E,\delta),\,\frac{Cb^2\sqrt{p-1}}{2bL\sqrt{M}}}\Big)^{1/p}\Big[Cb^2 + \sup_{x \in E}\max_{j \in \{1,2,\ldots,M\}}\big(\mathbb{E}\big[|\mathcal{Y}_{x,j} - \mathbb{E}[\mathcal{Y}_{x,j}]|^p\big]\big)^{1/p}\Big] \\
&\le \tfrac{2\sqrt{p-1}}{\sqrt{M}}\Big(\mathcal{C}^{(E,\delta),\,\frac{Cb\sqrt{p-1}}{2L\sqrt{M}}}\Big)^{1/p}[Cb^2 + b^2] = \Big(\mathcal{C}^{(E,\delta),\,\frac{Cb\sqrt{p-1}}{2L\sqrt{M}}}\Big)^{1/p}\left[\tfrac{2(C+1)b^2\sqrt{p-1}}{\sqrt{M}}\right].
\end{aligned} \tag{6.41}$$

This shows item (ii) and thus completes the proof of Lemma 6.3.1. $\qquad\square$

**Proposition 6.3.2.** *Let $d, \mathbf{d}, M \in \mathbb{N}$, $L, b \in (0, \infty)$, $\alpha \in \mathbb{R}$, $\beta \in (\alpha, \infty)$, $D \subseteq \mathbb{R}^d$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $\mathbb{X}_j = (X_j, Y_j) \colon \Omega \to (D \times \mathbb{R})$, $j \in \{1, 2, \ldots, M\}$, be i.i.d. random variables, let $f = (f_\theta)_{\theta \in [\alpha,\beta]^{\mathbf{d}}} \colon [\alpha, \beta]^{\mathbf{d}} \to C(D, \mathbb{R})$, assume for all $\theta, \vartheta \in [\alpha, \beta]^{\mathbf{d}}$, $j \in \{1, 2, \ldots, M\}$, $x \in D$ that $|f_\theta(X_j) - Y_j| \le b$ and $|f_\theta(x) - f_\vartheta(x)| \le L\|\theta - \vartheta\|_\infty$, let $\mathbf{R} \colon [\alpha, \beta]^{\mathbf{d}} \to [0, \infty)$ satisfy for all $\theta \in [\alpha, \beta]^{\mathbf{d}}$ that $\mathbf{R}(\theta) = \mathbb{E}[|f_\theta(X_1) - Y_1|^2]$, and let $\mathscr{R} \colon [\alpha, \beta]^{\mathbf{d}} \times \Omega \to [0, \infty)$ satisfy for all $\theta \in [\alpha, \beta]^{\mathbf{d}}$, $\omega \in \Omega$ that*

$$\mathscr{R}(\theta, \omega) = \frac{1}{M}\left[\sum_{j=1}^{M}|f_\theta(X_j(\omega)) - Y_j(\omega)|^2\right] \tag{6.42}$$

*(cf. Definition 3.1.16). Then*

(i) it holds that $\Omega \ni \omega \mapsto \sup_{\theta \in [\alpha,\beta]^{\mathbf{d}}} |\mathscr{R}(\theta,\omega) - \mathbf{R}(\theta)| \in [0,\infty]$ is $\mathcal{F}/\mathcal{B}([0,\infty])$-measurable and

(ii) it holds for all $p \in (0,\infty)$ that

$$
\left( \mathbb{E}\left[ \sup_{\theta \in [\alpha,\beta]^{\mathbf{d}}} |\mathscr{R}(\theta) - \mathbf{R}(\theta)|^p \right] \right)^{1/p}
$$
$$
\leq \inf_{C,\varepsilon \in (0,\infty)} \left[ \frac{2(C+1)b^2 \max\{1, [2\sqrt{M}L(\beta-\alpha)(Cb)^{-1}]^\varepsilon\} \sqrt{\max\{1,p,\mathbf{d}/\varepsilon\}}}{\sqrt{M}} \right]
$$
$$
\leq \inf_{C \in (0,\infty)} \left[ \frac{2(C+1)b^2 \sqrt{e \max\{1,p,\mathbf{d}\ln(4ML^2(\beta-\alpha)^2(Cb)^{-2})\}}}{\sqrt{M}} \right]. \tag{6.43}
$$

*Proof of Proposition 6.3.2.* Throughout this proof let $p \in (0,\infty)$, let $(\kappa_C)_{C \in (0,\infty)} \subseteq (0,\infty)$ satisfy for all $C \in (0,\infty)$ that $\kappa_C = {}^{2\sqrt{M}L(\beta-\alpha)}/_{(Cb)}$, let $\mathcal{X}_{\theta,j} \colon \Omega \to \mathbb{R}$, $j \in \{1,2,\dots,M\}$, $\theta \in [\alpha,\beta]^{\mathbf{d}}$, satisfy for all $\theta \in [\alpha,\beta]^{\mathbf{d}}$, $j \in \{1,2,\dots,M\}$ that $\mathcal{X}_{\theta,j} = f_\theta(X_j)$, and let $\delta \colon ([\alpha,\beta]^{\mathbf{d}}) \times ([\alpha,\beta]^{\mathbf{d}}) \to [0,\infty)$ satisfy for all $\theta, \vartheta \in [\alpha,\beta]^{\mathbf{d}}$ that $\delta(\theta,\vartheta) = \|\theta-\vartheta\|_\infty$. First of all, note that the assumption that $\forall\, \theta \in [\alpha,\beta]^{\mathbf{d}}$, $j \in \{1,2,\dots,M\} \colon |f_\theta(X_j) - Y_j| \leq b$ implies for all $\theta \in [\alpha,\beta]^{\mathbf{d}}$, $j \in \{1,2,\dots,M\}$ that

$$
|\mathcal{X}_{\theta,j} - Y_j| = |f_\theta(X_j) - Y_j| \leq b. \tag{6.44}
$$

In addition, the assumption that $\forall\, \theta, \vartheta \in [\alpha,\beta]^{\mathbf{d}}$, $x \in D \colon |f_\theta(x) - f_\vartheta(x)| \leq L\|\theta-\vartheta\|_\infty$ ensures for all $\theta, \vartheta \in [\alpha,\beta]^{\mathbf{d}}$, $j \in \{1,2,\dots,M\}$ that

$$
|\mathcal{X}_{\theta,j} - \mathcal{X}_{\vartheta,j}| = |f_\theta(X_j) - f_\vartheta(X_j)| \leq \sup_{x \in D} |f_\theta(x) - f_\vartheta(x)| \leq L\|\theta-\vartheta\|_\infty = L\delta(\theta,\vartheta). \tag{6.45}
$$

Combining this, (6.44), and the fact that for all $\theta \in [\alpha,\beta]^{\mathbf{d}}$ it holds that $(\mathcal{X}_{\theta,j}, Y_j)$, $j \in \{1,2,\dots,M\}$, are i.i.d. random variables with Lemma 6.3.1 (applied with $p \curvearrowright q$, $C \curvearrowright C$, $(E,\delta) \curvearrowright ([\alpha,\beta]^{\mathbf{d}},\delta)$, $(X_{x,j})_{x \in E, j \in \{1,2,\dots,M\}} \curvearrowright (\mathcal{X}_{\theta,j})_{\theta \in [\alpha,\beta]^{\mathbf{d}}, j \in \{1,2,\dots,M\}}$ for $q \in [2,\infty)$, $C \in (0,\infty)$ in the notation of Lemma 6.3.1) demonstrates that for all $C \in (0,\infty)$, $q \in [2,\infty)$ it holds that $\Omega \ni \omega \mapsto \sup_{\theta \in [\alpha,\beta]^{\mathbf{d}}} |\mathscr{R}(\theta,\omega) - \mathbf{R}(\theta)| \in [0,\infty]$ is $\mathcal{F}/\mathcal{B}([0,\infty])$-measurable and

$$
\left( \mathbb{E}\left[ \sup_{\theta \in [\alpha,\beta]^{\mathbf{d}}} |\mathscr{R}(\theta) - \mathbf{R}(\theta)|^q \right] \right)^{1/q} \leq \left( \mathcal{C}^{([\alpha,\beta]^{\mathbf{d}},\delta), \frac{Cb\sqrt{q-1}}{2L\sqrt{M}}} \right)^{1/q} \left[ \frac{2(C+1)b^2\sqrt{q-1}}{\sqrt{M}} \right] \tag{6.46}
$$

(cf. Definition 3.2.13). This finishes the proof of item (i). Next observe that item (ii) in Lemma 3.2.14 (applied with $d \curvearrowright \mathbf{d}$, $a \curvearrowright \alpha$, $b \curvearrowright \beta$, $r \curvearrowright r$ for $r \in (0,\infty)$ in the notation of Lemma 3.2.14) shows that for all $r \in (0,\infty)$ it holds that

$$
\begin{aligned}
\mathcal{C}^{([\alpha,\beta]^{\mathbf{d}},\delta),r} &\leq \mathbb{1}_{[0,r]}\left(\tfrac{\beta-\alpha}{2}\right) + \left(\tfrac{\beta-\alpha}{r}\right)^{\mathbf{d}} \mathbb{1}_{(r,\infty)}\left(\tfrac{\beta-\alpha}{2}\right) \\
&\leq \max\left\{1, \left(\tfrac{\beta-\alpha}{r}\right)^{\mathbf{d}}\right\} \left(\mathbb{1}_{[0,r]}\left(\tfrac{\beta-\alpha}{2}\right) + \mathbb{1}_{(r,\infty)}\left(\tfrac{\beta-\alpha}{2}\right)\right) \\
&= \max\left\{1, \left(\tfrac{\beta-\alpha}{r}\right)^{\mathbf{d}}\right\}.
\end{aligned} \tag{6.47}
$$

This yields for all $C \in (0,\infty)$, $q \in [2,\infty)$ that

$$
\begin{aligned}
\left( \mathcal{C}^{([\alpha,\beta]^{\mathbf{d}},\delta), \frac{Cb\sqrt{q-1}}{2L\sqrt{M}}} \right)^{1/q} &\leq \max\left\{1, \left(\tfrac{2(\beta-\alpha)L\sqrt{M}}{Cb\sqrt{q-1}}\right)^{\frac{\mathbf{d}}{q}}\right\} \\
&\leq \max\left\{1, \left(\tfrac{2(\beta-\alpha)L\sqrt{M}}{Cb}\right)^{\frac{\mathbf{d}}{q}}\right\} = \max\left\{1, (\kappa_C)^{\frac{\mathbf{d}}{q}}\right\}.
\end{aligned} \tag{6.48}
$$

Jensen's inequality and (6.46) hence prove that for all $C, \varepsilon \in (0, \infty)$ it holds that

$$
\begin{aligned}
&\left(\mathbb{E}\left[\sup_{\theta \in [\alpha,\beta]^{\mathbf{d}}} |\mathscr{R}(\theta) - \mathbf{R}(\theta)|^p\right]\right)^{1/p} \\
&\leq \left(\mathbb{E}\left[\sup_{\theta \in [\alpha,\beta]^{\mathbf{d}}} |\mathscr{R}(\theta) - \mathbf{R}(\theta)|^{\max\{2,p,\mathbf{d}/\varepsilon\}}\right]\right)^{\frac{1}{\max\{2,p,\mathbf{d}/\varepsilon\}}} \\
&\leq \max\left\{1, (\kappa_C)^{\frac{\mathbf{d}}{\max\{2,p,\mathbf{d}/\varepsilon\}}}\right\} \frac{2(C+1)b^2\sqrt{\max\{2,p,\mathbf{d}/\varepsilon\}-1}}{\sqrt{M}} \\
&= \max\left\{1, (\kappa_C)^{\min\{\mathbf{d}/2,\mathbf{d}/p,\varepsilon\}}\right\} \frac{2(C+1)b^2\sqrt{\max\{1,p-1,\mathbf{d}/\varepsilon-1\}}}{\sqrt{M}} \\
&\leq \frac{2(C+1)b^2\max\{1,(\kappa_C)^\varepsilon\}\sqrt{\max\{1,p,\mathbf{d}/\varepsilon\}}}{\sqrt{M}}.
\end{aligned}
\tag{6.49}
$$

Next note that the fact that $\forall\, a \in (1, \infty)\colon a^{1/(2\ln(a))} = e^{\ln(a)/(2\ln(a))} = e^{1/2} = \sqrt{e} \geq 1$ ensures that for all $C \in (0, \infty)$ with $\kappa_C > 1$ it holds that

$$
\begin{aligned}
&\inf_{\varepsilon \in (0,\infty)} \left[\frac{2(C+1)b^2\max\{1,(\kappa_C)^\varepsilon\}\sqrt{\max\{1,p,\mathbf{d}/\varepsilon\}}}{\sqrt{M}}\right] \\
&\leq \frac{2(C+1)b^2\max\{1,(\kappa_C)^{1/(2\ln(\kappa_C))}\}\sqrt{\max\{1,p,2\mathbf{d}\ln(\kappa_C)\}}}{\sqrt{M}} \\
&= \frac{2(C+1)b^2\sqrt{e\max\{1,p,\mathbf{d}\ln([\kappa_C]^2)\}}}{\sqrt{M}}.
\end{aligned}
\tag{6.50}
$$

In addition, observe that it holds for all $C \in (0, \infty)$ with $\kappa_C \leq 1$ that

$$
\begin{aligned}
&\inf_{\varepsilon \in (0,\infty)} \left[\frac{2(C+1)b^2\max\{1,(\kappa_C)^\varepsilon\}\sqrt{\max\{1,p,\mathbf{d}/\varepsilon\}}}{\sqrt{M}}\right] \\
&= \inf_{\varepsilon \in (0,\infty)} \left[\frac{2(C+1)b^2\sqrt{\max\{1,p,\mathbf{d}/\varepsilon\}}}{\sqrt{M}}\right] \leq \frac{2(C+1)b^2\sqrt{\max\{1,p\}}}{\sqrt{M}} \\
&\leq \frac{2(C+1)b^2\sqrt{e\max\{1,p,\mathbf{d}\ln([\kappa_C]^2)\}}}{\sqrt{M}}.
\end{aligned}
\tag{6.51}
$$

Combining (6.49) with (6.50) and (6.51) demonstrates that

$$
\begin{aligned}
&\left(\mathbb{E}\left[\sup_{\theta \in [\alpha,\beta]^{\mathbf{d}}} |\mathscr{R}(\theta) - \mathbf{R}(\theta)|^p\right]\right)^{1/p} \\
&\leq \inf_{C,\varepsilon \in (0,\infty)} \left[\frac{2(C+1)b^2\max\{1,(\kappa_C)^\varepsilon\}\sqrt{\max\{1,p,\mathbf{d}/\varepsilon\}}}{\sqrt{M}}\right] \\
&= \inf_{C,\varepsilon \in (0,\infty)} \left[\frac{2(C+1)b^2\max\{1,[2\sqrt{M}L(\beta-\alpha)(Cb)^{-1}]^\varepsilon\}\sqrt{\max\{1,p,\mathbf{d}/\varepsilon\}}}{\sqrt{M}}\right] \\
&\leq \inf_{C \in (0,\infty)} \left[\frac{2(C+1)b^2\sqrt{e\max\{1,p,\mathbf{d}\ln([\kappa_C]^2)\}}}{\sqrt{M}}\right] \\
&= \inf_{C \in (0,\infty)} \left[\frac{2(C+1)b^2\sqrt{e\max\{1,p,\mathbf{d}\ln(4ML^2(\beta-\alpha)^2(Cb)^{-2})\}}}{\sqrt{M}}\right].
\end{aligned}
\tag{6.52}
$$

This establishes item (ii) and thus completes the proof of Proposition 6.3.2. $\qquad\square$

**Corollary 6.3.3.** *Let $d, \mathbf{d}, \mathbf{L}, M \in \mathbb{N}$, $B, b \in [1, \infty)$, $u \in \mathbb{R}$, $v \in [u + 1, \infty)$, $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \ldots, \mathbf{l_L}) \in \mathbb{N}^{\mathbf{L}+1}$, $D \subseteq [-b, b]^d$, assume $\mathbf{l}_0 = d$, $\mathbf{l_L} = 1$, and $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1} + 1)$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $\mathbb{X}_j = (X_j, Y_j) \colon \Omega \to (D \times [u, v])$, $j \in \{1, 2, \ldots, M\}$, be i.i.d. random variables, let $\mathbf{R} \colon [-B, B]^{\mathbf{d}} \to [0, \infty)$ satisfy for all $\theta \in [-B, B]^{\mathbf{d}}$ that $\mathbf{R}(\theta) = \mathbb{E}[|\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_1) - Y_1|^2]$, and let $\mathscr{R} \colon [-B, B]^{\mathbf{d}} \times \Omega \to [0, \infty)$ satisfy for all $\theta \in [-B, B]^{\mathbf{d}}$, $\omega \in \Omega$ that*

$$\mathscr{R}(\theta, \omega) = \frac{1}{M}\left[\sum_{j=1}^{M} |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)|^2\right] \tag{6.53}$$

*(cf. Definition 2.1.27). Then*

(i) *it holds that $\Omega \ni \omega \mapsto \sup_{\theta \in [-B,B]^{\mathbf{d}}} |\mathscr{R}(\theta, \omega) - \mathbf{R}(\theta)| \in [0, \infty]$ is $\mathcal{F}/\mathcal{B}([0, \infty])$-measurable and*

(ii) *it holds for all $p \in (0, \infty)$ that*

$$\begin{aligned}
&\left(\mathbb{E}\left[\sup_{\theta \in [-B,B]^{\mathbf{d}}} |\mathscr{R}(\theta) - \mathbf{R}(\theta)|^p\right]\right)^{1/p} \\
&\leq \frac{9(v-u)^2 \mathbf{L}(\|\mathbf{l}\|_\infty + 1)\sqrt{\max\{p, \ln(4(Mb)^{1/\mathbf{L}}(\|\mathbf{l}\|_\infty + 1)B)\}}}{\sqrt{M}} \\
&\leq \frac{9(v-u)^2 \mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2 \max\{p, \ln(3MBb)\}}{\sqrt{M}}
\end{aligned} \tag{6.54}$$

*(cf. Definition 3.1.16).*

*Proof of Corollary 6.3.3.* Throughout this proof let $\mathfrak{d} \in \mathbb{N}$ be given by $\mathfrak{d} = \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1} + 1)$, let $L \in (0, \infty)$ be given by $L = b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}} B^{\mathbf{L}-1}$, let $f = (f_\theta)_{\theta \in [-B,B]^{\mathfrak{d}}} \colon [-B, B]^{\mathfrak{d}} \to C(D, \mathbb{R})$ satisfy for all $\theta \in [-B, B]^{\mathfrak{d}}$, $x \in D$ that $f_\theta(x) = \mathscr{N}_{u,v}^{\theta,\mathbf{l}}(x)$, let $\mathscr{R} \colon [-B, B]^{\mathfrak{d}} \to [0, \infty)$ satisfy for all $\theta \in [-B, B]^{\mathfrak{d}}$ that $\mathscr{R}(\theta) = \mathbb{E}[|f_\theta(X_1) - Y_1|^2] = \mathbb{E}[|\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_1) - Y_1|^2]$, and let $R \colon [-B, B]^{\mathfrak{d}} \times \Omega \to [0, \infty)$ satisfy for all $\theta \in [-B, B]^{\mathfrak{d}}$, $\omega \in \Omega$ that

$$R(\theta, \omega) = \frac{1}{M}\left[\sum_{j=1}^{M} |f_\theta(X_j(\omega)) - Y_j(\omega)|^2\right] = \frac{1}{M}\left[\sum_{j=1}^{M} |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)|^2\right]. \tag{6.55}$$

Note that the fact that $\forall\, \theta \in \mathbb{R}^{\mathfrak{d}}$, $x \in \mathbb{R}^d \colon \mathscr{N}_{u,v}^{\theta,\mathbf{l}}(x) \in [u, v]$ and the assumption that $\forall\, j \in \{1, 2, \ldots, M\} \colon Y_j(\Omega) \subseteq [u, v]$ imply for all $\theta \in [-B, B]^{\mathfrak{d}}$, $j \in \{1, 2, \ldots, M\}$ that

$$|f_\theta(X_j) - Y_j| = |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j) - Y_j| \leq \sup_{y_1, y_2 \in [u,v]} |y_1 - y_2| = v - u. \tag{6.56}$$

Moreover, the assumptions that $D \subseteq [-b, b]^d$, $\mathbf{l}_0 = d$, and $\mathbf{l_L} = 1$, Corollary 5.3.7 (applied with $a \curvearrowleft -b$, $b \curvearrowleft b$, $u \curvearrowleft u$, $v \curvearrowleft v$, $d \curvearrowleft \mathfrak{d}$, $L \curvearrowleft \mathbf{L}$, $l \curvearrowleft \mathbf{l}$ in the notation of Corollary 5.3.7), and the assumptions that $b \geq 1$ and $B \geq 1$ ensure that for all $\theta, \vartheta \in [-B, B]^{\mathfrak{d}}$, $x \in D$ it holds that

$$\begin{aligned}
|f_\theta(x) - f_\vartheta(x)| &\leq \sup_{y \in [-b,b]^d} |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(y) - \mathscr{N}_{u,v}^{\vartheta,\mathbf{l}}(y)| \\
&\leq \mathbf{L}\max\{1, b\}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}}(\max\{1, \|\theta\|_\infty, \|\vartheta\|_\infty\})^{\mathbf{L}-1}\|\theta - \vartheta\|_\infty \tag{6.57} \\
&\leq b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}} B^{\mathbf{L}-1}\|\theta - \vartheta\|_\infty = L\|\theta - \vartheta\|_\infty.
\end{aligned}$$

Furthermore, the facts that $\mathbf{d} \geq \mathfrak{d}$ and $\forall\, \theta = (\theta_1, \theta_2, \ldots, \theta_{\mathbf{d}}) \in \mathbb{R}^{\mathbf{d}} \colon \mathscr{N}_{u,v}^{\theta,\mathbf{l}} = \mathscr{N}_{u,v}^{(\theta_1, \theta_2, \ldots, \theta_{\mathfrak{d}}),\mathbf{l}}$ prove that for all $\omega \in \Omega$ it holds that

$$\sup_{\theta \in [-B,B]^{\mathbf{d}}} |\mathscr{R}(\theta, \omega) - \mathbf{R}(\theta)| = \sup_{\theta \in [-B,B]^{\mathfrak{d}}} |R(\theta, \omega) - \mathscr{R}(\theta)|. \tag{6.58}$$

Next observe that (6.56), (6.57), Proposition 6.3.2 (applied with $\mathbf{d} \curvearrowleft \mathfrak{d}$, $b \curvearrowleft v - u$, $\alpha \curvearrowleft -B$, $\beta \curvearrowleft B$, $\mathbf{R} \curvearrowleft \mathscr{R}$, $\mathscr{R} \curvearrowleft R$ in the notation of Proposition 6.3.2), and the facts that $v - u \geq (u + 1) - u = 1$ and $\mathfrak{d} \leq \mathbf{L}\|\mathbf{l}\|_\infty(\|\mathbf{l}\|_\infty + 1) \leq \mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2$ demonstrate that for all $p \in (0, \infty)$ it holds that $\Omega \ni \omega \mapsto \sup_{\theta \in [-B,B]^\mathfrak{d}} |R(\theta, \omega) - \mathscr{R}(\theta)| \in [0, \infty]$ is $\mathcal{F}/\mathcal{B}([0, \infty])$-measurable and

$$
\begin{aligned}
&\left(\mathbb{E}\left[\sup_{\theta \in [-B,B]^\mathfrak{d}} |R(\theta) - \mathscr{R}(\theta)|^p\right]\right)^{1/p} \\
&\leq \inf_{C \in (0,\infty)} \left[\frac{2(C+1)(v-u)^2\sqrt{e\max\{1, p, \mathfrak{d}\ln(4ML^2(2B)^2(C[v-u])^{-2})\}}}{\sqrt{M}}\right] \\
&\leq \inf_{C \in (0,\infty)} \left[\frac{2(C+1)(v-u)^2\sqrt{e\max\{1, p, \mathbf{L}(\|\mathbf{l}\|_\infty+1)^2\ln(2^4ML^2B^2C^{-2})\}}}{\sqrt{M}}\right].
\end{aligned}
\tag{6.59}
$$

This and (6.58) establish item (i). In addition, combining (6.58)–(6.59) with the fact that $2^6\mathbf{L}^2 \leq 2^6 \cdot 2^{2(\mathbf{L}-1)} = 2^{4+2\mathbf{L}} \leq 2^{4\mathbf{L}+2\mathbf{L}} = 2^{6\mathbf{L}}$ and the facts that $3 \geq e$, $B \geq 1$, $\mathbf{L} \geq 1$, $M \geq 1$, and $b \geq 1$ shows that for all $p \in (0, \infty)$ it holds that

$$
\begin{aligned}
&\left(\mathbb{E}\left[\sup_{\theta \in [-B,B]^\mathbf{d}} |\mathscr{R}(\theta) - \mathbf{R}(\theta)|^p\right]\right)^{1/p} = \left(\mathbb{E}\left[\sup_{\theta \in [-B,B]^\mathfrak{d}} |R(\theta) - \mathscr{R}(\theta)|^p\right]\right)^{1/p} \\
&\leq \frac{2(1/2 + 1)(v-u)^2\sqrt{e\max\{1, p, \mathbf{L}(\|\mathbf{l}\|_\infty+1)^2\ln(2^4ML^2B^22^2)\}}}{\sqrt{M}} \\
&= \frac{3(v-u)^2\sqrt{e\max\{p, \mathbf{L}(\|\mathbf{l}\|_\infty+1)^2\ln(2^6Mb^2\mathbf{L}^2(\|\mathbf{l}\|_\infty+1)^{2\mathbf{L}}B^{2\mathbf{L}})\}}}{\sqrt{M}} \\
&\leq \frac{3(v-u)^2\sqrt{e\max\{p, 3\mathbf{L}^2(\|\mathbf{l}\|_\infty+1)^2\ln([2^{6\mathbf{L}}Mb^2(\|\mathbf{l}\|_\infty+1)^{2\mathbf{L}}B^{2\mathbf{L}}]^{1/(3\mathbf{L})})\}}}{\sqrt{M}} \\
&\leq \frac{3(v-u)^2\sqrt{3\max\{p, 3\mathbf{L}^2(\|\mathbf{l}\|_\infty+1)^2\ln(2^2(Mb^2)^{1/(3\mathbf{L})}(\|\mathbf{l}\|_\infty+1)B)\}}}{\sqrt{M}} \\
&\leq \frac{9(v-u)^2\mathbf{L}(\|\mathbf{l}\|_\infty+1)\sqrt{\max\{p, \ln(4(Mb)^{1/\mathbf{L}}(\|\mathbf{l}\|_\infty+1)B)\}}}{\sqrt{M}}.
\end{aligned}
\tag{6.60}
$$

Furthermore, note that the fact that $\forall\, n \in \mathbb{N} \colon n \leq 2^{n-1}$ and the fact that $\|\mathbf{l}\|_\infty \geq 1$ imply that

$$
4(\|\mathbf{l}\|_\infty + 1) \leq 2^2 \cdot 2^{(\|\mathbf{l}\|_\infty+1)-1} = 2^3 \cdot 2^{(\|\mathbf{l}\|_\infty+1)-2} \leq 3^2 \cdot 3^{(\|\mathbf{l}\|_\infty+1)-2} = 3^{(\|\mathbf{l}\|_\infty+1)}.
\tag{6.61}
$$

This demonstrates for all $p \in (0, \infty)$ that

$$
\begin{aligned}
&\frac{9(v-u)^2\mathbf{L}(\|\mathbf{l}\|_\infty+1)\sqrt{\max\{p, \ln(4(Mb)^{1/\mathbf{L}}(\|\mathbf{l}\|_\infty+1)B)\}}}{\sqrt{M}} \\
&\leq \frac{9(v-u)^2\mathbf{L}(\|\mathbf{l}\|_\infty+1)\sqrt{\max\{p, (\|\mathbf{l}\|_\infty+1)\ln([3^{(\|\mathbf{l}\|_\infty+1)}(Mb)^{1/\mathbf{L}}B]^{1/(\|\mathbf{l}\|_\infty+1)})\}}}{\sqrt{M}} \\
&\leq \frac{9(v-u)^2\mathbf{L}(\|\mathbf{l}\|_\infty+1)^2\max\{p, \ln(3MBb)\}}{\sqrt{M}}.
\end{aligned}
\tag{6.62}
$$

Combining this with (6.60) shows item (ii). The proof of Corollary 6.3.3 is thus complete. $\qquad\square$

# Chapter 7

# Analysis of the overall error

## 7.1 Full strong error analysis for the training of ANNs

**Lemma 7.1.1.** *Let* $d, \mathbf{d}, \mathbf{L} \in \mathbb{N}$, $p \in [0, \infty)$, $u \in [-\infty, \infty)$, $v \in (u, \infty]$, $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \ldots, \mathbf{l}_\mathbf{L}) \in \mathbb{N}^{\mathbf{L}+1}$, $D \subseteq \mathbb{R}^d$, *assume* $\mathbf{l}_0 = d$, $\mathbf{l}_\mathbf{L} = 1$, *and* $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1} + 1)$, *let* $\mathcal{E} \colon D \to \mathbb{R}$ *be* $\mathcal{B}(D)/\mathcal{B}(\mathbb{R})$-*measurable, let* $(\Omega, \mathcal{F}, \mathbb{P})$ *be a probability space, and let* $X \colon \Omega \to D$, $\mathbf{k} \colon \Omega \to (\mathbb{N}_0)^2$, *and* $\Theta_{k,n} \colon \Omega \to \mathbb{R}^\mathbf{d}$, $k, n \in \mathbb{N}_0$, *be random variables. Then*

    *(i) it holds that* $\mathbb{R}^\mathbf{d} \times \mathbb{R}^d \ni (\theta, x) \mapsto \mathcal{N}_{u,v}^{\theta, \mathbf{l}}(x) \in \mathbb{R}$ *is* $(\mathcal{B}(\mathbb{R}^\mathbf{d}) \otimes \mathcal{B}(\mathbb{R}^d))/\mathcal{B}(\mathbb{R})$-*measurable,*

    *(ii) it holds that* $\Omega \ni \omega \mapsto \Theta_{\mathbf{k}(\omega)}(\omega) \in \mathbb{R}^\mathbf{d}$ *is* $\mathcal{F}/\mathcal{B}(\mathbb{R}^\mathbf{d})$-*measurable, and*

    *(iii) it holds that*

$$\Omega \ni \omega \mapsto \int_D |\mathcal{N}_{u,v}^{\Theta_{\mathbf{k}(\omega)}(\omega), \mathbf{l}}(x) - \mathcal{E}(x)|^p \, \mathbb{P}_X(\mathrm{d}x) \in [0, \infty] \qquad (7.1)$$

    *is* $\mathcal{F}/\mathcal{B}([0, \infty])$-*measurable*

*(cf. Definition 2.1.27).*

*Proof of Lemma 7.1.1.* First, observe that Corollary 5.3.7 (applied with $a \curvearrowleft -\|x\|_\infty$, $b \curvearrowleft \|x\|_\infty$, $u \curvearrowleft u$, $v \curvearrowleft v$, $d \curvearrowleft \mathbf{d}$, $L \curvearrowleft \mathbf{L}$, $l \curvearrowleft \mathbf{l}$ for $x \in \mathbb{R}^d$ in the notation of Corollary 5.3.7) demonstrates that for all $x \in \mathbb{R}^d$, $\theta, \vartheta \in \mathbb{R}^\mathbf{d}$ it holds that

$$\begin{aligned} |\mathcal{N}_{u,v}^{\theta, \mathbf{l}}(x) - \mathcal{N}_{u,v}^{\vartheta, \mathbf{l}}(x)| &\leq \sup_{y \in [-\|x\|_\infty, \|x\|_\infty]^{\mathbf{l}_0}} |\mathcal{N}_{u,v}^{\theta, \mathbf{l}}(y) - \mathcal{N}_{u,v}^{\vartheta, \mathbf{l}}(y)| \\ &\leq \mathbf{L} \max\{1, \|x\|_\infty\}(\|\mathbf{l}\|_\infty + 1)^\mathbf{L}(\max\{1, \|\theta\|_\infty, \|\vartheta\|_\infty\})^{\mathbf{L}-1} \|\theta - \vartheta\|_\infty \end{aligned} \qquad (7.2)$$

(cf. Definition 3.1.16). This implies for all $x \in \mathbb{R}^d$ that

$$\mathbb{R}^\mathbf{d} \ni \theta \mapsto \mathcal{N}_{u,v}^{\theta, \mathbf{l}}(x) \in \mathbb{R} \qquad (7.3)$$

is continuous. In addition, the fact that $\forall\, \theta \in \mathbb{R}^\mathbf{d} \colon \mathcal{N}_{u,v}^{\theta, \mathbf{l}} \in C(\mathbb{R}^d, \mathbb{R})$ ensures that for all $\theta \in \mathbb{R}^\mathbf{d}$ it holds that $\mathbb{R}^d \ni x \mapsto \mathcal{N}_{u,v}^{\theta, \mathbf{l}}(x) \in \mathbb{R}$ is $\mathcal{B}(\mathbb{R}^d)/\mathcal{B}(\mathbb{R})$-measurable. This, (7.3), the fact that $(\mathbb{R}^\mathbf{d}, \|\cdot\|_\infty|_{\mathbb{R}^\mathbf{d}})$ is a separable normed $\mathbb{R}$-vector space, and Lemma 5.2.6 show item (i). Next we prove item (ii). For this let $\Xi \colon \Omega \to \mathbb{R}^\mathbf{d}$ satisfy for all $\omega \in \Omega$

that $\Xi(\omega) = \Theta_{\mathbf{k}(\omega)}(\omega)$. Observe that the assumption that $\Theta_{k,n}\colon \Omega \to \mathbb{R}^{\mathbf{d}}$, $k, n \in \mathbb{N}_0$, and $\mathbf{k}\colon \Omega \to (\mathbb{N}_0)^2$ are random variables establishes that for all $U \in \mathcal{B}(\mathbb{R}^{\mathbf{d}})$ it holds that

$$
\begin{aligned}
\Xi^{-1}(U) &= \{\omega \in \Omega\colon \Xi(\omega) \in U\} = \{\omega \in \Omega\colon \Theta_{\mathbf{k}(\omega)}(\omega) \in U\} \\
&= \{\omega \in \Omega\colon \left[\exists\, k, n \in \mathbb{N}_0\colon ([\Theta_{k,n}(\omega) \in U] \wedge [\mathbf{k}(\omega) = (k,n)])\right]\} \\
&= \bigcup_{k=0}^{\infty} \bigcup_{n=0}^{\infty} \left(\{\omega \in \Omega\colon \Theta_{k,n}(\omega) \in U\} \cap \{\omega \in \Omega\colon \mathbf{k}(\omega) = (k,n)\}\right) \qquad (7.4) \\
&= \bigcup_{k=0}^{\infty} \bigcup_{n=0}^{\infty} \left([(\Theta_{k,n})^{-1}(U)] \cap [\mathbf{k}^{-1}(\{(k,n)\})]\right) \in \mathcal{F}.
\end{aligned}
$$

This implies item (ii). Moreover, note that item (i)–item (ii) yield that $\Omega \times \mathbb{R}^d \ni (\omega, x) \mapsto \mathscr{N}_{u,v}^{\Theta_{\mathbf{k}(\omega)}(\omega),\mathbf{l}}(x) \in \mathbb{R}$ is $(\mathcal{F} \otimes \mathcal{B}(\mathbb{R}^d))/\mathcal{B}(\mathbb{R})$-measurable. This and the assumption that $\mathcal{E}\colon D \to \mathbb{R}$ is $\mathcal{B}(D)/\mathcal{B}(\mathbb{R})$-measurable demonstrate that $\Omega \times D \ni (\omega, x) \mapsto |\mathscr{N}_{u,v}^{\Theta_{\mathbf{k}(\omega)}(\omega),\mathbf{l}}(x) - \mathcal{E}(x)|^p \in [0,\infty)$ is $(\mathcal{F} \otimes \mathcal{B}(D))/\mathcal{B}([0,\infty))$-measurable. Tonelli's theorem hence establishes item (iii). The proof of Lemma 7.1.1 is thus complete. $\qquad\square$

**Proposition 7.1.2.** *Let* $d, \mathbf{d}, \mathbf{L}, M, K, N \in \mathbb{N}$, $b, c \in [1, \infty)$, $B \in [c, \infty)$, $u \in \mathbb{R}$, $v \in (u, \infty)$, $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \ldots, \mathbf{l}_{\mathbf{L}}) \in \mathbb{N}^{\mathbf{L}+1}$, $\mathbf{T} \subseteq \{0, 1, \ldots, N\}$, $D \subseteq [-b, b]^d$, *assume* $0 \in \mathbf{T}$, $\mathbf{l}_0 = d$, $\mathbf{l}_{\mathbf{L}} = 1$, *and* $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1} + 1)$, *let* $(\Omega, \mathcal{F}, \mathbb{P})$ *be a probability space, let* $X_j\colon \Omega \to D$, $j \in \mathbb{N}$, *and* $Y_j\colon \Omega \to [u, v]$, $j \in \mathbb{N}$, *be functions, assume that* $(X_j, Y_j)$, $j \in \{1, 2, \ldots, M\}$, *are i.i.d. random variables, let* $\mathcal{E}\colon D \to [u, v]$ *be* $\mathcal{B}(D)/\mathcal{B}([u, v])$-*measurable, assume that it holds* $\mathbb{P}$-*a.s. that* $\mathcal{E}(X_1) = \mathbb{E}[Y_1|X_1]$, *let* $\Theta_{k,n}\colon \Omega \to \mathbb{R}^{\mathbf{d}}$, $k, n \in \mathbb{N}_0$, *and* $\mathbf{k}\colon \Omega \to (\mathbb{N}_0)^2$ *be random variables, assume* $\left(\bigcup_{k=1}^{\infty} \Theta_{k,0}(\Omega)\right) \subseteq [-B, B]^{\mathbf{d}}$, *assume that* $\Theta_{k,0}$, $k \in \{1, 2, \ldots, K\}$, *are i.i.d., assume that* $\Theta_{1,0}$ *is continuous uniformly distributed on* $[-c, c]^{\mathbf{d}}$, *and let* $\mathscr{R}\colon \mathbb{R}^{\mathbf{d}} \times \Omega \to [0, \infty)$ *satisfy for all* $\theta \in \mathbb{R}^{\mathbf{d}}$, $\omega \in \Omega$ *that*

$$
\mathscr{R}(\theta, \omega) = \frac{1}{M}\left[\sum_{j=1}^{M} |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)|^2\right] \qquad and \qquad (7.5)
$$

$$
\mathbf{k}(\omega) \in \operatorname{argmin}_{(k,n)\in\{1,2,\ldots,K\}\times\mathbf{T},\, \|\Theta_{k,n}(\omega)\|_\infty \leq B} \mathscr{R}(\Theta_{k,n}(\omega), \omega) \qquad (7.6)
$$

*(cf. Definitions 2.1.27 and 3.1.16). Then it holds for all* $p \in (0, \infty)$ *that*

$$
\left(\mathbb{E}\left[\left(\int_D |\mathscr{N}_{u,v}^{\Theta_{\mathbf{k}},\mathbf{l}}(x) - \mathcal{E}(x)|^2\, \mathbb{P}_{X_1}(\mathrm{d}x)\right)^p\right]\right)^{1/p}
$$

$$
\begin{aligned}
&\leq \left[\inf_{\theta\in[-c,c]^{\mathbf{d}}} \sup_{x\in D} |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(x) - \mathcal{E}(x)|^2\right] + \frac{4(v-u)b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}} c^{\mathbf{L}} \max\{1, p\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}} \\
&\quad + \frac{18 \max\{1, (v-u)^2\}\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2 \max\{p, \ln(3MBb)\}}{\sqrt{M}} \\
&\leq \left[\inf_{\theta\in[-c,c]^{\mathbf{d}}} \sup_{x\in D} |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(x) - \mathcal{E}(x)|^2\right] \\
&\quad + \frac{20 \max\{1, (v-u)^2\}b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}+1} B^{\mathbf{L}} \max\{p, \ln(3M)\}}{\min\{\sqrt{M}, K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}\}}
\end{aligned} \qquad (7.7)
$$

*(cf. item (iii) in Lemma 7.1.1).*

*Proof of Proposition 7.1.2.* Throughout this proof let $\mathbf{R}\colon \mathbb{R}^{\mathbf{d}} \to [0, \infty)$ satisfy for all $\theta \in \mathbb{R}^{\mathbf{d}}$ that $\mathbf{R}(\theta) = \mathbb{E}[|\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_1) - Y_1|^2]$. First of all, observe that the assumption that $\left(\bigcup_{k=1}^{\infty} \Theta_{k,0}(\Omega)\right) \subseteq [-B, B]^{\mathbf{d}}$, the assumption that $0 \in \mathbf{T}$, and Proposition 4.2.1 show that

for all $\vartheta \in [-B, B]^{\mathbf{d}}$ it holds that

$$
\begin{aligned}
\int_D & |\mathcal{N}_{u,v}^{\Theta_{\mathbf{k}},\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) \\
&\le \left[\sup_{x\in D}|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)|^2\right] + 2\left[\sup_{\theta\in[-B,B]^{\mathbf{d}}}|\mathscr{R}(\theta) - \mathbf{R}(\theta)|\right] \\
&\quad + \min_{(k,n)\in\{1,2,\dots,K\}\times\mathbf{T}, \|\Theta_{k,n}\|_\infty\le B}|\mathscr{R}(\Theta_{k,n}) - \mathscr{R}(\vartheta)| \\
&\le \left[\sup_{x\in D}|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)|^2\right] + 2\left[\sup_{\theta\in[-B,B]^{\mathbf{d}}}|\mathscr{R}(\theta) - \mathbf{R}(\theta)|\right] \\
&\quad + \min_{k\in\{1,2,\dots,K\}, \|\Theta_{k,0}\|_\infty\le B}|\mathscr{R}(\Theta_{k,0}) - \mathscr{R}(\vartheta)| \\
&= \left[\sup_{x\in D}|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)|^2\right] + 2\left[\sup_{\theta\in[-B,B]^{\mathbf{d}}}|\mathscr{R}(\theta) - \mathbf{R}(\theta)|\right] \\
&\quad + \min_{k\in\{1,2,\dots,K\}}|\mathscr{R}(\Theta_{k,0}) - \mathscr{R}(\vartheta)|.
\end{aligned}
\tag{7.8}
$$

Minkowski's inequality hence establishes that for all $p \in [1, \infty)$, $\vartheta \in [-c, c]^{\mathbf{d}} \subseteq [-B, B]^{\mathbf{d}}$ it holds that

$$
\begin{aligned}
&\left(\mathbb{E}\left[\left(\int_D|\mathcal{N}_{u,v}^{\Theta_{\mathbf{k}},\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x)\right)^p\right]\right)^{1/p} \\
&\le \left(\mathbb{E}\left[\sup_{x\in D}|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)|^{2p}\right]\right)^{1/p} + 2\left(\mathbb{E}\left[\sup_{\theta\in[-B,B]^{\mathbf{d}}}|\mathscr{R}(\theta) - \mathbf{R}(\theta)|^p\right]\right)^{1/p} \\
&\quad + \left(\mathbb{E}\left[\min_{k\in\{1,2,\dots,K\}}|\mathscr{R}(\Theta_{k,0}) - \mathscr{R}(\vartheta)|^p\right]\right)^{1/p} \\
&\le \left[\sup_{x\in D}|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)|^2\right] + 2\left(\mathbb{E}\left[\sup_{\theta\in[-B,B]^{\mathbf{d}}}|\mathscr{R}(\theta) - \mathbf{R}(\theta)|^p\right]\right)^{1/p} \\
&\quad + \sup_{\theta\in[-c,c]^{\mathbf{d}}}\left(\mathbb{E}\left[\min_{k\in\{1,2,\dots,K\}}|\mathscr{R}(\Theta_{k,0}) - \mathscr{R}(\theta)|^p\right]\right)^{1/p}
\end{aligned}
\tag{7.9}
$$

(cf. item (i) in Corollary 6.3.3 and item (i) in Corollary 5.3.9). Next note that Corollary 6.3.3 (applied with $v \curvearrowleft \max\{u+1, v\}$, $\mathbf{R} \curvearrowleft \mathbf{R}|_{[-B,B]^{\mathbf{d}}}$, $\mathscr{R} \curvearrowleft \mathscr{R}|_{[-B,B]^{\mathbf{d}}\times\Omega}$ in the notation of Corollary 6.3.3) proves that for all $p \in (0, \infty)$ it holds that

$$
\begin{aligned}
&\left(\mathbb{E}\left[\sup_{\theta\in[-B,B]^{\mathbf{d}}}|\mathscr{R}(\theta) - \mathbf{R}(\theta)|^p\right]\right)^{1/p} \\
&\le \frac{9(\max\{u+1,v\} - u)^2\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2 \max\{p, \ln(3MBb)\}}{\sqrt{M}} \\
&= \frac{9\max\{1, (v-u)^2\}\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2 \max\{p, \ln(3MBb)\}}{\sqrt{M}}.
\end{aligned}
\tag{7.10}
$$

In addition, observe that Corollary 5.3.9 (applied with $\mathfrak{d} \curvearrowleft \sum_{i=1}^{\mathbf{L}}\mathbf{l}_i(\mathbf{l}_{i-1} + 1)$, $B \curvearrowleft c$, $(\Theta_k)_{k\in\{1,2,\dots,K\}} \curvearrowleft (\Omega \ni \omega \mapsto \mathbb{1}_{\{\Theta_{k,0}\in[-c,c]^{\mathbf{d}}\}}(\omega)\Theta_{k,0}(\omega) \in [-c,c]^{\mathbf{d}})_{k\in\{1,2,\dots,K\}}$, $\mathscr{R} \curvearrowleft \mathscr{R}|_{[-c,c]^{\mathbf{d}}\times\Omega}$ in the notation of Corollary 5.3.9) implies that for all $p \in (0, \infty)$ it holds that

$$
\begin{aligned}
&\sup_{\theta\in[-c,c]^{\mathbf{d}}}\left(\mathbb{E}\left[\min_{k\in\{1,2,\dots,K\}}|\mathscr{R}(\Theta_{k,0}) - \mathscr{R}(\theta)|^p\right]\right)^{1/p} \\
&= \sup_{\theta\in[-c,c]^{\mathbf{d}}}\left(\mathbb{E}\left[\min_{k\in\{1,2,\dots,K\}}|\mathscr{R}(\mathbb{1}_{\{\Theta_{k,0}\in[-c,c]^{\mathbf{d}}\}}\Theta_{k,0}) - \mathscr{R}(\theta)|^p\right]\right)^{1/p} \\
&\le \frac{4(v-u)b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}}c^{\mathbf{L}}\max\{1, p\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}}.
\end{aligned}
\tag{7.11}
$$

Combining this, (7.9), (7.10), and the fact that $\ln(3MBb) \ge 1$ with Jensen's inequality

demonstrates that for all $p \in (0, \infty)$ it holds that

$$\left(\mathbb{E}\left[\left(\int_D |\mathcal{N}_{u,v}^{\Theta_{\mathbf{k}},\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x)\right)^p\right]\right)^{1/p}$$

$$\leq \left(\mathbb{E}\left[\left(\int_D |\mathcal{N}_{u,v}^{\Theta_{\mathbf{k}},\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x)\right)^{\max\{1,p\}}\right]\right)^{\frac{1}{\max\{1,p\}}}$$

$$\leq \left[\inf_{\theta \in [-c,c]^{\mathbf{d}}} \sup_{x \in D} |\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(x) - \mathcal{E}(x)|^2\right]$$

$$\quad + \sup_{\theta \in [-c,c]^{\mathbf{d}}} \left(\mathbb{E}\left[\min_{k \in \{1,2,\ldots,K\}} |\mathcal{R}(\Theta_{k,0}) - \mathcal{R}(\theta)|^{\max\{1,p\}}\right]\right)^{\frac{1}{\max\{1,p\}}} \tag{7.12}$$

$$\quad + 2\left(\mathbb{E}\left[\sup_{\theta \in [-B,B]^{\mathbf{d}}} |\mathcal{R}(\theta) - \mathbf{R}(\theta)|^{\max\{1,p\}}\right]\right)^{\frac{1}{\max\{1,p\}}}$$

$$\leq \left[\inf_{\theta \in [-c,c]^{\mathbf{d}}} \sup_{x \in D} |\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(x) - \mathcal{E}(x)|^2\right] + \frac{4(v-u)b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}} c^{\mathbf{L}} \max\{1,p\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}}$$

$$\quad + \frac{18 \max\{1, (v-u)^2\}\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2 \max\{p, \ln(3MBb)\}}{\sqrt{M}}.$$

Moreover, note that the fact that $\forall\, x \in [0, \infty)\colon x + 1 \leq e^x \leq 3^x$ and the facts that $Bb \geq 1$ and $M \geq 1$ ensure that

$$\ln(3MBb) \leq \ln(3M3^{Bb-1}) = \ln(3^{Bb}M) = Bb\ln([3^{Bb}M]^{1/(Bb)}) \leq Bb\ln(3M). \tag{7.13}$$

The facts that $\|\mathbf{l}\|_\infty + 1 \geq 2$, $B \geq c \geq 1$, $\ln(3M) \geq 1$, $b \geq 1$, and $\mathbf{L} \geq 1$ hence show that for all $p \in (0, \infty)$ it holds that

$$\frac{4(v-u)b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}} c^{\mathbf{L}} \max\{1,p\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}}$$

$$\quad + \frac{18 \max\{1, (v-u)^2\}\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2 \max\{p, \ln(3MBb)\}}{\sqrt{M}}$$

$$\leq \frac{2(\|\mathbf{l}\|_\infty + 1) \max\{1, (v-u)^2\} b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}} B^{\mathbf{L}} \max\{p, \ln(3M)\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}} \tag{7.14}$$

$$\quad + \frac{18 \max\{1, (v-u)^2\} b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2 B \max\{p, \ln(3M)\}}{\sqrt{M}}$$

$$\leq \frac{20 \max\{1, (v-u)^2\} b\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}+1} B^{\mathbf{L}} \max\{p, \ln(3M)\}}{\min\{\sqrt{M}, K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}\}}.$$

This and (7.12) complete the proof of Proposition 7.1.2. $\qquad\square$

**Lemma 7.1.3.** *Let $a, x, p \in (0, \infty)$, $M, c \in [1, \infty)$, $B \in [c, \infty)$. Then*

*(i) it holds that $ax^p \leq \exp\!\left(a^{1/p}\frac{px}{e}\right)$ and*

*(ii) it holds that $\ln(3MBc) \leq \frac{23B}{18}\ln(eM)$.*

*Proof of Lemma 7.1.3.* First, note that the fact that $\forall\, y \in \mathbb{R}\colon y + 1 \leq e^y$ demonstrates that

$$ax^p = (a^{1/p}x)^p = \left[e\!\left(a^{1/p}\tfrac{x}{e} - 1 + 1\right)\right]^p \leq \left[e\exp\!\left(a^{1/p}\tfrac{x}{e} - 1\right)\right]^p = \exp\!\left(a^{1/p}\tfrac{px}{e}\right). \tag{7.15}$$

This proves item (i).

Second, observe that item (i) and the fact that $2\sqrt{3}/e \leq 23/18$ ensure that

$$3B^2 \leq \exp\!\left(\sqrt{3}\tfrac{2B}{e}\right) = \exp\!\left(\tfrac{2\sqrt{3}B}{e}\right) \leq \exp\!\left(\tfrac{23B}{18}\right). \tag{7.16}$$

The facts that $B \geq c \geq 1$ and $M \geq 1$ hence imply that

$$\ln(3MBc) \leq \ln(3B^2M) \leq \ln([eM]^{23B/18}) = \tfrac{23B}{18}\ln(eM). \qquad (7.17)$$

This establishes item (ii). The proof of Lemma 7.1.3 is thus complete. $\qquad \square$

**Theorem 7.1.4.** *Let* $d, \mathbf{d}, \mathbf{L}, M, K, N \in \mathbb{N}$, $A \in (0, \infty)$, $L, a, u \in \mathbb{R}$, $b \in (a, \infty)$, $v \in (u, \infty)$, $c \in [\max\{1, L, |a|, |b|, 2|u|, 2|v|\}, \infty)$, $B \in [c, \infty)$, $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \dots, \mathbf{l_L}) \in \mathbb{N}^{\mathbf{L}+1}$, $\mathbf{T} \subseteq \{0, 1, \dots, N\}$, *assume* $0 \in \mathbf{T}$, $\mathbf{L} \geq {}^{A\mathbb{1}_{(6^d, \infty)}(A)}/_{(2d)} + 1$, $\mathbf{l}_0 = d$, $\mathbf{l}_1 \geq A\mathbb{1}_{(6^d, \infty)}(A)$, $\mathbf{l_L} = 1$, *and* $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1} + 1)$, *assume for all* $i \in \{2, 3, \dots\} \cap [0, \mathbf{L})$ *that* $\mathbf{l}_i \geq \mathbb{1}_{(6^d, \infty)}(A)\max\{A/d - 2i + 3, 2\}$, *let* $(\Omega, \mathcal{F}, \mathbb{P})$ *be a probability space, let* $X_j \colon \Omega \to [a, b]^d$, $j \in \mathbb{N}$, *and* $Y_j \colon \Omega \to [u, v]$, $j \in \mathbb{N}$, *be functions, assume that* $(X_j, Y_j)$, $j \in \{1, 2, \dots, M\}$, *are i.i.d. random variables, let* $\mathcal{E} \colon [a, b]^d \to [u, v]$ *satisfy* $\mathbb{P}$*-a.s. that* $\mathcal{E}(X_1) = \mathbb{E}[Y_1 | X_1]$, *assume for all* $x, y \in [a, b]^d$ *that* $|\mathcal{E}(x) - \mathcal{E}(y)| \leq L\|x - y\|_1$, *let* $\Theta_{k,n} \colon \Omega \to \mathbb{R}^{\mathbf{d}}$, $k, n \in \mathbb{N}_0$, *and* $\mathbf{k} \colon \Omega \to (\mathbb{N}_0)^2$ *be random variables, assume* $\left(\bigcup_{k=1}^{\infty} \Theta_{k,0}(\Omega)\right) \subseteq [-B, B]^{\mathbf{d}}$, *assume that* $\Theta_{k,0}$, $k \in \{1, 2, \dots, K\}$, *are i.i.d., assume that* $\Theta_{1,0}$ *is continuous uniformly distributed on* $[-c, c]^{\mathbf{d}}$, *and let* $\mathscr{R} \colon \mathbb{R}^{\mathbf{d}} \times \Omega \to [0, \infty)$ *satisfy for all* $\theta \in \mathbb{R}^{\mathbf{d}}$, $\omega \in \Omega$ *that*

$$\mathscr{R}(\theta, \omega) = \frac{1}{M}\left[\sum_{j=1}^{M} |\mathscr{N}_{u,v}^{\theta, \mathbf{l}}(X_j(\omega)) - Y_j(\omega)|^2\right] \qquad and \qquad (7.18)$$

$$\mathbf{k}(\omega) \in \operatorname{argmin}_{(k,n) \in \{1,2,\dots,K\} \times \mathbf{T}, \|\Theta_{k,n}(\omega)\|_\infty \leq B} \mathscr{R}(\Theta_{k,n}(\omega), \omega) \qquad (7.19)$$

*(cf. Definitions 2.1.27 and 3.1.16). Then it holds for all* $p \in (0, \infty)$ *that*

$$
\begin{aligned}
&\left(\mathbb{E}\left[\left(\int_{[a,b]^d} |\mathscr{N}_{u,v}^{\Theta_{\mathbf{k}}, \mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x)\right)^p\right]\right)^{1/p} \\
&\leq \frac{9d^2 L^2 (b-a)^2}{A^{2/d}} + \frac{4(v-u)\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}} c^{\mathbf{L}+1} \max\{1, p\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}} \\
&\quad + \frac{18\max\{1, (v-u)^2\}\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2 \max\{p, \ln(3MBc)\}}{\sqrt{M}} \\
&\leq \frac{36d^2 c^4}{A^{2/d}} + \frac{4\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}} c^{\mathbf{L}+2} \max\{1, p\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}} + \frac{23B^3 \mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2 \max\{p, \ln(eM)\}}{\sqrt{M}}
\end{aligned}
$$
$$(7.20)$$

*(cf. item (iii) in Lemma 7.1.1).*

*Proof of Theorem 7.1.4.* First of all, note that the assumption that $\forall\, x, y \in [a, b]^d \colon |\mathcal{E}(x) - \mathcal{E}(y)| \leq L\|x - y\|_1$ ensures that $\mathcal{E} \colon [a, b]^d \to [u, v]$ is $\mathcal{B}([a, b]^d)/\mathcal{B}([u, v])$-measurable. The fact that $\max\{1, |a|, |b|\} \leq c$ and Proposition 7.1.2 (applied with $b \curvearrowleft \max\{1, |a|, |b|\}$, $D \curvearrowleft [a, b]^d$ in the notation of Proposition 7.1.2) hence show that for all $p \in (0, \infty)$ it

holds that

$$
\begin{aligned}
&\left(\mathbb{E}\left[\left(\int_{[a,b]^d}|\mathcal{N}_{u,v}^{\Theta_\mathbf{k},\mathbf{l}}(x) - \mathcal{E}(x)|^2\,\mathbb{P}_{X_1}(\mathrm{d}x)\right)^p\right]\right)^{1/p} \\
&\leq \left[\inf_{\theta\in[-c,c]^\mathbf{d}} \sup_{x\in[a,b]^d}|\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(x) - \mathcal{E}(x)|^2\right] \\
&\quad + \frac{4(v-u)\max\{1,|a|,|b|\}\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^\mathbf{L} c^\mathbf{L}\max\{1,p\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}} \\
&\quad + \frac{18\max\{1,(v-u)^2\}\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2\max\{p,\ln(3MB\max\{1,|a|,|b|\})\}}{\sqrt{M}} \qquad (7.21)\\
&\leq \left[\inf_{\theta\in[-c,c]^\mathbf{d}} \sup_{x\in[a,b]^d}|\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(x) - \mathcal{E}(x)|^2\right] + \frac{4(v-u)\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^\mathbf{L} c^{\mathbf{L}+1}\max\{1,p\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}} \\
&\quad + \frac{18\max\{1,(v-u)^2\}\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2\max\{p,\ln(3MBc)\}}{\sqrt{M}}.
\end{aligned}
$$

Furthermore, observe that Proposition 3.2.29 (applied with $f \curvearrowright \mathcal{E}$ in the notation of Proposition 3.2.29) proves that there exists $\vartheta \in \mathbb{R}^\mathbf{d}$ such that $\|\vartheta\|_\infty \leq \max\{1, L, |a|, |b|, 2[\sup_{x\in[a,b]^d}|\mathcal{E}(x)|]\}$ and

$$
\sup_{x\in[a,b]^d}|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)| \leq \frac{3dL(b-a)}{A^{1/d}}. \qquad (7.22)
$$

The fact that $\forall\, x \in [a,b]^d\colon \mathcal{E}(x) \in [u,v]$ hence implies that

$$
\|\vartheta\|_\infty \leq \max\{1, L, |a|, |b|, 2|u|, 2|v|\} \leq c. \qquad (7.23)
$$

This and (7.22) demonstrate that

$$
\begin{aligned}
&\inf_{\theta\in[-c,c]^\mathbf{d}} \sup_{x\in[a,b]^d}|\mathcal{N}_{u,v}^{\theta,\mathbf{l}}(x) - \mathcal{E}(x)|^2 \\
&\leq \sup_{x\in[a,b]^d}|\mathcal{N}_{u,v}^{\vartheta,\mathbf{l}}(x) - \mathcal{E}(x)|^2 \\
&\leq \left[\frac{3dL(b-a)}{A^{1/d}}\right]^2 = \frac{9d^2L^2(b-a)^2}{A^{2/d}}.
\end{aligned} \qquad (7.24)
$$

Combining this with (7.21) establishes that for all $p \in (0,\infty)$ it holds that

$$
\begin{aligned}
&\left(\mathbb{E}\left[\left(\int_{[a,b]^d}|\mathcal{N}_{u,v}^{\Theta_\mathbf{k},\mathbf{l}}(x) - \mathcal{E}(x)|^2\,\mathbb{P}_{X_1}(\mathrm{d}x)\right)^p\right]\right)^{1/p} \\
&\leq \frac{9d^2L^2(b-a)^2}{A^{2/d}} + \frac{4(v-u)\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^\mathbf{L} c^{\mathbf{L}+1}\max\{1,p\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}} \qquad (7.25)\\
&\quad + \frac{18\max\{1,(v-u)^2\}\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2\max\{p,\ln(3MBc)\}}{\sqrt{M}}.
\end{aligned}
$$

Moreover, note that the facts that $\max\{1, L, |a|, |b|\} \leq c$ and $(b-a)^2 \leq (|a| + |b|)^2 \leq 2(a^2 + b^2)$ yield that

$$
9L^2(b-a)^2 \leq 18c^2(a^2 + b^2) \leq 18c^2(c^2 + c^2) = 36c^4. \qquad (7.26)
$$

In addition, the fact that $B \geq c \geq 1$, the fact that $M \geq 1$, and item (ii) in Lemma 7.1.3 ensure that $\ln(3MBc) \leq \frac{23B}{18}\ln(eM)$. This, (7.26), the fact that $(v-u) \leq 2\max\{|u|,|v|\} = \max\{2|u|,2|v|\} \leq c \leq B$, and the fact that $B \geq 1$ prove that for all $p \in (0,\infty)$ it holds

that

$$\frac{9d^2L^2(b-a)^2}{A^{2/d}} + \frac{4(v-u)\mathbf{L}(\|\mathbf{l}\|_\infty+1)^{\mathbf{L}}c^{\mathbf{L}+1}\max\{1,p\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}}$$

$$+ \frac{18\max\{1,(v-u)^2\}\mathbf{L}(\|\mathbf{l}\|_\infty+1)^2\max\{p,\ln(3MBc)\}}{\sqrt{M}}$$

$$\leq \frac{36d^2c^4}{A^{2/d}} + \frac{4\mathbf{L}(\|\mathbf{l}\|_\infty+1)^{\mathbf{L}}c^{\mathbf{L}+2}\max\{1,p\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}} + \frac{23B^3\mathbf{L}(\|\mathbf{l}\|_\infty+1)^2\max\{p,\ln(eM)\}}{\sqrt{M}}.$$

$$(7.27)$$

Combining this with (7.25) shows (7.20). The proof of Theorem 7.1.4 is thus complete. $\square$

**Corollary 7.1.5.** *Let* $d, \mathbf{d}, \mathbf{L}, M, K, N \in \mathbb{N}$, $L, a, u \in \mathbb{R}$, $b \in (a,\infty)$, $v \in (u,\infty)$, $c \in [\max\{1, L, |a|, |b|, 2|u|, 2|v|\}, \infty)$, $B \in [c,\infty)$, $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \ldots, \mathbf{l_L}) \in \mathbb{N}^{\mathbf{L}+1}$, $\mathbf{T} \subseteq \{0,1,\ldots, N\}$, assume* $0 \in \mathbf{T}$, $\mathbf{l}_0 = d$, $\mathbf{l_L} = 1$, *and* $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}}\mathbf{l}_i(\mathbf{l}_{i-1}+1)$, *let* $(\Omega, \mathcal{F}, \mathbb{P})$ *be a probability space, let* $X_j\colon \Omega \to [a,b]^d$, $j \in \mathbb{N}$, *and* $Y_j\colon \Omega \to [u,v]$, $j \in \mathbb{N}$, *be functions, assume that* $(X_j, Y_j)$, $j \in \{1,2,\ldots,M\}$, *are i.i.d. random variables, let* $\mathcal{E}\colon [a,b]^d \to [u,v]$ *satisfy* $\mathbb{P}$-*a.s. that* $\mathcal{E}(X_1) = \mathbb{E}[Y_1|X_1]$, *assume for all* $x, y \in [a,b]^d$ *that* $|\mathcal{E}(x) - \mathcal{E}(y)| \leq L\|x-y\|_1$, *let* $\Theta_{k,n}\colon \Omega \to \mathbb{R}^{\mathbf{d}}$, $k, n \in \mathbb{N}_0$, *and* $\mathbf{k}\colon \Omega \to (\mathbb{N}_0)^2$ *be random variables, assume* $\left(\bigcup_{k=1}^{\infty}\Theta_{k,0}(\Omega)\right) \subseteq [-B,B]^{\mathbf{d}}$, *assume that* $\Theta_{k,0}$, $k \in \{1,2,\ldots,K\}$, *are i.i.d., assume that* $\Theta_{1,0}$ *is continuous uniformly distributed on* $[-c,c]^{\mathbf{d}}$, *and let* $\mathscr{R}\colon \mathbb{R}^{\mathbf{d}} \times \Omega \to [0,\infty)$ *satisfy for all* $\theta \in \mathbb{R}^{\mathbf{d}}$, $\omega \in \Omega$ *that*

$$\mathscr{R}(\theta,\omega) = \frac{1}{M}\left[\sum_{j=1}^{M}|\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j(\omega)) - Y_j(\omega)|^2\right] \qquad and \qquad (7.28)$$

$$\mathbf{k}(\omega) \in \operatorname{argmin}_{(k,n)\in\{1,2,\ldots,K\}\times\mathbf{T}, \|\Theta_{k,n}(\omega)\|_\infty \leq B}\mathscr{R}(\Theta_{k,n}(\omega),\omega) \qquad (7.29)$$

*(cf. Definitions 2.1.27 and 3.1.16). Then it holds for all* $p \in (0,\infty)$ *that*

$$\left(\mathbb{E}\left[\left(\int_{[a,b]^d}|\mathscr{N}_{u,v}^{\Theta_{\mathbf{k}},\mathbf{l}}(x) - \mathcal{E}(x)|^2\,\mathbb{P}_{X_1}(dx)\right)^{p/2}\right]\right)^{1/p}$$

$$\leq \frac{3dL(b-a)}{[\min(\{\mathbf{L}\}\cup\{\mathbf{l}_i\colon i \in \mathbb{N}\cap[0,\mathbf{L})\})]^{1/d}} + \frac{2[(v-u)\mathbf{L}(\|\mathbf{l}\|_\infty+1)^{\mathbf{L}}c^{\mathbf{L}+1}\max\{1,p/2\}]^{1/2}}{K^{[(2\mathbf{L})^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}}$$

$$+ \frac{3\max\{1,v-u\}(\|\mathbf{l}\|_\infty+1)[\mathbf{L}\max\{p,2\ln(3MBc)\}]^{1/2}}{M^{1/4}} \qquad (7.30)$$

$$\leq \frac{6dc^2}{[\min(\{\mathbf{L}\}\cup\{\mathbf{l}_i\colon i \in \mathbb{N}\cap[0,\mathbf{L})\})]^{1/d}} + \frac{2\mathbf{L}(\|\mathbf{l}\|_\infty+1)^{\mathbf{L}}c^{\mathbf{L}+1}\max\{1,p\}}{K^{[(2\mathbf{L})^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}}$$

$$+ \frac{5B^2\mathbf{L}(\|\mathbf{l}\|_\infty+1)\max\{p,\ln(eM)\}}{M^{1/4}}$$

*(cf. item (iii) in Lemma 7.1.1).*

*Proof of Corollary 7.1.5.* Throughout this proof let $A \in (0,\infty)$ be given by

$$A = \min(\{\mathbf{L}\}\cup\{\mathbf{l}_i\colon i \in \mathbb{N}\cap[0,\mathbf{L})\}). \qquad (7.31)$$

Note that (7.31) ensures that

$$\mathbf{L} \geq A = A - 1 + 1 \geq (A-1)\mathbb{1}_{[2,\infty)}(A) + 1$$

$$\geq \left(A - \tfrac{A}{2}\right)\mathbb{1}_{[2,\infty)}(A) + 1 = \frac{A\mathbb{1}_{[2,\infty)}(A)}{2} + 1 \geq \frac{A\mathbb{1}_{[6^d,\infty)}(A)}{2d} + 1. \qquad (7.32)$$

Moreover, note that the assumption that $\mathbf{l_L} = 1$ and (7.31) imply that

$$\mathbf{l}_1 = \mathbf{l}_1 \mathbb{1}_{\{1\}}(\mathbf{L}) + \mathbf{l}_1 \mathbb{1}_{[2,\infty)}(\mathbf{L}) \geq \mathbb{1}_{\{1\}}(\mathbf{L}) + A\mathbb{1}_{[2,\infty)}(\mathbf{L}) = A \geq A\mathbb{1}_{(6^d,\infty)}(A). \qquad (7.33)$$

Furthermore, observe that (7.31) shows that for all $i \in \{2, 3, \ldots\} \cap [0, \mathbf{L})$ it holds that

$$\begin{aligned}
\mathbf{l}_i \geq A &\geq A\mathbb{1}_{[2,\infty)}(A) \geq \mathbb{1}_{[2,\infty)}(A) \max\{A - 1, 2\} = \mathbb{1}_{[2,\infty)}(A) \max\{A - 4 + 3, 2\} \\
&\geq \mathbb{1}_{[2,\infty)}(A) \max\{A - 2i + 3, 2\} \geq \mathbb{1}_{(6^d,\infty)}(A) \max\{A/d - 2i + 3, 2\}.
\end{aligned} \qquad (7.34)$$

Combining (7.32)–(7.34) and Theorem 7.1.4 (applied with $p \curvearrowleft p/2$ for $p \in (0, \infty)$ in the notation of Theorem 7.1.4) establishes that for all $p \in (0, \infty)$ it holds that

$$\begin{aligned}
&\left(\mathbb{E}\left[\left(\int_{[a,b]^d} |\mathcal{N}_{u,v}^{\Theta_\mathbf{k},\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x)\right)^{p/2}\right]\right)^{2/p} \\
&\leq \frac{9d^2 L^2(b-a)^2}{A^{2/d}} + \frac{4(v-u)\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^\mathbf{L} c^{\mathbf{L}+1} \max\{1, p/2\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}} \\
&\quad + \frac{18\max\{1, (v-u)^2\}\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2 \max\{p/2, \ln(3MBc)\}}{\sqrt{M}} \qquad (7.35) \\
&\leq \frac{36d^2 c^4}{A^{2/d}} + \frac{4\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^\mathbf{L} c^{\mathbf{L}+2} \max\{1, p/2\}}{K^{[\mathbf{L}^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}} + \frac{23B^3\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2 \max\{p/2, \ln(eM)\}}{\sqrt{M}}.
\end{aligned}$$

This, (7.31), and the facts that $\mathbf{L} \geq 1$, $c \geq 1$, $B \geq 1$, and $\ln(eM) \geq 1$ demonstrate that for all $p \in (0, \infty)$ it holds that

$$\begin{aligned}
&\left(\mathbb{E}\left[\left(\int_{[a,b]^d} |\mathcal{N}_{u,v}^{\Theta_\mathbf{k},\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x)\right)^{p/2}\right]\right)^{1/p} \\
&\leq \frac{3dL(b-a)}{[\min(\{\mathbf{L}\} \cup \{\mathbf{l}_i \colon i \in \mathbb{N} \cap [0, \mathbf{L})\})]^{1/d}} + \frac{2[(v-u)\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^\mathbf{L} c^{\mathbf{L}+1} \max\{1, p/2\}]^{1/2}}{K^{[(2\mathbf{L})^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}} \\
&\quad + \frac{3\max\{1, v-u\}(\|\mathbf{l}\|_\infty + 1)[\mathbf{L}\max\{p, 2\ln(3MBc)\}]^{1/2}}{M^{1/4}} \\
&\leq \frac{6dc^2}{[\min(\{\mathbf{L}\} \cup \{\mathbf{l}_i \colon i \in \mathbb{N} \cap [0, \mathbf{L})\})]^{1/d}} + \frac{2[\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^\mathbf{L} c^{\mathbf{L}+2} \max\{1, p/2\}]^{1/2}}{K^{[(2\mathbf{L})^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}} \qquad (7.36) \\
&\quad + \frac{5B^3[\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^2 \max\{p/2, \ln(eM)\}]^{1/2}}{M^{1/4}} \\
&\leq \frac{6dc^2}{[\min(\{\mathbf{L}\} \cup \{\mathbf{l}_i \colon i \in \mathbb{N} \cap [0, \mathbf{L})\})]^{1/d}} + \frac{2\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^\mathbf{L} c^{\mathbf{L}+1} \max\{1, p\}}{K^{[(2\mathbf{L})^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}} \\
&\quad + \frac{5B^2\mathbf{L}(\|\mathbf{l}\|_\infty + 1) \max\{p, \ln(eM)\}}{M^{1/4}}.
\end{aligned}$$

The proof of Corollary 7.1.5 is thus complete. $\qquad\square$

## 7.2 Full strong error analysis for the training of ANNs with optimisation via stochastic gradient descent with random initialisation

**Corollary 7.2.1.** *Let* $d, \mathbf{d}, \mathbf{L}, M, K, N \in \mathbb{N}$, $L, a, u \in \mathbb{R}$, $b \in (a, \infty)$, $v \in (u, \infty)$, $c \in [\max\{1, L, |a|, |b|, 2|u|, 2|v|\}, \infty)$, $B \in [c, \infty)$, $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \ldots, \mathbf{l_L}) \in \mathbb{N}^{\mathbf{L}+1}$, $\mathbf{T} \subseteq \{0, 1, \ldots,$

$N\}$, $(\mathbf{J}_n)_{n\in\mathbb{N}} \subseteq \mathbb{N}$, $(\gamma_n)_{n\in\mathbb{N}} \subseteq \mathbb{R}$, *assume* $0 \in \mathbf{T}$, $\mathbf{l}_0 = d$, $\mathbf{l_L} = 1$, *and* $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1}+1)$, *let* $(\Omega, \mathcal{F}, \mathbb{P})$ *be a probability space, let* $X_j^{k,n}\colon \Omega \to [a,b]^d$, $k,n,j \in \mathbb{N}_0$, *and* $Y_j^{k,n}\colon \Omega \to [u,v]$, $k,n,j \in \mathbb{N}_0$, *be functions, assume that* $(X_j^{0,0}, Y_j^{0,0})$, $j \in \{1,2,\dots,M\}$, *are i.i.d. random variables, let* $\mathcal{E}\colon [a,b]^d \to [u,v]$ *satisfy* $\mathbb{P}$-*a.s. that* $\mathcal{E}(X_1^{0,0}) = \mathbb{E}[Y_1^{0,0}|X_1^{0,0}]$, *assume for all* $x,y \in [a,b]^d$ *that* $|\mathcal{E}(x)-\mathcal{E}(y)| \leq L\|x-y\|_1$, *let* $\Theta_{k,n}\colon \Omega \to \mathbb{R}^{\mathbf{d}}$, $k,n \in \mathbb{N}_0$, *and* $\mathbf{k}\colon \Omega \to (\mathbb{N}_0)^2$ *be random variables, assume* $\left(\bigcup_{k=1}^{\infty} \Theta_{k,0}(\Omega)\right) \subseteq [-B,B]^{\mathbf{d}}$, *assume that* $\Theta_{k,0}$, $k \in \{1,2,\dots,K\}$, *are i.i.d., assume that* $\Theta_{1,0}$ *is continuous uniformly distributed on* $[-c,c]^{\mathbf{d}}$, *let* $\mathscr{R}_J^{k,n}\colon \mathbb{R}^{\mathbf{d}} \times \Omega \to [0,\infty)$, $k,n,J \in \mathbb{N}_0$, *and* $\mathcal{G}^{k,n}\colon \mathbb{R}^{\mathbf{d}} \times \Omega \to \mathbb{R}^{\mathbf{d}}$, $k,n \in \mathbb{N}$, *satisfy for all* $k,n \in \mathbb{N}$, $\omega \in \Omega$, $\theta \in \{\vartheta \in \mathbb{R}^{\mathbf{d}}\colon (\mathscr{R}_{\mathbf{J}_n}^{k,n}(\cdot,\omega)\colon \mathbb{R}^{\mathbf{d}} \to [0,\infty)$ *is differentiable at* $\vartheta)\}$ *that* $\mathcal{G}^{k,n}(\theta,\omega) = (\nabla_\theta \mathscr{R}_{\mathbf{J}_n}^{k,n})(\theta,\omega)$, *assume for all* $k,n \in \mathbb{N}$ *that* $\Theta_{k,n} = \Theta_{k,n-1} - \gamma_n \mathcal{G}^{k,n}(\Theta_{k,n-1})$, *and assume for all* $k,n \in \mathbb{N}_0$, $J \in \mathbb{N}$, $\theta \in \mathbb{R}^{\mathbf{d}}$, $\omega \in \Omega$ *that*

$$\mathscr{R}_J^{k,n}(\theta,\omega) = \frac{1}{J}\left[\sum_{j=1}^J |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j^{k,n}(\omega)) - Y_j^{k,n}(\omega)|^2\right] \qquad and \qquad (7.37)$$

$$\mathbf{k}(\omega) \in \mathrm{argmin}_{(l,m)\in\{1,2,\dots,K\}\times\mathbf{T},\, \|\Theta_{l,m}(\omega)\|_\infty \leq B}\, \mathscr{R}_M^{0,0}(\Theta_{l,m}(\omega),\omega) \qquad (7.38)$$

*(cf. Definitions 2.1.27 and 3.1.16). Then it holds for all* $p \in (0,\infty)$ *that*

$$\left(\mathbb{E}\left[\left(\int_{[a,b]^d} |\mathscr{N}_{u,v}^{\Theta_{\mathbf{k}},\mathbf{l}}(x) - \mathcal{E}(x)|^2\, \mathbb{P}_{X_1^{0,0}}(\mathrm{d}x)\right)^{p/2}\right]\right)^{1/p}$$

$$\leq \frac{3dL(b-a)}{[\min(\{\mathbf{L}\}\cup\{\mathbf{l}_i\colon i\in\mathbb{N}\cap[0,\mathbf{L})\})]^{1/d}} + \frac{2[(v-u)\mathbf{L}(\|\mathbf{l}\|_\infty+1)^{\mathbf{L}}c^{\mathbf{L}+1}\max\{1,p/2\}]^{1/2}}{K^{[(2\mathbf{L})^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}}$$

$$+ \frac{3\max\{1,v-u\}(\|\mathbf{l}\|_\infty+1)[\mathbf{L}\max\{p,2\ln(3MBc)\}]^{1/2}}{M^{1/4}} \qquad (7.39)$$

$$\leq \frac{6dc^2}{[\min(\{\mathbf{L}\}\cup\{\mathbf{l}_i\colon i\in\mathbb{N}\cap[0,\mathbf{L})\})]^{1/d}} + \frac{2\mathbf{L}(\|\mathbf{l}\|_\infty+1)^{\mathbf{L}}c^{\mathbf{L}+1}\max\{1,p\}}{K^{[(2\mathbf{L})^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}}$$

$$+ \frac{5B^2\mathbf{L}(\|\mathbf{l}\|_\infty+1)\max\{p,\ln(eM)\}}{M^{1/4}}$$

*(cf. item (iii) in Lemma 7.1.1).*

*Proof of Corollary 7.2.1.* Observe that Corollary 7.1.5 (applied with $(X_j)_{j\in\mathbb{N}} \curvearrowleft (X_j^{0,0})_{j\in\mathbb{N}}$, $(Y_j)_{j\in\mathbb{N}} \curvearrowleft (Y_j^{0,0})_{j\in\mathbb{N}}$, $\mathscr{R} \curvearrowleft \mathscr{R}_M^{0,0}$ in the notation of Corollary 7.1.5) shows (7.39). The proof of Corollary 7.2.1 is thus complete. $\square$

**Corollary 7.2.2.** *Let* $d, \mathbf{d}, \mathbf{L}, M, K, N \in \mathbb{N}$, $L, a, u \in \mathbb{R}$, $b \in (a,\infty)$, $v \in (u,\infty)$, $c \in [\max\{1,L,|a|,|b|,2|u|,2|v|\},\infty)$, $B \in [c,\infty)$, $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \dots, \mathbf{l_L}) \in \mathbb{N}^{\mathbf{L}+1}$, $\mathbf{T} \subseteq \{0,1,\dots, N\}$, $(\mathbf{J}_n)_{n\in\mathbb{N}} \subseteq \mathbb{N}$, $(\gamma_n)_{n\in\mathbb{N}} \subseteq \mathbb{R}$, assume $0 \in \mathbf{T}$, $\mathbf{l}_0 = d$, $\mathbf{l_L} = 1$, and $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1}+1)$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $X_j^{k,n}\colon \Omega \to [a,b]^d$, $k,n,j \in \mathbb{N}_0$, and $Y_j^{k,n}\colon \Omega \to [u,v]$, $k,n,j \in \mathbb{N}_0$, be functions, assume that $(X_j^{0,0}, Y_j^{0,0})$, $j \in \{1,2,\dots,M\}$, are i.i.d. random variables, let $\mathcal{E}\colon [a,b]^d \to [u,v]$ satisfy $\mathbb{P}$-a.s. that $\mathcal{E}(X_1^{0,0}) = \mathbb{E}[Y_1^{0,0}|X_1^{0,0}]$, assume for all $x,y \in [a,b]^d$ that $|\mathcal{E}(x)-\mathcal{E}(y)| \leq L\|x-y\|_1$, let $\Theta_{k,n}\colon \Omega \to \mathbb{R}^{\mathbf{d}}$, $k,n \in \mathbb{N}_0$, and $\mathbf{k}\colon \Omega \to (\mathbb{N}_0)^2$ be random variables, assume $\left(\bigcup_{k=1}^{\infty} \Theta_{k,0}(\Omega)\right) \subseteq [-B,B]^{\mathbf{d}}$, assume that $\Theta_{k,0}$, $k \in \{1,2,\dots,K\}$, are i.i.d., assume that $\Theta_{1,0}$ is continuous uniformly distributed on $[-c,c]^{\mathbf{d}}$, let $\mathscr{R}_J^{k,n}\colon \mathbb{R}^{\mathbf{d}} \times \Omega \to [0,\infty)$, $k,n,J \in \mathbb{N}_0$, and $\mathcal{G}^{k,n}\colon \mathbb{R}^{\mathbf{d}} \times \Omega \to \mathbb{R}^{\mathbf{d}}$, $k,n \in \mathbb{N}$, satisfy for all $k,n \in \mathbb{N}$, $\omega \in \Omega$, $\theta \in \{\vartheta \in \mathbb{R}^{\mathbf{d}}\colon (\mathscr{R}_{\mathbf{J}_n}^{k,n}(\cdot,\omega)\colon \mathbb{R}^{\mathbf{d}} \to [0,\infty)$ is differentiable at $\vartheta)\}$ that $\mathcal{G}^{k,n}(\theta,\omega) = (\nabla_\theta \mathscr{R}_{\mathbf{J}_n}^{k,n})(\theta,\omega)$, assume for all $k,n \in \mathbb{N}$ that $\Theta_{k,n} = \Theta_{k,n-1} - \gamma_n \mathcal{G}^{k,n}(\Theta_{k,n-1})$,

and assume for all $k, n \in \mathbb{N}_0$, $J \in \mathbb{N}$, $\theta \in \mathbb{R}^{\mathbf{d}}$, $\omega \in \Omega$ that

$$\mathscr{R}_J^{k,n}(\theta,\omega) = \frac{1}{J}\left[\sum_{j=1}^{J} |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j^{k,n}(\omega)) - Y_j^{k,n}(\omega)|^2\right] \qquad and \tag{7.40}$$

$$\mathbf{k}(\omega) \in \operatorname{argmin}_{(l,m)\in\{1,2,\ldots,K\}\times\mathbf{T}, \|\Theta_{l,m}(\omega)\|_\infty \leq B} \mathscr{R}_M^{0,0}(\Theta_{l,m}(\omega),\omega) \tag{7.41}$$

*(cf. Definitions 2.1.27 and 3.1.16). Then*

$$\mathbb{E}\left[\int_{[a,b]^d} |\mathscr{N}_{u,v}^{\Theta_{\mathbf{k}},\mathbf{l}}(x) - \mathcal{E}(x)| \, \mathbb{P}_{X_1^{0,0}}(dx)\right] \leq \frac{2[(v-u)\mathbf{L}(\|\mathbf{l}\|_\infty+1)^{\mathbf{L}}c^{\mathbf{L}+1}]^{1/2}}{K^{[(2\mathbf{L})^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}}$$

$$+ \frac{3dL(b-a)}{[\min\{\mathbf{L},\mathbf{l}_1,\mathbf{l}_2,\ldots,\mathbf{l}_{\mathbf{L}-1}\}]^{1/d}} + \frac{3\max\{1,v-u\}(\|\mathbf{l}\|_\infty+1)[2\mathbf{L}\ln(3MBc)]^{1/2}}{M^{1/4}} \tag{7.42}$$

$$\leq \frac{6dc^2}{[\min\{\mathbf{L},\mathbf{l}_1,\mathbf{l}_2,\ldots,\mathbf{l}_{\mathbf{L}-1}\}]^{1/d}} + \frac{5B^2\mathbf{L}(\|\mathbf{l}\|_\infty+1)\ln(eM)}{M^{1/4}} + \frac{2\mathbf{L}(\|\mathbf{l}\|_\infty+1)^{\mathbf{L}}c^{\mathbf{L}+1}}{K^{[(2\mathbf{L})^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}}$$

*(cf. item (iii) in Lemma 7.1.1).*

*Proof of Corollary 7.2.2.* Note that Jensen's inequality implies that

$$\mathbb{E}\left[\int_{[a,b]^d} |\mathscr{N}_{u,v}^{\Theta_{\mathbf{k}},\mathbf{l}}(x) - \mathcal{E}(x)| \, \mathbb{P}_{X_1^{0,0}}(dx)\right] \leq \mathbb{E}\left[\left(\int_{[a,b]^d} |\mathscr{N}_{u,v}^{\Theta_{\mathbf{k}},\mathbf{l}}(x) - \mathcal{E}(x)|^2 \, \mathbb{P}_{X_1^{0,0}}(dx)\right)^{1/2}\right]. \tag{7.43}$$

This and Corollary 7.2.1 (applied with $p \curvearrowleft 1$ in the notation of Corollary 7.2.1) complete the proof of Corollary 7.2.2. $\qquad\square$

**Corollary 7.2.3.** *Let $d, \mathbf{d}, \mathbf{L}, M, K, N \in \mathbb{N}$, $L \in \mathbb{R}$, $c \in [\max\{2, L\}, \infty)$, $B \in [c, \infty)$, $\mathbf{l} = (\mathbf{l}_0, \mathbf{l}_1, \ldots, \mathbf{l}_{\mathbf{L}}) \in \mathbb{N}^{\mathbf{L}+1}$, $\mathbf{T} \subseteq \{0, 1, \ldots, N\}$, $(\mathbf{J}_n)_{n\in\mathbb{N}} \subseteq \mathbb{N}$, $(\gamma_n)_{n\in\mathbb{N}} \subseteq \mathbb{R}$, assume $0 \in \mathbf{T}$, $\mathbf{l}_0 = d$, $\mathbf{l}_{\mathbf{L}} = 1$, and $\mathbf{d} \geq \sum_{i=1}^{\mathbf{L}} \mathbf{l}_i(\mathbf{l}_{i-1}+1)$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $X_j^{k,n}: \Omega \to [0,1]^d$, $k, n, j \in \mathbb{N}_0$, and $Y_j^{k,n}: \Omega \to [0,1]$, $k, n, j \in \mathbb{N}_0$, be functions, assume that $(X_j^{0,0}, Y_j^{0,0})$, $j \in \{1, 2, \ldots, M\}$, are i.i.d. random variables, let $\mathcal{E}: [0,1]^d \to [0,1]$ satisfy $\mathbb{P}$-a.s. that $\mathcal{E}(X_1^{0,0}) = \mathbb{E}[Y_1^{0,0}|X_1^{0,0}]$, assume for all $x, y \in [0,1]^d$ that $|\mathcal{E}(x) - \mathcal{E}(y)| \leq L\|x - y\|_1$, let $\Theta_{k,n}: \Omega \to \mathbb{R}^{\mathbf{d}}$, $k, n \in \mathbb{N}_0$, and $\mathbf{k}: \Omega \to (\mathbb{N}_0)^2$ be random variables, assume $\left(\bigcup_{k=1}^\infty \Theta_{k,0}(\Omega)\right) \subseteq [-B, B]^{\mathbf{d}}$, assume that $\Theta_{k,0}$, $k \in \{1, 2, \ldots, K\}$, are i.i.d., assume that $\Theta_{1,0}$ is continuous uniformly distributed on $[-c,c]^{\mathbf{d}}$, let $\mathscr{R}_J^{k,n}: \mathbb{R}^{\mathbf{d}} \times \Omega \to [0,\infty)$, $k, n, J \in \mathbb{N}_0$, and $\mathcal{G}^{k,n}: \mathbb{R}^{\mathbf{d}} \times \Omega \to \mathbb{R}^{\mathbf{d}}$, $k, n \in \mathbb{N}$, satisfy for all $k, n \in \mathbb{N}$, $\omega \in \Omega$, $\theta \in \{\vartheta \in \mathbb{R}^{\mathbf{d}}: (\mathscr{R}_{\mathbf{J}_n}^{k,n}(\cdot,\omega): \mathbb{R}^{\mathbf{d}} \to [0,\infty)$ is differentiable at $\vartheta\}$ that $\mathcal{G}^{k,n}(\theta,\omega) = (\nabla_\theta \mathscr{R}_{\mathbf{J}_n}^{k,n})(\theta,\omega)$, assume for all $k, n \in \mathbb{N}$ that $\Theta_{k,n} = \Theta_{k,n-1} - \gamma_n \mathcal{G}^{k,n}(\Theta_{k,n-1})$, and assume for all $k, n \in \mathbb{N}_0$, $J \in \mathbb{N}$, $\theta \in \mathbb{R}^{\mathbf{d}}$, $\omega \in \Omega$ that*

$$\mathscr{R}_J^{k,n}(\theta,\omega) = \frac{1}{J}\left[\sum_{j=1}^{J} |\mathscr{N}_{u,v}^{\theta,\mathbf{l}}(X_j^{k,n}(\omega)) - Y_j^{k,n}(\omega)|^2\right] \qquad and \tag{7.44}$$

$$\mathbf{k}(\omega) \in \operatorname{argmin}_{(l,m)\in\{1,2,\ldots,K\}\times\mathbf{T}, \|\Theta_{l,m}(\omega)\|_\infty \leq B} \mathscr{R}_M^{0,0}(\Theta_{l,m}(\omega),\omega) \tag{7.45}$$

*(cf. Definitions 2.1.27 and 3.1.16). Then*

$$\mathbb{E}\left[\int_{[0,1]^d} |\mathscr{N}_{u,v}^{\Theta_{\mathbf{k}},\mathbf{l}}(x) - \mathcal{E}(x)| \, \mathbb{P}_{X_1^{0,0}}(dx)\right]$$

$$\leq \frac{3dL}{[\min\{\mathbf{L},\mathbf{l}_1,\mathbf{l}_2,\ldots,\mathbf{l}_{\mathbf{L}-1}\}]^{1/d}} + \frac{3(\|\mathbf{l}\|_\infty+1)[2\mathbf{L}\ln(3MBc)]^{1/2}}{M^{1/4}} + \frac{2[\mathbf{L}(\|\mathbf{l}\|_\infty+1)^{\mathbf{L}}c^{\mathbf{L}+1}]^{1/2}}{K^{[(2\mathbf{L})^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}}$$

$$\leq \frac{dc^3}{[\min\{\mathbf{L},\mathbf{l}_1,\mathbf{l}_2,\ldots,\mathbf{l}_{\mathbf{L}-1}\}]^{1/d}} + \frac{B^3\mathbf{L}(\|\mathbf{l}\|_\infty+1)\ln(eM)}{M^{1/4}} + \frac{\mathbf{L}(\|\mathbf{l}\|_\infty+1)^{\mathbf{L}}c^{\mathbf{L}+1}}{K^{[(2\mathbf{L})^{-1}(\|\mathbf{l}\|_\infty+1)^{-2}]}} \tag{7.46}$$

*(cf. item (iii) in Lemma 7.1.1).*

*Proof of Corollary 7.2.3.* Observe that Corollary 7.2.2 (applied with $a \curvearrowright 0$, $u \curvearrowright 0$, $b \curvearrowright 1$, $v \curvearrowright 1$ in the notation of Corollary 7.2.2), the facts that $B \geq c \geq \max\{2, L\}$ and $M \geq 1$, and item (ii) in Lemma 7.1.3 show that

$$
\mathbb{E}\Big[\int_{[0,1]^d} |\mathcal{N}_{u,v}^{\Theta_{\mathbf{k}},\mathbf{l}}(x) - \mathcal{E}(x)| \, \mathbb{P}_{X_1^{0,0}}(\mathrm{d}x)\Big]
$$
$$
\leq \frac{3dL}{[\min\{\mathbf{L}, \mathbf{l}_1, \mathbf{l}_2, \ldots, \mathbf{l}_{\mathbf{L}-1}\}]^{1/d}} + \frac{3(\|\mathbf{l}\|_\infty + 1)[2\mathbf{L}\ln(3MBc)]^{1/2}}{M^{1/4}} + \frac{2[\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}}c^{\mathbf{L}+1}]^{1/2}}{K^{[(2\mathbf{L})^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}}
$$
$$
\leq \frac{dc^3}{[\min\{\mathbf{L}, \mathbf{l}_1, \mathbf{l}_2, \ldots, \mathbf{l}_{\mathbf{L}-1}\}]^{1/d}} + \frac{(\|\mathbf{l}\|_\infty + 1)[23B\mathbf{L}\ln(eM)]^{1/2}}{M^{1/4}} + \frac{[\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}}c^{2\mathbf{L}+2}]^{1/2}}{K^{[(2\mathbf{L})^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}}
$$
$$
\leq \frac{dc^3}{[\min\{\mathbf{L}, \mathbf{l}_1, \mathbf{l}_2, \ldots, \mathbf{l}_{\mathbf{L}-1}\}]^{1/d}} + \frac{B^3\mathbf{L}(\|\mathbf{l}\|_\infty + 1)\ln(eM)}{M^{1/4}} + \frac{\mathbf{L}(\|\mathbf{l}\|_\infty + 1)^{\mathbf{L}}c^{\mathbf{L}+1}}{K^{[(2\mathbf{L})^{-1}(\|\mathbf{l}\|_\infty + 1)^{-2}]}}. \tag{7.47}
$$

The proof of Corollary 7.2.3 is thus complete. $\qquad\square$

# Chapter 8

# Stochastic gradient descent type optimization methods

This chapter reviews and studies stochastic gradient descent (SGD) type optimization methods such as the classical plain vanilla SGD optimization method (see Section 8.1) as well as more sophisticated SGD type optimization methods including SGD type optimization methods with momenta (cf. Sections 8.2, 8.3, and 8.7 below) and SGD type optimization methods with adaptive modifications of the learning rate (cf. Sections 8.4–8.7 below). We also refer to the overview article Ruder [26] and the reference list in [16] for further references on SGD type optimization methods.

## 8.1 The stochastic gradient descent optimization method

**Definition 8.1.1** (Stochastic gradient descent optimization method). *Let $d \in \mathbb{N}$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \infty)$, $(J_n)_{n \in \mathbb{N}} \subseteq \mathbb{N}$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $(S, \mathcal{S})$ be a measurable space, let $\xi \colon \Omega \to \mathbb{R}^d$ and $X_{n,j} \colon \Omega \to S$, $j \in \{1, 2, \dots, J_n\}$, $n \in \mathbb{N}$, be random variables, and let $F = (F(\theta, x))_{(\theta, x) \in \mathbb{R}^d \times S} \colon \mathbb{R}^d \times S \to \mathbb{R}$ and $G \colon \mathbb{R}^d \times S \to \mathbb{R}^d$ satisfy for all $x \in S$, $\theta \in \{v \in \mathbb{R}^d \colon F(\cdot, x) \text{ is differentiable at } v\}$ that*

$$G(\theta, x) = (\nabla_\theta F)(\theta, x). \tag{8.1}$$

*Then we say that $\Theta$ is the stochastic gradient descent process on $((\Omega, \mathcal{F}, \mathbb{P}), (S, \mathcal{S}))$ for the loss function $F$ with generalized gradient $G$, learning rates $(\gamma_n)_{n \in \mathbb{N}}$, batch sizes $(J_n)_{n \in \mathbb{N}}$, initial value $\xi$, and data $(X_{n,j})_{j \in \{1, 2, \dots, J_n\}, n \in \mathbb{N}}$ (we say that $\Theta$ is the stochastic gradient descent process for the loss function $F$ with learning rates $(\gamma_n)_{n \in \mathbb{N}}$, batch sizes $(J_n)_{n \in \mathbb{N}}$, initial value $\xi$, and data $(X_{n,j})_{j \in \{1, 2, \dots, J_n\}, n \in \mathbb{N}}$) if and only if it holds that $\Theta \colon \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ is the function from $\mathbb{N}_0 \times \Omega$ to $\mathbb{R}^d$ which satisfies for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \gamma_n \left[ \frac{1}{J_n} \sum_{j=1}^{J_n} G(\Theta_{n-1}, X_{n,j}) \right]. \tag{8.2}$$

### 8.1.1 Properties of the learning rates of the SGD optimization method

#### 8.1.1.1 Bias-variance decomposition of the mean square error

**Lemma 8.1.2** (Bias-variance decomposition of the mean square error). *Let $d \in \mathbb{N}$, $\vartheta \in \mathbb{R}^d$, let $\langle\!\langle \cdot, \cdot \rangle\!\rangle \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ be a scalar product, let $|\!|\!|\cdot|\!|\!| \colon \mathbb{R}^d \to [0, \infty)$ be the function*

which satisfies for all $v \in \mathbb{R}^d$ that $\||v\|| = \sqrt{\langle\!\langle v, v \rangle\!\rangle}$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, and let $Z \colon \Omega \to \mathbb{R}^d$ be a random variable with $\mathbb{E}[\||Z\||] < \infty$. Then

$$\mathbb{E}\big[\||Z - \vartheta\||^2\big] = \mathbb{E}\big[\||Z - \mathbb{E}[Z]\||^2\big] + \||\mathbb{E}[Z] - \vartheta\||^2. \tag{8.3}$$

*Proof of Lemma 8.1.2.* Observe that the assumption that $\mathbb{E}[\||Z\||] < \infty$ and the Cauchy-Schwarz inequality ensure that

$$
\begin{aligned}
\mathbb{E}\big[|\langle\!\langle Z - \mathbb{E}[Z], \mathbb{E}[Z] - \vartheta \rangle\!\rangle|\big] &\le \mathbb{E}\big[\||Z - \mathbb{E}[Z]\|| \, \||\mathbb{E}[Z] - \vartheta\||\big] \\
&\le (\mathbb{E}[\||Z\||] + \||\mathbb{E}[Z]\||) \, \||\mathbb{E}[Z] - \vartheta\|| < \infty.
\end{aligned}
\tag{8.4}
$$

The linearity of the expectation hence shows that

$$
\begin{aligned}
\mathbb{E}\big[\||Z - \vartheta\||^2\big] &= \mathbb{E}\big[\||(Z - \mathbb{E}[Z]) + (\mathbb{E}[Z] - \vartheta)\||^2\big] \\
&= \mathbb{E}\big[\||Z - \mathbb{E}[Z]\||^2 + 2\langle\!\langle Z - \mathbb{E}[Z], \mathbb{E}[Z] - \vartheta \rangle\!\rangle + \||\mathbb{E}[Z] - \vartheta\||^2\big] \\
&= \mathbb{E}\big[\||Z - \mathbb{E}[Z]\||^2\big] + 2\langle\!\langle \mathbb{E}[Z] - \mathbb{E}[Z], \mathbb{E}[Z] - \vartheta \rangle\!\rangle + \||\mathbb{E}[Z] - \vartheta\||^2 \\
&= \mathbb{E}\big[\||Z - \mathbb{E}[Z]\||^2\big] + \||\mathbb{E}[Z] - \vartheta\||^2.
\end{aligned}
\tag{8.5}
$$

The proof of Lemma 8.1.2 is thus complete. $\qquad\square$

### 8.1.1.2 On the stochasticity in the SGD optimization method

In this section we present Lemma 8.1.7, Corollary 8.1.8, and Example 8.1.9. Our proof of Lemma 8.1.7 employs the auxiliary results in Lemmas 8.1.3–8.1.6 below. Lemma 8.1.3 recalls an elementary and well known property for the expectation of the product of independent random variables (see, e.g., Klenke [19, Theorem 5.4]). In the elementary Lemma 8.1.6 we prove under suitable hypotheses the measurability of certain derivatives of a function. A result similar to Lemma 8.1.6 can, e.g., be found in [16, Lemma 4.4].

**Lemma 8.1.3.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and let $X, Y \colon \Omega \to \mathbb{R}$ be independent random variables with $\mathbb{E}[|X| + |Y|] < \infty$. Then*

(i) *it holds that $\mathbb{E}\big[|XY|\big] = \mathbb{E}\big[|X|\big]\mathbb{E}\big[|Y|\big] < \infty$ and*

(ii) *it holds that $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$.*

*Proof of Lemma 8.1.3.* Note that the fact that $(X, Y)(\mathbb{P}) = (X(\mathbb{P})) \otimes (Y(\mathbb{P}))$, the integral transformation theorem, Fubini's theorem, and the assumption that $\mathbb{E}[|X| + |Y|] < \infty$ assure that

$$
\begin{aligned}
\mathbb{E}\big[|XY|\big] &= \int_\Omega |X(\omega)Y(\omega)| \, \mathbb{P}(d\omega) \\
&= \int_{\mathbb{R} \times \mathbb{R}} |xy| \, \big((X, Y)(\mathbb{P})\big)(\mathrm{d}x, \mathrm{d}y) \\
&= \int_{\mathbb{R}} \left[\int_{\mathbb{R}} |xy| \, (X(\mathbb{P}))(\mathrm{d}x)\right] (Y(\mathbb{P}))(\mathrm{d}y) \\
&= \int_{\mathbb{R}} |y| \left[\int_{\mathbb{R}} |x| \, (X(\mathbb{P}))(\mathrm{d}x)\right] (Y(\mathbb{P}))(\mathrm{d}y) \\
&= \left[\int_{\mathbb{R}} |x| \, (X(\mathbb{P}))(\mathrm{d}x)\right]\left[\int_{\mathbb{R}} |y| \, (Y(\mathbb{P}))(\mathrm{d}y)\right] \\
&= \mathbb{E}\big[|X|\big]\mathbb{E}\big[|Y|\big] < \infty.
\end{aligned}
\tag{8.6}
$$

This proves item (i). In addition, observe that item (i), the fact that $(X, Y)(\mathbb{P}) = (X(\mathbb{P})) \otimes (Y(\mathbb{P}))$, the integral transformation theorem, and Fubini's theorem demonstrate that

$$
\begin{aligned}
\mathbb{E}[XY] &= \int_\Omega X(\omega) Y(\omega) \, \mathbb{P}(d\omega) \\
&= \int_{\mathbb{R} \times \mathbb{R}} xy \, ((X, Y)(\mathbb{P}))(\mathrm{d}x, \mathrm{d}y) \\
&= \int_{\mathbb{R}} \left[ \int_{\mathbb{R}} xy \, (X(\mathbb{P}))(\mathrm{d}x) \right] (Y(\mathbb{P}))(\mathrm{d}y) \\
&= \int_{\mathbb{R}} y \left[ \int_{\mathbb{R}} x \, (X(\mathbb{P}))(\mathrm{d}x) \right] (Y(\mathbb{P}))(\mathrm{d}y) \\
&= \left[ \int_{\mathbb{R}} x \, (X(\mathbb{P}))(\mathrm{d}x) \right] \left[ \int_{\mathbb{R}} y \, (Y(\mathbb{P}))(\mathrm{d}y) \right] \\
&= \mathbb{E}[X]\mathbb{E}[Y].
\end{aligned}
\tag{8.7}
$$

This establishes item (ii). The proof of Lemma 8.1.3 is thus complete. $\qquad\square$

**Lemma 8.1.4.** *Let $d, n \in \mathbb{N}$, let $\langle\!\langle \cdot, \cdot \rangle\!\rangle \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ be a scalar product, let $|\!|\!|\cdot|\!|\!| \colon \mathbb{R}^d \to [0, \infty)$ be the function which satisfies for all $v \in \mathbb{R}^d$ that $|\!|\!|v|\!|\!| = \sqrt{\langle\!\langle v, v \rangle\!\rangle}$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, and let $X_1, X_2, \ldots, X_n \colon \Omega \to \mathbb{R}^d$ be independent random variables which satisfy $\mathbb{E}[|\!|\!|X_1|\!|\!| + |\!|\!|X_2|\!|\!| + \ldots + |\!|\!|X_n|\!|\!|] < \infty$. Then it holds that*

$$
\mathbb{E}\left[ |\!|\!|\textstyle\sum_{k=1}^n (X_k - \mathbb{E}[X_k])|\!|\!|^2 \right] = \sum_{k=1}^n \mathbb{E}\left[ |\!|\!|X_k - \mathbb{E}[X_k]|\!|\!|^2 \right].
\tag{8.8}
$$

*Proof of Lemma 8.1.4.* First, note that Lemma 8.1.3 and the assumption that $\mathbb{E}[|\!|\!|X_1|\!|\!| + |\!|\!|X_2|\!|\!| + \ldots + |\!|\!|X_n|\!|\!|] < \infty$ ensure that for all $k_1, k_2 \in \{1, 2, \ldots, n\}$ with $k_1 \neq k_2$ it holds that

$$
\mathbb{E}\left[ |\langle\!\langle X_{k_1} - \mathbb{E}[X_{k_1}], X_{k_2} - \mathbb{E}[X_{k_2}] \rangle\!\rangle| \right] \leq \mathbb{E}\left[ |\!|\!|X_{k_1} - \mathbb{E}[X_{k_1}]|\!|\!| \, |\!|\!|X_{k_2} - \mathbb{E}[X_{k_2}]|\!|\!| \right] < \infty
\tag{8.9}
$$

and

$$
\begin{aligned}
&\mathbb{E}\left[ \langle\!\langle X_{k_1} - \mathbb{E}[X_{k_1}], X_{k_2} - \mathbb{E}[X_{k_2}] \rangle\!\rangle \right] \\
&= \langle \mathbb{E}[X_{k_1} - \mathbb{E}[X_{k_1}]], \mathbb{E}[X_{k_2} - \mathbb{E}[X_{k_2}]] \rangle \\
&= \langle \mathbb{E}[X_{k_1}] - \mathbb{E}[X_{k_1}], \mathbb{E}[X_{k_2}] - \mathbb{E}[X_{k_2}] \rangle = 0.
\end{aligned}
\tag{8.10}
$$

Hence, we obtain that

$$
\begin{aligned}
&\mathbb{E}\Big[\big\|\textstyle\sum_{k=1}^{n}(X_k-\mathbb{E}[X_k])\big\|^2\Big]\\
&=\mathbb{E}\big[\big\langle\textstyle\sum_{k_1=1}^{n}(X_{k_1}-\mathbb{E}[X_{k_1}]),\sum_{k_2=1}^{n}(X_{k_2}-\mathbb{E}[X_{k_2}])\big\rangle\big]\\
&=\mathbb{E}\Big[\textstyle\sum_{k_1,k_2=1}^{n}\langle X_{k_1}-\mathbb{E}[X_{k_1}],X_{k_2}-\mathbb{E}[X_{k_2}]\rangle\Big]\\
&=\mathbb{E}\Bigg[\bigg(\sum_{k=1}^{n}\|X_k-\mathbb{E}[X_k]\|^2\bigg)+\Bigg(\sum_{\substack{k_1,k_2=1\\k_1\neq k_2}}^{n}\langle\!\langle X_{k_1}-\mathbb{E}[X_{k_1}],X_{k_2}-\mathbb{E}[X_{k_2}]\rangle\!\rangle\Bigg)\Bigg]\\
&=\bigg(\sum_{k=1}^{n}\mathbb{E}\big[\|X_k-\mathbb{E}[X_k]\|^2\big]\bigg)+\Bigg(\sum_{\substack{k_1,k_2=1\\k_1\neq k_2}}^{n}\mathbb{E}\big[\langle\!\langle X_{k_1}-\mathbb{E}[X_{k_1}],X_{k_2}-\mathbb{E}[X_{k_2}]\rangle\!\rangle\big]\Bigg)\\
&=\sum_{k=1}^{n}\mathbb{E}\big[\|X_k-\mathbb{E}[X_k]\|^2\big].
\end{aligned}
\tag{8.11}
$$

The proof of Lemma 8.1.4 is thus complete. $\qquad\square$

**Lemma 8.1.5** (Factorization lemma for independent random variables). *Let $(\Omega,\mathcal{F},\mathbb{P})$ be a probability space, let $(\mathbb{X},\mathcal{X})$ and $(\mathbb{Y},\mathcal{Y})$ be measurable spaces, let $X\colon\Omega\to\mathbb{X}$ be $\mathcal{F}/\mathcal{X}$-measurable, let $Y\colon\Omega\to\mathbb{Y}$ be $\mathcal{F}/\mathcal{Y}$-measurable, assume that $X$ and $Y$ are independent, let $\Phi\colon\mathbb{X}\times\mathbb{Y}\to[0,\infty]$ be $(\mathcal{X}\otimes\mathcal{Y})/\mathcal{B}([0,\infty])$-measurable, and let $\phi\colon\mathbb{Y}\to[0,\infty]$ be the function which satisfies for all $y\in\mathbb{Y}$ that $\phi(y)=\mathbb{E}\big[\Phi(X,y)\big]$. Then*

(i) *it holds that the function $\phi$ is $\mathcal{Y}/\mathcal{B}([0,\infty])$-measurable and*

(ii) *it holds that*

$$
\mathbb{E}\big[\Phi(X,Y)\big]=\mathbb{E}\big[\phi(Y)\big].
\tag{8.12}
$$

*Proof of Lemma 8.1.5.* First, note that Fubini's theorem (cf., e.g., Klenke [19, (14.6) in Theorem 14.16]), the assumption that the function $X\colon\Omega\to\mathbb{X}$ is $\mathcal{F}/\mathcal{X}$-measurable, and the assumption that the function $\Phi\colon\mathbb{X}\times\mathbb{Y}\to[0,\infty]$ is $(\mathcal{X}\otimes\mathcal{Y})/\mathcal{B}([0,\infty])$-measurable demonstrate that the function

$$
\mathbb{Y}\ni y\mapsto\phi(y)=\mathbb{E}\big[\Phi(X,y)\big]=\int_{\Omega}\Phi(X(\omega),y)\,\mathbb{P}(\mathrm{d}\omega)\in[0,\infty]
\tag{8.13}
$$

is $\mathcal{Y}/\mathcal{B}([0,\infty])$-measurable. This proves item (i). Next observe that the integral transformation theorem, the fact that $(X,Y)(\mathbb{P})=(X(\mathbb{P}))\otimes(Y(\mathbb{P}))$, and Fubini's theorem prove that

$$
\begin{aligned}
\mathbb{E}\big[\Phi(X,Y)\big]&=\int_{\Omega}\Phi(X(\omega),Y(\omega))\,\mathbb{P}(\mathrm{d}\omega)\\
&=\int_{\mathbb{X}\times\mathbb{Y}}\Phi(x,y)\,\big((X,Y)(\mathbb{P})\big)(\mathrm{d}x,\mathrm{d}y)\\
&=\int_{\mathbb{Y}}\bigg[\int_{\mathbb{X}}\Phi(x,y)\,(X(\mathbb{P}))(\mathrm{d}x)\bigg]\,(Y(\mathbb{P}))(\mathrm{d}y)\\
&=\int_{\mathbb{Y}}\mathbb{E}\big[\Phi(X,y)\big]\,(Y(\mathbb{P}))(\mathrm{d}y)\\
&=\int_{\mathbb{Y}}\phi(y)\,(Y(\mathbb{P}))(\mathrm{d}y)=\mathbb{E}\big[\phi(Y)\big].
\end{aligned}
\tag{8.14}
$$

This establishes item (ii). The proof of Lemma 8.1.5 is thus complete. $\qquad\square$

**Lemma 8.1.6.** *Let $d \in \mathbb{N}$, let $(S, \mathcal{S})$ be a measurable space, let $F = (F(\theta, x))_{\theta \in \mathbb{R}^d, x \in S}$: $\mathbb{R}^d \times S \to \mathbb{R}$ be $(\mathcal{B}(\mathbb{R}^d) \otimes \mathcal{S})/\mathcal{B}(\mathbb{R})$-measurable, and assume for every $x \in S$ that the function $\mathbb{R}^d \ni \theta \mapsto F(\theta, x) \in \mathbb{R}$ is differentiable. Then it holds that the function*

$$\mathbb{R}^d \times S \ni (\theta, x) \mapsto (\nabla_\theta F)(\theta, x) \in \mathbb{R}^d \tag{8.15}$$

*is $(\mathcal{B}(\mathbb{R}^d) \otimes \mathcal{S})/\mathcal{B}(\mathbb{R}^d)$-measurable.*

*Proof of Lemma 8.1.6.* Throughout this proof let $G = (G_1, \ldots, G_d)\colon \mathbb{R}^d \times S \to \mathbb{R}^d$ be the function which satisfies for all $\theta \in \mathbb{R}^d$, $x \in S$ that

$$G(\theta, x) = (\nabla_\theta F)(\theta, x). \tag{8.16}$$

The assumption that the function $F\colon \mathbb{R}^d \times S \to \mathbb{R}$ is $(\mathcal{B}(\mathbb{R}^d) \otimes \mathcal{S})/\mathcal{B}(\mathbb{R})$-measurable implies that for all $i \in \{1, \ldots, d\}$, $h \in \mathbb{R} \backslash \{0\}$ it holds that the function

$$\mathbb{R}^d \times S \ni (\theta, x) = ((\theta_1, \ldots, \theta_d), x) \mapsto \left( \tfrac{F((\theta_1, \ldots, \theta_{i-1}, \theta_i + h, \theta_{i+1}, \ldots, \theta_d), x) - F(\theta, x)}{h} \right) \in \mathbb{R} \tag{8.17}$$

is $(\mathcal{B}(\mathbb{R}^d) \otimes \mathcal{S})/\mathcal{B}(\mathbb{R})$-measurable. The fact that for all $i \in \{1, \ldots, d\}$, $\theta = (\theta_1, \ldots, \theta_d) \in \mathbb{R}^d$, $x \in S$ it holds that

$$G_i(\theta, x) = \lim_{n \to \infty} \left( \tfrac{F((\theta_1, \ldots, \theta_{i-1}, \theta_i + 2^{-n}, \theta_{i+1}, \ldots, \theta_d), x) - F(\theta, x)}{2^{-n}} \right) \tag{8.18}$$

hence ensures that for all $i \in \{1, \ldots, d\}$ it holds that the function $G_i\colon \mathbb{R}^d \times S \to \mathbb{R}$ is $(\mathcal{B}(\mathbb{R}^d) \otimes \mathcal{S})/\mathcal{B}(\mathbb{R})$-measurable. This implies that $G$ is $(\mathcal{B}(\mathbb{R}^d) \otimes \mathcal{S})/\mathcal{B}(\mathbb{R}^d)$-measurable. The proof of Lemma 8.1.6 is thus complete. $\qquad\square$

**Lemma 8.1.7.** *Let $d \in \mathbb{N}$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \infty)$, $(J_n)_{n \in \mathbb{N}} \subseteq \mathbb{N}$, let $\langle \cdot, \cdot \rangle\colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ be a scalar product, let $\|\cdot\|\colon \mathbb{R}^d \to [0, \infty)$ be the function which satisfies for all $v \in \mathbb{R}^d$ that $\|v\| = \sqrt{\langle v, v \rangle}$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $\xi\colon \Omega \to \mathbb{R}^d$ be a random variable, let $(S, \mathcal{S})$ be a measurable space, let $X_{n,j}\colon \Omega \to S$, $j \in \{1, 2, \ldots, J_n\}$, $n \in \mathbb{N}$, be i.i.d. random variables, assume that $\xi$ and $(X_{n,j})_{j \in \{1,2,\ldots,J_n\}, n \in \mathbb{N}}$ are independent, let $F = (F(\theta, x))_{(\theta, x) \in \mathbb{R}^d \times S}\colon \mathbb{R}^d \times S \to \mathbb{R}$ be $(\mathcal{B}(\mathbb{R}^d) \otimes \mathcal{S})/\mathcal{B}(\mathbb{R})$-measurable, assume for all $x \in S$ that $(\mathbb{R}^d \ni \theta \mapsto F(\theta, x) \in \mathbb{R}) \in C^1(\mathbb{R}^d, \mathbb{R})$, assume for all $\theta \in \mathbb{R}^d$ that $\mathbb{E}[\|(\nabla_\theta F)(\theta, X_{1,1})\|] < \infty$ (cf. Lemma 8.1.6), let $\mathcal{V}\colon \mathbb{R}^d \to [0, \infty]$ be the function which satisfies for all $\theta \in \mathbb{R}^d$ that*

$$\mathcal{V}(\theta) = \mathbb{E}\Big[ \big\| (\nabla_\theta F)(\theta, X_{1,1}) - \mathbb{E}[(\nabla_\theta F)(\theta, X_{1,1})] \big\|^2 \Big], \tag{8.19}$$

*and let $\Theta\colon \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ be the stochastic process which satisfies for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad \text{and} \qquad \Theta_n = \Theta_{n-1} - \gamma_n \left[ \frac{1}{J_n} \sum_{j=1}^{J_n} (\nabla_\theta F)(\Theta_{n-1}, X_{n,j}) \right]. \tag{8.20}$$

*Then it holds for all $n \in \mathbb{N}$, $\vartheta \in \mathbb{R}^d$ that*

$$\big( \mathbb{E}[\|\Theta_n - \vartheta\|^2] \big)^{1/2} \geq \frac{\gamma_n}{(J_n)^{1/2}} \big( \mathbb{E}[\mathcal{V}(\Theta_{n-1})] \big)^{1/2}. \tag{8.21}$$

*Proof of Lemma 8.1.7.* Throughout this proof let $\phi_n \colon \mathbb{R}^d \to [0, \infty]$, $n \in \mathbb{N}$, be the functions which satisfy for all $n \in \mathbb{N}$, $\theta \in \mathbb{R}^d$ that

$$\phi_n(\theta) = \mathbb{E}\left[\left\|\theta - \tfrac{\gamma_n}{J_n}\left[\textstyle\sum_{j=1}^{J_n}(\nabla_\theta F)(\theta, X_{n,j})\right] - \vartheta\right\|^2\right]. \tag{8.22}$$

Observe that Lemma 8.1.2 ensures that for all $\vartheta \in \mathbb{R}^d$ and all random variables $Z \colon \Omega \to \mathbb{R}^d$ with $\mathbb{E}[\|Z\|] < \infty$ it holds that

$$\begin{aligned}
\mathbb{E}\big[\|Z - \vartheta\|^2\big] &= \mathbb{E}\big[\|Z - \mathbb{E}[Z]\|^2\big] + \|\mathbb{E}[Z] - \vartheta\|^2 \\
&\geq \mathbb{E}\big[\|Z - \mathbb{E}[Z]\|^2\big].
\end{aligned} \tag{8.23}$$

Hence, we obtain for all $n \in \mathbb{N}$, $\theta \in \mathbb{R}^d$ that

$$\begin{aligned}
\phi_n(\theta) &= \mathbb{E}\left[\left\|\tfrac{\gamma_n}{J_n}\left[\textstyle\sum_{j=1}^{J_n}(\nabla_\theta F)(\theta, X_{n,j})\right] - (\theta - \vartheta)\right\|^2\right] \\
&\geq \mathbb{E}\left[\left\|\tfrac{\gamma_n}{J_n}\left[\textstyle\sum_{j=1}^{J_n}(\nabla_\theta F)(\theta, X_{n,j})\right] - \mathbb{E}\left[\tfrac{\gamma_n}{J_n}\left[\textstyle\sum_{j=1}^{J_n}(\nabla_\theta F)(\theta, X_{n,j})\right]\right]\right\|^2\right] \\
&= \tfrac{(\gamma_n)^2}{(J_n)^2}\, \mathbb{E}\left[\left\|\textstyle\sum_{j=1}^{J_n}\big((\nabla_\theta F)(\theta, X_{n,j}) - \mathbb{E}[(\nabla_\theta F)(\theta, X_{n,j})]\big)\right\|^2\right].
\end{aligned} \tag{8.24}$$

Lemma 8.1.4, the fact that $X_{n,j} \colon \Omega \to S$, $j \in \{1, 2, \ldots, J_n\}, n \in \mathbb{N}$, are i.i.d. random variables, and the fact that for all $n \in \mathbb{N}, j \in \{1, 2, \ldots, J_n\}, \theta \in \mathbb{R}^d$ it holds that

$$\mathbb{E}\big[\|(\nabla_\theta F)(\theta, X_{n,j})\|\big] = \mathbb{E}\big[\|(\nabla_\theta F)(\theta, X_{1,1})\|\big] < \infty \tag{8.25}$$

hence demonstrates that for all $n \in \mathbb{N}$, $\theta \in \mathbb{R}^d$ it holds that

$$\begin{aligned}
\phi_n(\theta) &\geq \tfrac{(\gamma_n)^2}{(J_n)^2}\left[\sum_{j=1}^{J_n} \mathbb{E}\Big[\big\|(\nabla_\theta F)(\theta, X_{n,j}) - \mathbb{E}[(\nabla_\theta F)(\theta, X_{n,j})]\big\|^2\Big]\right] \\
&= \tfrac{(\gamma_n)^2}{(J_n)^2}\left[\sum_{j=1}^{J_n} \mathbb{E}\Big[\big\|(\nabla_\theta F)(\theta, X_{1,1}) - \mathbb{E}[(\nabla_\theta F)(\theta, X_{1,1})]\big\|^2\Big]\right] \\
&= \tfrac{(\gamma_n)^2}{(J_n)^2}\left[\sum_{j=1}^{J_n} \mathcal{V}(\theta)\right] = \tfrac{(\gamma_n)^2}{(J_n)^2}\big[J_n \mathcal{V}(\theta)\big] = \left(\tfrac{(\gamma_n)^2}{J_n}\right)\mathcal{V}(\theta).
\end{aligned} \tag{8.26}$$

In addition, observe that (8.20), (8.22), the fact that for all $n \in \mathbb{N}$ it holds that $\Theta_{n-1}$ and $X_n$ are independent random variables, and Lemma 8.1.5 assure that for all $n \in \mathbb{N}$, $\vartheta \in \mathbb{R}^d$ it holds that

$$\begin{aligned}
\mathbb{E}\big[\|\Theta_n - \vartheta\|^2\big] &= \mathbb{E}\left[\left\|\Theta_{n-1} - \tfrac{\gamma_n}{J_n}\left[\textstyle\sum_{j=1}^{J_n}(\nabla_\theta F)(\Theta_{n-1}, X_{n,j})\right] - \vartheta\right\|^2\right] \\
&= \mathbb{E}\big[\phi_n(\Theta_{n-1})\big].
\end{aligned} \tag{8.27}$$

Combining this with (8.26) proves that for all $n \in \mathbb{N}$, $\vartheta \in \mathbb{R}^d$ it holds that

$$\mathbb{E}\big[\|\Theta_n - \vartheta\|^2\big] \geq \mathbb{E}\left[\left(\tfrac{(\gamma_n)^2}{J_n}\right)\mathcal{V}(\Theta_{n-1})\right] = \left(\tfrac{(\gamma_n)^2}{J_n}\right)\mathbb{E}\big[\mathcal{V}(\Theta_{n-1})\big]. \tag{8.28}$$

This establishes (8.21). The proof of Lemma 8.1.7 is thus complete. $\qquad\square$

**Corollary 8.1.8.** *Let $d \in \mathbb{N}$, $\varepsilon \in [0, \infty)$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \infty)$, $(J_n)_{n \in \mathbb{N}} \subseteq \mathbb{N}$, let $\langle \cdot, \cdot \rangle \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ be a scalar product, let $\|\cdot\| \colon \mathbb{R}^d \to [0, \infty)$ be the function which satisfies for all $v \in \mathbb{R}^d$ that $\|v\| = \sqrt{\langle v, v \rangle}$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $\xi \colon \Omega \to \mathbb{R}^d$ be a random variable, let $(S, \mathcal{S})$ be a measurable space, let $X_{n,j} \colon \Omega \to S$, $j \in \{1, 2, \ldots, J_n\}$, $n \in \mathbb{N}$, be i.i.d. random variables, assume that $\xi$ and $(X_{n,j})_{j \in \{1,2,\ldots,J_n\}, n \in \mathbb{N}}$ are independent, let $F = (F(\theta, x))_{(\theta,x) \in \mathbb{R}^d \times S} \colon \mathbb{R}^d \times S \to \mathbb{R}$ be $(\mathcal{B}(\mathbb{R}^d) \otimes \mathcal{S})/\mathcal{B}(\mathbb{R})$-measurable, assume for all $x \in S$ that $(\mathbb{R}^d \ni \theta \mapsto F(\theta, x) \in \mathbb{R}) \in C^1(\mathbb{R}^d, \mathbb{R})$, assume for all $\theta \in \mathbb{R}^d$ that $\mathbb{E}\big[\|(\nabla_\theta F)(\theta, X_{1,1})\|\big] < \infty$ (cf. Lemma 8.1.6) and*

$$\left(\mathbb{E}\Big[\big\|(\nabla_\theta F)(\theta, X_{1,1}) - \mathbb{E}\big[(\nabla_\theta F)(\theta, X_{1,1})\big]\big\|^2\Big]\right)^{1/2} \geq \varepsilon, \tag{8.29}$$

*and let $\Theta \colon \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ be the stochastic process which satisfies for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad \text{and} \qquad \Theta_n = \Theta_{n-1} - \gamma_n \left[\frac{1}{J_n} \sum_{j=1}^{J_n} (\nabla_\theta F)(\Theta_{n-1}, X_{n,j})\right]. \tag{8.30}$$

*Then*

*(i) it holds for all $n \in \mathbb{N}$, $\vartheta \in \mathbb{R}^d$ that*

$$\left(\mathbb{E}\big[\|\Theta_n - \vartheta\|^2\big]\right)^{1/2} \geq \varepsilon\left(\frac{\gamma_n}{(J_n)^{1/2}}\right) \tag{8.31}$$

*and*

*(ii) it holds for all $\vartheta \in \mathbb{R}^d$ that*

$$\liminf_{n \to \infty}\left(\mathbb{E}\big[\|\Theta_n - \vartheta\|^2\big]\right)^{1/2} \geq \varepsilon\left(\liminf_{n \to \infty}\left[\frac{\gamma_n}{(J_n)^{1/2}}\right]\right). \tag{8.32}$$

*Proof of Corollary 8.1.8.* Throughout this proof let $\mathcal{V} \colon \mathbb{R}^d \to [0, \infty]$ be the function which satisfies for all $\theta \in \mathbb{R}^d$ that

$$\mathcal{V}(\theta) = \mathbb{E}\Big[\big\|(\nabla_\theta F)(\theta, X_{1,1}) - \mathbb{E}\big[(\nabla_\theta F)(\theta, X_{1,1})\big]\big\|^2\Big]. \tag{8.33}$$

Note that (8.29) assures that for all $\theta \in \mathbb{R}^d$ it holds that

$$\mathcal{V}(\theta) \geq \varepsilon^2. \tag{8.34}$$

Lemma 8.1.7 therefore demonstrates that for all $n \in \mathbb{N}$, $\vartheta \in \mathbb{R}^d$ it holds that

$$\left(\mathbb{E}\big[\|\Theta_n - \vartheta\|^2\big]\right)^{1/2} \geq \frac{\gamma_n}{(J_n)^{1/2}}\left(\mathbb{E}\big[\mathcal{V}(\Theta_{n-1})\big]\right)^{1/2} \geq \left[\frac{\gamma_n}{(J_n)^{1/2}}\right](\varepsilon^2)^{1/2} = \frac{\gamma_n \varepsilon}{(J_n)^{1/2}}. \tag{8.35}$$

This proves item (i). Furthermore, note that item (i) implies item (ii). The proof of Corollary 8.1.8 is thus complete. $\qquad \square$

**Example 8.1.9** (A lower bound for the SGD optimization method)**.** *Let $d \in \mathbb{N}$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \infty)$, $(J_n)_{n \in \mathbb{N}} \subseteq \mathbb{N}$, let $\|\cdot\| \colon \mathbb{R}^d \to [0, \infty)$ be the d-dimensional Euclidean norm, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $\xi \colon \Omega \to \mathbb{R}^d$ be a random variable, let $X_{n,j} \colon \Omega \to \mathbb{R}^d$, $j \in \{1, 2, \ldots, J_n\}$, $n \in \mathbb{N}$, be i.i.d. random variables with $\mathbb{E}[\|X_{1,1}\|] < \infty$, assume that $\xi$*

and $(X_{n,j})_{j \in \{1,2,\ldots,J_n\}, n \in \mathbb{N}}$ *are independent, let* $F = (F(\theta, x))_{(\theta, x) \in \mathbb{R}^d \times \mathbb{R}^d} \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ *be the function which satisfies for all* $\theta, x \in \mathbb{R}^d$ *that*

$$F(\theta, x) = \tfrac{1}{2}\|\theta - x\|^2, \tag{8.36}$$

*and let* $\Theta \colon \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ *be the stochastic process which satisfies for all* $n \in \mathbb{N}$ *that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \gamma_n \left[ \frac{1}{J_n} \sum_{j=1}^{J_n} (\nabla_\theta F)(\Theta_{n-1}, X_{n,j}) \right]. \tag{8.37}$$

*Then*

(i) *it holds for all* $\theta \in \mathbb{R}^d$ *that*

$$\mathbb{E}\big[\|(\nabla_\theta F)(\theta, X_{1,1})\|\big] < \infty, \tag{8.38}$$

(ii) *it holds for all* $\theta \in \mathbb{R}^d$ *that*

$$\mathbb{E}\left[\big\|(\nabla_\theta F)(\theta, X_{1,1}) - \mathbb{E}\big[(\nabla_\theta F)(\theta, X_{1,1})\big]\big\|^2\right] = \mathbb{E}\big[\|X_{1,1} - \mathbb{E}[X_{1,1}]\|^2\big], \tag{8.39}$$

*and*

(iii) *it holds for all* $n \in \mathbb{N}$, $\vartheta \in \mathbb{R}^d$ *that*

$$\big(\mathbb{E}\big[\|\Theta_n - \vartheta\|^2\big]\big)^{1/2} \geq \big(\mathbb{E}\big[\|X_{1,1} - \mathbb{E}[X_{1,1}]\|^2\big]\big)^{1/2}\left[\frac{\gamma_n}{(J_n)^{1/2}}\right]. \tag{8.40}$$

*Proof of Example 8.1.9.* First, note that (8.36) and Lemma 13.2.4 imply that for all $\theta, x \in \mathbb{R}^d$ it holds that

$$(\nabla_\theta F)(\theta, x) = \tfrac{1}{2}(2(\theta - x)) = \theta - x. \tag{8.41}$$

The assumption that $\mathbb{E}[\|X_{1,1}\|] < \infty$ therefore assures that for all $\theta \in \mathbb{R}^d$ it holds that

$$\mathbb{E}\big[\|(\nabla_\theta F)(\theta, X_{1,1})\|\big] = \mathbb{E}\big[\|\theta - X_{1,1}\|\big] \leq \|\theta\| + \mathbb{E}\big[\|X_{1,1}\|\big] < \infty. \tag{8.42}$$

This establishes item (i). Moreover, observe that (8.41) and item (i) ensure that for all $\theta \in \mathbb{R}^d$ it holds that

$$\begin{aligned}
&\mathbb{E}\big[\|(\nabla_\theta F)(\theta, X_{1,1}) - \mathbb{E}[(\nabla_\theta F)(\theta, X_{1,1})]\|^2\big] \\
&= \mathbb{E}\big[\|(\theta - X_{1,1}) - \mathbb{E}[\theta - X_{1,1}]\|^2\big] = \mathbb{E}\big[\|X_{1,1} - \mathbb{E}[X_{1,1}]\|^2\big].
\end{aligned} \tag{8.43}$$

This proves item (ii). In addition, note that item (i) in Corollary 8.1.8 and items (i)–(ii) establish item (iii). The proof of Example 8.1.9 is thus complete. $\qquad\square$

### 8.1.1.3   A lower bound for the natural logarithm

In the next auxiliary result, Lemma 8.1.10 below, we recall a well known lower bound for the natural logarithm.

**Lemma 8.1.10** (A lower bound for the natural logarithm)**.** *It holds for all* $x \in (0, \infty)$ *that*

$$\ln(x) \geq \frac{(x-1)}{x}. \tag{8.44}$$

*Proof of Lemma 8.1.10.* First, note that the fundamental theorem of calculus ensures that for all $x \in [1, \infty)$ it holds that

$$\ln(x) = \ln(x) - \ln(1) = \int_1^x \frac{1}{t}\, dt \geq \int_1^x \frac{1}{x}\, dt = \frac{(x-1)}{x}. \tag{8.45}$$

Moreover, observe hat the fundamental theorem of calculus ensures that for all $x \in (0, 1]$ it holds that

$$\begin{aligned}
\ln(x) &= \ln(x) - \ln(1) = -(\ln(1) - \ln(x)) = -\left[\int_x^1 \frac{1}{t}\, dt\right] \\
&= \int_x^1 \left(-\frac{1}{t}\right) dt \geq \int_x^1 \left(-\frac{1}{x}\right) dt = (1-x)\left(-\frac{1}{x}\right) = \frac{(x-1)}{x}.
\end{aligned} \tag{8.46}$$

Combining this and (8.45) establishes (8.44). The proof of Lemma 8.1.10 is thus complete. $\qquad\square$

### 8.1.1.4 Summable learning rates

**Lemma 8.1.11** (Gradient descent fails to converge for a summable sequence of learning rates)**.** *Let $d \in \mathbb{N}$, $\vartheta \in \mathbb{R}^d$, $\xi \in \mathbb{R}^d/\{\vartheta\}$, $\alpha \in (0, \infty)$, $(\gamma_n)_{n\in\mathbb{N}} \subseteq (0, \infty)\backslash\{1/\alpha\}$ satisfy $\sum_{n=1}^\infty \gamma_n < \infty$, let $\|\cdot\|\colon \mathbb{R}^d \to [0, \infty)$ be the d-dimensional Euclidean norm, let $f\colon \mathbb{R}^d \to \mathbb{R}$ be the function which satisfies for all $\theta \in \mathbb{R}^d$ that*

$$f(\theta) = \tfrac{\alpha}{2}\|\theta - \vartheta\|^2, \tag{8.47}$$

*and let $\Theta\colon \mathbb{N}_0 \to \mathbb{R}^d$ be the function which satisfies for all $n \in \mathbb{N}$ that $\Theta_0 = \xi$ and*

$$\Theta_n = \Theta_{n-1} - \gamma_n(\nabla f)(\Theta_{n-1}). \tag{8.48}$$

*Then*

*(i) it holds for all $n \in \mathbb{N}_0$ that*

$$\Theta_n - \vartheta = \left[\prod_{k=1}^n (1 - \gamma_k\alpha)\right](\xi - \vartheta), \tag{8.49}$$

*(ii) it holds that*

$$\liminf_{n\to\infty}\left[\prod_{k=1}^n |1 - \gamma_k\alpha|\right] > 0, \tag{8.50}$$

*and*

*(iii) it holds that*

$$\liminf_{n\to\infty}\|\Theta_n - \vartheta\| > 0. \tag{8.51}$$

*Proof of Lemma 8.1.11.* Throughout this proof let $m \in \mathbb{N}$ satisfy for all $k \in \mathbb{N} \cap [m, \infty)$ that $\gamma_k < 1/(2\alpha)$. Observe that Lemma 13.2.4 implies that for all $\theta \in \mathbb{R}^d$ it holds that

$$(\nabla f)(\theta) = \tfrac{\alpha}{2}(2(\theta - \vartheta)) = \alpha(\theta - \vartheta). \tag{8.52}$$

Therefore, we obtain for all $n \in \mathbb{N}$ that

$$
\begin{aligned}
\Theta_n - \vartheta &= \Theta_{n-1} - \gamma_n (\nabla f)(\Theta_{n-1}) - \vartheta \\
&= \Theta_{n-1} - \gamma_n \alpha (\Theta_{n-1} - \vartheta) - \vartheta \\
&= (1 - \gamma_n \alpha)(\Theta_{n-1} - \vartheta).
\end{aligned}
\tag{8.53}
$$

Induction hence proves that for all $n \in \mathbb{N}$ it holds that

$$
\Theta_n - \vartheta = \left[ \prod_{k=1}^{n} (1 - \gamma_k \alpha) \right] (\Theta_0 - \vartheta),
\tag{8.54}
$$

This and the assumption that $\Theta_0 = \xi$ establish item (i). Next observe that the fact that for all $k \in \mathbb{N}$ it holds that $\gamma_k \alpha \neq 1$ ensures that

$$
\prod_{k=1}^{m-1} |1 - \gamma_k \alpha| > 0.
\tag{8.55}
$$

Moreover, note that the fact that for all $k \in \mathbb{N} \cap [m, \infty)$ it holds that $\gamma_k \alpha \in (0, 1/2)$ assures that for all $k \in \mathbb{N} \cap [m, \infty)$ it holds that

$$
(1 - \gamma_k \alpha) \in (1/2, 1).
\tag{8.56}
$$

This, Lemma 8.1.10, and the assumption that $\sum_{n=1}^{\infty} \gamma_n < \infty$ demonstrate that for all $n \in \mathbb{N} \cap [m, \infty)$ it holds that

$$
\begin{aligned}
\ln \left( \prod_{k=m}^{n} |1 - \gamma_k \alpha| \right) &= \sum_{k=m}^{n} \ln(1 - \gamma_k \alpha) \\
&\geq \sum_{k=m}^{n} \frac{(1 - \gamma_k \alpha) - 1}{(1 - \gamma_k \alpha)} = \sum_{k=m}^{n} \left[ -\frac{\gamma_k \alpha}{(1 - \gamma_k \alpha)} \right] \\
&\geq \sum_{k=m}^{n} \left[ -\frac{\gamma_k \alpha}{\left(\frac{1}{2}\right)} \right] = -2\alpha \left[ \sum_{k=m}^{n} \gamma_k \right] \geq -2\alpha \left[ \sum_{k=1}^{\infty} \gamma_k \right] > -\infty.
\end{aligned}
\tag{8.57}
$$

Combining this with (8.55) proves that for all $n \in \mathbb{N} \cap [m, \infty)$ it holds that

$$
\begin{aligned}
\prod_{k=1}^{n} |1 - \gamma_k \alpha| &= \left[ \prod_{k=1}^{m-1} |1 - \gamma_k \alpha| \right] \exp \left( \ln \left( \prod_{k=m}^{n} |1 - \gamma_k \alpha| \right) \right) \\
&\geq \left[ \prod_{k=1}^{m-1} |1 - \gamma_k \alpha| \right] \exp \left( -2\alpha \left[ \sum_{k=1}^{\infty} \gamma_k \right] \right) > 0.
\end{aligned}
\tag{8.58}
$$

Therefore, we obtain that

$$
\liminf_{n \to \infty} \left[ \prod_{k=1}^{n} |1 - \gamma_k \alpha| \right] \geq \left[ \prod_{k=1}^{m-1} |1 - \gamma_k \alpha| \right] \exp \left( -2\alpha \left[ \sum_{k=1}^{\infty} \gamma_k \right] \right) > 0.
\tag{8.59}
$$

This proves item (ii). Furthermore, observe that items (i)–(ii) and the assumption that $\xi \neq \vartheta$ demonstrate that

$$
\begin{aligned}
\liminf_{n\to\infty} \|\Theta_n - \vartheta\| &= \liminf_{n\to\infty} \left\| \left[ \prod_{k=1}^{n} (1 - \gamma_k \alpha) \right] (\xi - \vartheta) \right\| \\
&= \liminf_{n\to\infty} \left( \left| \prod_{k=1}^{n} (1 - \gamma_k \alpha) \right| \|\xi - \vartheta\| \right) \\
&= \|\xi - \vartheta\| \left( \liminf_{n\to\infty} \left[ \prod_{k=1}^{n} |1 - \gamma_k \alpha| \right] \right) > 0.
\end{aligned}
\tag{8.60}
$$

This establishes item (iii). The proof of Lemma 8.1.11 is thus complete. $\qquad\square$

### 8.1.2 Example of a stochastic gradient descent process

Example 8.1.13 below, in particular, provides an error analysis for the SGD optimization method in the case of one specific stochastic optimization problem (see (8.61) below). More general error analyses for the SGD optimization method can, e.g., be found in [16, 17] and the references mentioned therein (cf. Subsection 8.1.4 below).

**Lemma 8.1.12** (Example of a stochastic gradient descent process). *Let $d \in \mathbb{N}$, let $\|\cdot\| \colon \mathbb{R}^d \to [0, \infty)$ be the $d$-dimensional Euclidean norm, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $X_n \colon \Omega \to \mathbb{R}^d$, $n \in \mathbb{N}$, be i.i.d. random variables with $\mathbb{E}[\|X_1\|^2] < \infty$, let $F = (F(\theta, x))_{(\theta,x) \in \mathbb{R}^d \times \mathbb{R}^d} \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ and $f \colon \mathbb{R}^d \to \mathbb{R}$ be the functions which satisfy for all $\theta, x \in \mathbb{R}^d$ that*

$$
F(\theta, x) = \tfrac{1}{2} \|\theta - x\|^2 \qquad and \qquad f(\theta) = \mathbb{E}\big[ F(\theta, X_1) \big],
\tag{8.61}
$$

*and let $\Theta \colon \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ be the stochastic process which satisfies for all $n \in \mathbb{N}$ that $\Theta_0 = 0$ and*

$$
\Theta_n = \Theta_{n-1} - \tfrac{1}{n} (\nabla_\theta F)(\Theta_{n-1}, X_n).
\tag{8.62}
$$

*Then*

*(i) it holds that $\{\theta \in \mathbb{R}^d \colon f(\theta) = \inf_{w \in \mathbb{R}^d} f(w)\} = \{\mathbb{E}[X_1]\}$,*

*(ii) it holds for all $n \in \mathbb{N}$ that $\Theta_n = \tfrac{1}{n}(X_1 + \ldots + X_n)$,*

*(iii) it holds for all $n \in \mathbb{N}$ that*

$$
\big( \mathbb{E}\big[ \|\Theta_n - \mathbb{E}[X_1]\|^2 \big] \big)^{1/2} = \big( \mathbb{E}\big[ \|X_1 - \mathbb{E}[X_1]\|^2 \big] \big)^{1/2} n^{-1/2},
\tag{8.63}
$$

*and*

*(iv) it holds for all $n \in \mathbb{N}$ that*

$$
\mathbb{E}[f(\Theta_n)] - f(\mathbb{E}[X_1]) = \tfrac{1}{2} \mathbb{E}\big[ \|X_1 - \mathbb{E}[X_1]\|^2 \big] n^{-1}.
\tag{8.64}
$$

*Proof of Lemma 8.1.12.* Throughout this proof let $\langle \cdot, \cdot \rangle \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ be the $d$-dimensional Euclidean scalar product. Note that the assumption that $\mathbb{E}[\|X_1\|^2] < \infty$ and Lemma 8.1.2 ensure that for all $\theta \in \mathbb{R}^d$ it holds that

$$
\begin{aligned}
f(\theta) &= \mathbb{E}\big[ F(\theta, X_1) \big] = \tfrac{1}{2} \mathbb{E}\big[ \|X_1 - \theta\|^2 \big] \\
&= \tfrac{1}{2} \big( \mathbb{E}\big[ \|X_1 - \mathbb{E}[X_1]\|^2 \big] + \|\theta - \mathbb{E}[X_1]\|^2 \big).
\end{aligned}
\tag{8.65}
$$

This proves item (i). Moreover, note that Lemma 13.2.4 ensures that for all $\theta, x \in \mathbb{R}^d$ it holds that

$$(\nabla_\theta F)(\theta, x) = \tfrac{1}{2}(2(\theta - x)) = \theta - x. \tag{8.66}$$

This and (8.62) assure that for all $n \in \mathbb{N}$ it holds that

$$\Theta_n = \Theta_{n-1} - \tfrac{1}{n}(\Theta_{n-1} - X_n) = \left(1 - \tfrac{1}{n}\right)\Theta_{n-1} + \tfrac{1}{n}X_n = \tfrac{(n-1)}{n}\Theta_{n-1} + \tfrac{1}{n}X_n. \tag{8.67}$$

Next we claim that for all $n \in \mathbb{N}$ it holds that

$$\Theta_n = \tfrac{1}{n}(X_1 + \ldots + X_n). \tag{8.68}$$

We now prove (8.68) by induction on $n \in \mathbb{N}$. For the base case $n = 1$ note that (8.67) implies that

$$\Theta_1 = \left(\tfrac{0}{1}\right)\Theta_0 + X_1 = \left(\tfrac{1}{1}\right)(X_1). \tag{8.69}$$

This establishes (8.68) in the base case $n = 1$. For the induction step note that (8.67) ensures that for all $n \in \{2, 3, \ldots\}$ with $\Theta_{n-1} = \tfrac{1}{(n-1)}(X_1 + \ldots + X_{n-1})$ it holds that

$$\begin{aligned}
\Theta_n &= \tfrac{(n-1)}{n}\Theta_{n-1} + \tfrac{1}{n}X_n = \left[\tfrac{(n-1)}{n}\right]\left[\tfrac{1}{(n-1)}\right](X_1 + \ldots + X_{n-1}) + \tfrac{1}{n}X_n \\
&= \tfrac{1}{n}(X_1 + \ldots + X_{n-1}) + \tfrac{1}{n}X_n = \tfrac{1}{n}(X_1 + \ldots + X_n).
\end{aligned} \tag{8.70}$$

Induction thus proves (8.68). Next observe that (8.68) establishes item (ii). Moreover, note that Lemma 8.1.4, item (ii), and the fact that $(X_n)_{n \in \mathbb{N}}$ are i.i.d. random variables with $\mathbb{E}[\|X_1\|] < \infty$ ensure that for all $n \in \mathbb{N}$ it holds that

$$\begin{aligned}
\mathbb{E}\big[\|\Theta_n - \mathbb{E}[X_1]\|^2\big] &= \mathbb{E}\big[\|\tfrac{1}{n}(X_1 + \ldots + X_n) - \mathbb{E}[X_1]\|^2\big] \\
&= \mathbb{E}\left[\left\|\tfrac{1}{n}\left[\sum_{k=1}^{n}(X_k - \mathbb{E}[X_1])\right]\right\|^2\right] \\
&= \tfrac{1}{n^2}\left(\mathbb{E}\left[\left\|\sum_{k=1}^{n}(X_k - \mathbb{E}[X_k])\right\|^2\right]\right) \\
&= \tfrac{1}{n^2}\left[\sum_{k=1}^{n}\mathbb{E}\big[\|X_k - \mathbb{E}[X_k]\|^2\big]\right] \\
&= \tfrac{1}{n^2}\big[n\,\mathbb{E}\big[\|X_1 - \mathbb{E}[X_1]\|^2\big]\big] \\
&= \tfrac{\mathbb{E}[\|X_1 - \mathbb{E}[X_1]\|^2]}{n}.
\end{aligned} \tag{8.71}$$

This implies item (iii). It thus remains to prove item (iv). For this note that (8.65) and (8.71) assure that for all $n \in \mathbb{N}$ it holds that

$$\begin{aligned}
&\mathbb{E}[f(\Theta_n)] - f(\mathbb{E}[X_1]) \\
&= \mathbb{E}\big[\tfrac{1}{2}\big(\mathbb{E}\big[\|\mathbb{E}[X_1] - X_1\|^2\big] + \|\Theta_n - \mathbb{E}[X_1]\|^2\big)\big] \\
&\quad - \tfrac{1}{2}\big(\mathbb{E}\big[\|\mathbb{E}[X_1] - X_1\|^2\big] + \|\mathbb{E}[X_1] - \mathbb{E}[X_1]\|^2\big) \\
&= \tfrac{1}{2}\mathbb{E}\big[\|\Theta_n - \mathbb{E}[X_1]\|^2\big] \\
&= \tfrac{1}{2}\mathbb{E}\big[\|X_1 - \mathbb{E}[X_1]\|^2\big]n^{-1}.
\end{aligned} \tag{8.72}$$

This establishes item (iv). The proof of Lemma 8.1.12 is thus complete. $\qquad \square$

## 8.1.3 Examples for stochastic optimization problems

**Example 8.1.13** (Sums of optimiziation problems). *Let $d, N \in \mathbb{N}$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \infty)$, $\xi \in \mathbb{R}^d$, let $f_k \colon \mathbb{R}^d \to \mathbb{R}$, $k \in \{1, 2, \ldots, N\}$, be differentiable functions, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $\mathsf{k}_n \colon \Omega \to \{1, 2, \ldots, N\}$, $n \in \mathbb{N}$, be independent $\mathcal{U}_{\{1,2,\ldots,N\}}$-distributed random variables, let $\Theta \colon \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ be the stochastic process which satisfies for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \gamma_n (\nabla f_{\mathsf{k}_n})(\Theta_{n-1}), \tag{8.73}$$

*and let $F \colon \mathbb{R}^d \times \{1, 2, \ldots, N\} \to \mathbb{R}$ be the function which satisfies for all $\theta \in \mathbb{R}^d$, $k \in \{1, 2, \ldots, N\}$ that*

$$F(\theta, k) = f_k(\theta). \tag{8.74}$$

*Then*

*(i) it holds that $\Theta$ is the stochastic gradient descent process for the loss function $F$ with learning rates $(\gamma_n)_{n \in \mathbb{N}}$, batch sizes $\mathbb{N} \ni n \mapsto 1 \in \mathbb{N}$, initial value $\xi$, and data $(\mathsf{k}_n)_{n \in \mathbb{N}}$ (cf. Definition 8.1.1) and*

*(ii) it holds for all $\theta \in \mathbb{R}^d$ that*

$$\mathbb{E}\big[F(\theta, \mathsf{k}_1)\big] = \frac{1}{N}\left[\sum_{k=1}^{N} f_k(\theta)\right]. \tag{8.75}$$

*Proof of Example 8.1.13.* First, note that (8.74) ensures that for all $n \in \mathbb{N}$ it holds that

$$\Theta_n = \Theta_{n-1} - \gamma_n (\nabla f_{\mathsf{k}_n})(\Theta_{n-1}) = \Theta_{n-1} - \gamma_n (\nabla_\theta F)(\Theta_{n-1}, \mathsf{k}_n). \tag{8.76}$$

Combining this with the assumption that $\Theta_0 = \xi$ proves item (i). Moreover, observe that (8.74) and the assumption that $\mathsf{k}_1$ is a $\mathcal{U}_{\{1,2,\ldots,N\}}$-distributed random variable demonstrate that

$$\mathbb{E}\big[F(\theta, \mathsf{k}_1)\big] = \frac{1}{N}\left[\sum_{k=1}^{N} F(\theta, k)\right] = \frac{1}{N}\left[\sum_{k=1}^{N} f_k(\theta)\right]. \tag{8.77}$$

This establishes item (ii). The proof of Example 8.1.13 is thus complete. $\square$

**Example 8.1.14** (Objective functions induced by data). *Let $d, N, \mathcal{I}, \mathcal{O} \in \mathbb{N}$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \infty)$, $\xi \in \mathbb{R}^d$, $x_1, x_2, \ldots, x_N \in \mathbb{R}^{\mathcal{I}}$, let $\|\cdot\| \colon \mathbb{R}^{\mathcal{O}} \to [0, \infty)$ be the $\mathcal{O}$-dimensional Euclidean norm, let $\Phi \colon \mathbb{R}^{\mathcal{I}} \to \mathbb{R}^{\mathcal{O}}$ be a function, let $u = (u_\theta(x))_{(\theta,x) \in \mathbb{R}^d \times \mathbb{R}^{\mathcal{I}}} \colon \mathbb{R}^d \times \mathbb{R}^{\mathcal{I}} \to \mathbb{R}^{\mathcal{O}}$ be a function which satisfies for every $x \in \mathbb{R}^{\mathcal{I}}$ that the function $\mathbb{R}^d \ni \theta \mapsto u_\theta(x) \in \mathbb{R}^{\mathcal{O}}$ is differentiable, let $F \colon \mathbb{R}^d \times \{1, 2, \ldots, N\} \to \mathbb{R}$ be the function which satisfies for all $\theta \in \mathbb{R}^d$, $k \in \{1, 2, \ldots, N\}$ that*

$$F(\theta, k) = \|u_\theta(x_k) - \Phi(x_k)\|^2, \tag{8.78}$$

*let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $\mathsf{k}_n \colon \Omega \to \{1, 2, \ldots, N\}$, $n \in \mathbb{N}$, be independent $\mathcal{U}_{\{1,2,\ldots,N\}}$-distributed random variables, and let $\Theta \colon \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ be the stochastic process which satisfies for all $n \in \mathbb{N}$ that $\Theta_0 = \xi$ and*

$$\Theta_n = \Theta_{n-1} - \gamma_n (\nabla_\theta F)(\Theta_{n-1}, \mathsf{k}_n) \tag{8.79}$$

*Then*

(i) *it holds that $\Theta$ is the stochastic gradient descent process for the loss function $F$ with learning rates $(\gamma_n)_{n\in\mathbb{N}}$, batch sizes $\mathbb{N} \ni n \mapsto 1 \in \mathbb{N}$, initial value $\xi$, and data $(\mathsf{k}_n)_{n\in\mathbb{N}}$ (cf. Definition 8.1.1) and*

(ii) *it holds for all $\theta \in \mathbb{R}^d$ that*

$$\mathbb{E}\big[F(\theta, \mathsf{k}_1)\big] = \frac{1}{N}\left[\sum_{k=1}^{N}\|u_\theta(x_k) - \Phi(x_k)\|^2\right]. \tag{8.80}$$

*Proof of Example 8.1.14.* Throughout this proof let $f_k \colon \mathbb{R}^d \to \mathbb{R}$, $k \in \{1, 2, \ldots, N\}$, be the functions which satisfy for all $\theta \in \mathbb{R}^d$, $k \in \{1, 2, \ldots, N\}$ that

$$f_k(\theta) = \|u_\theta(x_k) - \Phi(x_k)\|^2. \tag{8.81}$$

Note that Example 8.1.13 (applied with $f_k \curvearrowleft f_k$ for $k \in \{1, 2, \ldots, N\}$ in the notation of Example 8.1.13) establishes items (i)–(ii). The proof of Example 8.1.14 is thus complete. $\square$

### 8.1.4 Convergence rates in dependence of learning rates

The next result, Theorem 8.1.15 below, specifies strong and weak convergence rates for the SGD optimization method in dependence on the asymptotic behavior of the sequence of learning rates. The statement and the proof of Theorem 8.1.15 can be found in [17, Theorem 1.1].

**Theorem 8.1.15** (Convergence rates in dependence of learning rates). *Let $d \in \mathbb{N}$, $\alpha, \gamma, \nu \in (0, \infty)$, $\xi \in \mathbb{R}^d$, let $\langle\cdot, \cdot\rangle \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ be the d-dimensional Euclidean scalar product, let $\|\cdot\| \colon \mathbb{R}^d \to [0, \infty)$ be the d-dimensional Euclidean norm, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $X_n \colon \Omega \to \mathbb{R}^d$, $n \in \mathbb{N}$, be i.i.d. random variables with $\mathbb{E}[\|X_1\|^2] < \infty$ and $\mathbb{P}(X_1 = \mathbb{E}[X_1]) < 1$, let $(r_{\varepsilon,i})_{\varepsilon\in(0,\infty),i\in\{0,1\}} \subseteq \mathbb{R}$ satisfy for all $\varepsilon \in (0, \infty)$, $i \in \{0, 1\}$ that*

$$r_{\varepsilon,i} = \begin{cases} \nu/2 & : \nu < 1 \\ \min\{1/2, \gamma\alpha + (-1)^i\varepsilon\} & : \nu = 1 \\ 0 & : \nu > 1, \end{cases} \tag{8.82}$$

*let $F = (F(\theta, x))_{(\theta,x)\in\mathbb{R}^d\times\mathbb{R}^d} \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ and $f \colon \mathbb{R}^d \to \mathbb{R}$ be the functions which satisfy for all $\theta, x \in \mathbb{R}^d$ that*

$$F(\theta, x) = \tfrac{\alpha}{2}\|\theta - x\|^2 \qquad \text{and} \qquad f(\theta) = \mathbb{E}\big[F(\theta, X_1)\big], \tag{8.83}$$

*and let $\Theta \colon \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ be the stochastic process which satisfies for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad \text{and} \qquad \Theta_n = \Theta_{n-1} - \tfrac{\gamma}{n^\nu}(\nabla_\theta F)(\Theta_{n-1}, X_n). \tag{8.84}$$

*Then*

(i) *there exists a unique $\vartheta \in \mathbb{R}^d$ such that $\{\theta \in \mathbb{R}^d \colon f(\theta) = \inf_{w\in\mathbb{R}^d} f(w)\} = \{\vartheta\}$,*

(ii) *for every $\varepsilon \in (0, \infty)$ there exist $c_0, c_1 \in (0, \infty)$ such that for all $n \in \mathbb{N}$ it holds that*

$$c_0 n^{-r_{\varepsilon,0}} \leq \big(\mathbb{E}\big[\|\Theta_n - \vartheta\|^2\big]\big)^{1/2} \leq c_1 n^{-r_{\varepsilon,1}}, \tag{8.85}$$

*and*

(iii) *for every $\varepsilon \in (0, \infty)$ there exist $C_0, C_1 \in (0, \infty)$ such that for all $n \in \mathbb{N}$ it holds that*

$$C_0 n^{-2r_{\varepsilon,0}} \leq \mathbb{E}[f(\Theta_n)] - f(\vartheta) \leq C_1 n^{-2r_{\varepsilon,1}}. \tag{8.86}$$

## 8.2 The stochastic gradient descent optimization method with classical momentum

In this section we present the SGD optimization method with classical momentum. The idea for classical momentum was first introduced by Polyak for the (deterministic) GD optimization method (see Polyak [25] and Section 14.2 above).

**Definition 8.2.1** (Momentum stochastic gradient descent optimization method). *Let $d \in \mathbb{N}$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \infty)$, $(J_n)_{n \in \mathbb{N}} \subseteq \mathbb{N}$, $(\alpha_n)_{n \in \mathbb{N}} \subseteq [0, 1]$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $(S, \mathcal{S})$ be a measurable space, let $\xi \colon \Omega \to \mathbb{R}^d$ and $X_{n,j} \colon \Omega \to S$, $j \in \{1, 2, \ldots, J_n\}$, $n \in \mathbb{N}$, be random variables, and let $F = (F(\theta, x))_{(\theta, x) \in \mathbb{R}^d \times S} \colon \mathbb{R}^d \times S \to \mathbb{R}$ and $G \colon \mathbb{R}^d \times S \to \mathbb{R}^d$ be functions which satisfy for all $x \in S$, $\theta \in \{v \in \mathbb{R}^d \colon F(\cdot, x) \text{ is differentiable at } v\}$ that*

$$G(\theta, x) = (\nabla_\theta F)(\theta, x). \tag{8.87}$$

*Then we say that $\Theta$ is the momentum stochastic gradient descent process on $((\Omega, \mathcal{F}, \mathbb{P}), (S, \mathcal{S}))$ for the loss function $F$ with generalized gradient $G$, learning rates $(\gamma_n)_{n \in \mathbb{N}}$, batch sizes $(J_n)_{n \in \mathbb{N}}$, momentum decay factors $(\alpha_n)_{n \in \mathbb{N}}$, initial value $\xi$, and data $(X_{n,j})_{j \in \{1,2,\ldots,J_n\}, n \in \mathbb{N}}$ (we say that $\Theta$ is the momentum stochastic gradient descent process for the loss function $F$ with learning rates $(\gamma_n)_{n \in \mathbb{N}}$, batch sizes $(J_n)_{n \in \mathbb{N}}$, momentum decay factors $(\alpha_n)_{n \in \mathbb{N}}$, initial value $\xi$, and data $(X_{n,j})_{j \in \{1,2,\ldots,J_n\}, n \in \mathbb{N}}$) if and only if $\Theta \colon \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ is the function from $\mathbb{N}_0 \times \Omega$ to $\mathbb{R}^d$ which satisfies that there exists a function $\mathbf{m} \colon \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ such that for all $n \in \mathbb{N}$ it holds that*

$$\Theta_0 = \xi, \qquad \mathbf{m}_0 = 0, \tag{8.88}$$

$$\mathbf{m}_n = \alpha_n \mathbf{m}_{n-1} + (1 - \alpha_n) \left[ \frac{1}{J_n} \sum_{j=1}^{J_n} G(\Theta_{n-1}, X_{n,j}) \right], \tag{8.89}$$

$$\text{and} \qquad \Theta_n = \Theta_{n-1} - \gamma_n \mathbf{m}_n. \tag{8.90}$$

## 8.3 The stochastic gradient descent optimization method with Nesterov momentum

Nesterov accelerated stochastic gradient descent (NAG) builds on the idea of classical momentum and attemps to provide some kind of foresight to the scheme. This idea was first introduced by Nesterov as an adaption of the deterministic momentum GD optimization method (see Nesterov [22]).

**Definition 8.3.1** (Nesterov accelerated stochastic gradient descent optimization method). *Let $d \in \mathbb{N}$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \infty)$, $(J_n)_{n \in \mathbb{N}} \subseteq \mathbb{N}$, $(\alpha_n)_{n \in \mathbb{N}} \subseteq [0, 1]$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $(S, \mathcal{S})$ be a measurable space, let $\xi \colon \Omega \to \mathbb{R}^d$ and $X_{n,j} \colon \Omega \to S$, $j \in \{1, 2, \ldots, J_n\}$, $n \in \mathbb{N}$, be random variables, and let $F = (F(\theta, x))_{(\theta, x) \in \mathbb{R}^d \times S} \colon \mathbb{R}^d \times S \to \mathbb{R}$ and $G \colon \mathbb{R}^d \times S \to \mathbb{R}^d$ be functions which satisfy for all $x \in S$, $\theta \in \{v \in \mathbb{R}^d \colon F(\cdot, x) \text{ is differentiable at } v\}$ that*

$$G(\theta, x) = (\nabla_\theta F)(\theta, x). \tag{8.91}$$

*Then we say that $\Theta$ is the Nesterov accelerated stochastic gradient descent process on $((\Omega, \mathcal{F}, \mathbb{P}), (S, \mathcal{S}))$ for the loss function $F$ with generalized gradient $G$, learning rates $(\gamma_n)_{n \in \mathbb{N}}$, batch sizes $(J_n)_{n \in \mathbb{N}}$, momentum decay factors $(\alpha_n)_{n \in \mathbb{N}}$, initial value $\xi$, and data*

$(X_{n,j})_{j\in\{1,2,\ldots,J_n\},n\in\mathbb{N}}$ *(we say that $\Theta$ is the Nesterov accelerated stochastic gradient descent process for the loss function $F$ with learning rates $(\gamma_n)_{n\in\mathbb{N}}$, batch sizes $(J_n)_{n\in\mathbb{N}}$, momentum decay rates $(\alpha_n)_{n\in\mathbb{N}}$, initial value $\xi$, and data $(X_{n,j})_{j\in\{1,2,\ldots,J_n\},n\in\mathbb{N}}$) if and only if $\Theta\colon \mathbb{N}_0\times\Omega\to\mathbb{R}^d$ is the function from $\mathbb{N}_0\times\Omega$ to $\mathbb{R}^d$ which satisfies that there exists a function $\mathbf{m}\colon \mathbb{N}_0\times\Omega\to\mathbb{R}^d$ such that for all $n\in\mathbb{N}$ it holds that*

$$\Theta_0 = \xi, \qquad \mathbf{m}_0 = 0, \tag{8.92}$$

$$\mathbf{m}_n = \alpha_n\mathbf{m}_{n-1} + (1-\alpha_n)\left[\frac{1}{J_n}\sum_{j=1}^{J_n}G\big(\Theta_{n-1} - \gamma_n\alpha_n\mathbf{m}_{n-1},\, X_{n,j}\big)\right], \tag{8.93}$$

$$\text{and} \qquad \Theta_n = \Theta_{n-1} - \gamma_n\mathbf{m}_n. \tag{8.94}$$

# 8.4 The adaptive stochastic gradient descent optimization method (Adagrad)

**Definition 8.4.1** (Adagrad stochastic gradient descent optimization method)*. Let $d\in\mathbb{N}$, $(\gamma_n)_{n\in\mathbb{N}}\subseteq[0,\infty)$, $(J_n)_{n\in\mathbb{N}}\subseteq\mathbb{N}$, $\varepsilon\in(0,\infty)$, let $(\Omega,\mathcal{F},\mathbb{P})$ be a probability space, let $(S,\mathcal{S})$ be a measurable space, let $\xi\colon \Omega\to\mathbb{R}^d$ and $X_{n,j}\colon \Omega\to S$, $j\in\{1,2,\ldots,J_n\}$, $n\in\mathbb{N}$, be random variables, and let $F = (F(\theta,x))_{(\theta,x)\in\mathbb{R}^d\times S}\colon \mathbb{R}^d\times S\to\mathbb{R}$ and $G = (G_1,\ldots,G_d)\colon \mathbb{R}^d\times S\to\mathbb{R}^d$ be functions which satisfy for all $x\in S$, $\theta\in\{v\in\mathbb{R}^d\colon F(\cdot,x)$ is differentiable at $v\}$ that*

$$G(\theta,x) = (\nabla_\theta F)(\theta,x). \tag{8.95}$$

*Then we say that $\Theta$ is the Adagrad stochastic gradient descent process on $((\Omega,\mathcal{F},\mathbb{P}), (S,\mathcal{S}))$ for the loss function $F$ with generalized gradient $G$, learning rates $(\gamma_n)_{n\in\mathbb{N}}$, batch sizes $(J_n)_{n\in\mathbb{N}}$, regularizing factor $\varepsilon$, initial value $\xi$, and data $(X_{n,j})_{j\in\{1,2,\ldots,J_n\},n\in\mathbb{N}}$ (we say that $\Theta$ is the Adagrad stochastic gradient descent process for the loss function $F$ with learning rates $(\gamma_n)_{n\in\mathbb{N}}$, batch sizes $(J_n)_{n\in\mathbb{N}}$, regularizing factor $\varepsilon$, initial value $\xi$, and data $(X_{n,j})_{j\in\{1,2,\ldots,J_n\},n\in\mathbb{N}}$) if and only if it holds that $\Theta = (\Theta^{(1)},\ldots,\Theta^{(d)})\colon \mathbb{N}_0\times\Omega\to\mathbb{R}^d$ is the function from $\mathbb{N}_0\times\Omega$ to $\mathbb{R}^d$ which satisfies for all $n\in\mathbb{N}$, $i\in\{1,2,\ldots,d\}$ that $\Theta_0 = \xi$ and*

$$
\begin{aligned}
\Theta_n^{(i)} = \Theta_{n-1}^{(i)} \\
- \gamma_n\left(\varepsilon + \sum_{k=1}^{n}\left[\tfrac{1}{J_k}\sum_{j=1}^{J_k}G_i(\Theta_{k-1},X_{k,j})\right]^2\right)^{-1/2}\left[\frac{1}{J_n}\sum_{j=1}^{J_n}G_i(\Theta_{n-1},X_{n,j})\right].
\end{aligned}
\tag{8.96}
$$

# 8.5 The root mean square propagation stochastic gradient descent optimization method (RMSprop)

**Definition 8.5.1** (RMSprop stochastic gradient descent optimization method)*. Let $d\in\mathbb{N}$, $(\gamma_n)_{n\in\mathbb{N}}\subseteq[0,\infty)$, $(J_n)_{n\in\mathbb{N}}\subseteq\mathbb{N}$, $(\beta_n)_{n\in\mathbb{N}}\subseteq[0,1]$, $\varepsilon\in(0,\infty)$, let $(\Omega,\mathcal{F},\mathbb{P})$ be a probability space, let $(S,\mathcal{S})$ be a measurable space, let $\xi\colon \Omega\to\mathbb{R}^d$ and $X_{n,j}\colon \Omega\to S$, $j\in\{1,2,\ldots,J_n\}$, $n\in\mathbb{N}$, be random variables, and let $F = (F(\theta,x))_{(\theta,x)\in\mathbb{R}^d\times S}\colon \mathbb{R}^d\times S\to\mathbb{R}$ and $G = (G_1,\ldots,G_d)\colon \mathbb{R}^d\times S\to\mathbb{R}^d$ be functions which satisfy for all $x\in S$, $\theta\in\{v\in\mathbb{R}^d\colon F(\cdot,x)$ is differentiable at $v\}$ that*

$$G(\theta,x) = (\nabla_\theta F)(\theta,x). \tag{8.97}$$

*Then we say that $\Theta$ is the RMSprop stochastic gradient descent process on $((\Omega, \mathcal{F}, \mathbb{P}),$ $(S, \mathcal{S}))$ for the loss function $F$ with generalized gradient $G$, learning rates $(\gamma_n)_{n \in \mathbb{N}}$, batch sizes $(J_n)_{n \in \mathbb{N}}$, second moment decay factors $(\beta_n)_{n \in \mathbb{N}}$, regularizing factor $\varepsilon$, initial value $\xi$, and data $(X_{n,j})_{j \in \{1,2,\ldots,J_n\}, n \in \mathbb{N}}$ (we say that $\Theta$ is the RMSprop stochastic gradient descent process for the loss function $F$ with learning rates $(\gamma_n)_{n \in \mathbb{N}}$, batch sizes $(J_n)_{n \in \mathbb{N}}$, second moment decay factors $(\beta_n)_{n \in \mathbb{N}}$, regularizing factor $\varepsilon$, initial value $\xi$, and data $(X_{n,j})_{j \in \{1,2,\ldots,J_n\}, n \in \mathbb{N}}$) if and only if it holds that $\Theta = (\Theta^{(1)}, \ldots, \Theta^{(d)}) \colon \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ is the function from $\mathbb{N}_0 \times \Omega$ to $\mathbb{R}^d$ which satisfies that there exists a function $\mathbb{M} = (\mathbb{M}^{(1)}, \ldots, \mathbb{M}^{(d)}) \colon \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ such that for all $n \in \mathbb{N}$, $i \in \{1, 2, \ldots, d\}$ it holds that*

$$\Theta_0 = \xi, \qquad \mathbb{M}_0 = 0, \tag{8.98}$$

$$\mathbb{M}_n^{(i)} = \beta_n \, \mathbb{M}_{n-1}^{(i)} + (1 - \beta_n) \left[ \frac{1}{J_n} \sum_{j=1}^{J_n} G_i(\Theta_{n-1}, X_{n,j}) \right]^2, \tag{8.99}$$

$$and \qquad \Theta_n^{(i)} = \Theta_{n-1}^{(i)} - \frac{\gamma_n}{\left[ \varepsilon + \mathbb{M}_n^{(i)} \right]^{1/2}} \left[ \frac{1}{J_n} \sum_{j=1}^{J_n} G_i(\Theta_{n-1}, X_{n,j}) \right]. \tag{8.100}$$

Hinton et al. [14] suggests the choice that for all $n \in \mathbb{N}$ it holds that

$$\beta_n = 0.9 \tag{8.101}$$

as default values for the second moment decay factors $(\beta_n)_{n \in \mathbb{N}} \subseteq [0, 1]$ in Definition 8.5.1. This default value in used several machine learning libraries that implement RMSprop (see, e.g., Tensorflow [28] and Lasagne [20]).

## 8.6 The Adadelta stochastic gradient descent optimization method

The Adadelta SGD optimization method was proposed in Zeiler [30]. It is a extension of RMSprop SGD optimization method. Like the RMSprop SGD optimization method, the Adadelta SGD optimization method adapts the learning rate for every component separately. To do this, the Adadelta SGD optimization method uses two exponentially decaying averages: one over the squares of the past partial derivatives and another one over the squares of the past increments (cf. Definition 8.6.1 below).

**Definition 8.6.1** (Adadelta stochastic gradient descent optimization method). *Let $d \in \mathbb{N}$, $(J_n)_{n \in \mathbb{N}} \subseteq \mathbb{N}$, $(\beta_n)_{n \in \mathbb{N}}, (\delta_n)_{n \in \mathbb{N}} \subseteq [0, 1]$, $\varepsilon \in (0, \infty)$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $(S, \mathcal{S})$ be a measurable space, let $\xi \colon \Omega \to \mathbb{R}^d$ and $X_{n,j} \colon \Omega \to S$, $j \in \{1, 2, \ldots, J_n\}$, $n \in \mathbb{N}$, be random variables, and let $F = (F(\theta, x))_{(\theta, x) \in \mathbb{R}^d \times S} \colon \mathbb{R}^d \times S \to \mathbb{R}$ and $G = (G_1, \ldots, G_d) \colon \mathbb{R}^d \times S \to \mathbb{R}^d$ be functions which satisfy for all $x \in S$, $\theta \in \{v \in \mathbb{R}^d \colon F(\cdot, x)$ is differentiable at $v\}$ that*

$$G(\theta, x) = (\nabla_\theta F)(\theta, x). \tag{8.102}$$

*Then we say that $\Theta$ is the Adadelta stochastic gradient descent process on $((\Omega, \mathcal{F}, \mathbb{P}),$ $(S, \mathcal{S}))$ for the loss function $F$ with generalized gradient $G$, batch sizes $(J_n)_{n \in \mathbb{N}}$, second moment decay factors $(\beta_n)_{n \in \mathbb{N}}$, delta decay factors $(\delta_n)_{n \in \mathbb{N}}$, regularizing factor $\varepsilon$, initial value $\xi$, and data $(X_{n,j})_{j \in \{1,2,\ldots,J_n\}, n \in \mathbb{N}}$ (we say that $\Theta$ is the Adadelta stochastic gradient descent process for the loss function $F$ with batch sizes $(J_n)_{n \in \mathbb{N}}$, second moment*

*decay factors $(\beta_n)_{n\in\mathbb{N}}$, delta decay factors $(\delta_n)_{n\in\mathbb{N}}$, regularizing factor $\varepsilon$, initial value $\xi$, and data $(X_{n,j})_{j\in\{1,2,\dots,J_n\},n\in\mathbb{N}}$) if and only if it holds that $\Theta = (\Theta^{(1)}, \dots, \Theta^{(d)}): \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ is the function from $\mathbb{N}_0 \times \Omega$ to $\mathbb{R}^d$ which satisfies that there exist functions $\mathbb{M} = (\mathbb{M}^{(1)}, \dots, \mathbb{M}^{(d)}), \Delta = (\Delta^{(1)}, \dots, \Delta^{(d)}): \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ such that for all $n \in \mathbb{N}$, $i \in \{1, 2, \dots, d\}$ it holds that*

$$\Theta_0 = \xi, \qquad \mathbb{M}_0 = 0, \qquad \Delta_0 = 0, \tag{8.103}$$

$$\mathbb{M}_n^{(i)} = \beta_n\, \mathbb{M}_{n-1}^{(i)} + (1 - \beta_n) \left[ \frac{1}{J_n} \sum_{j=1}^{J_n} G_i(\Theta_{n-1}, X_{n,j}) \right]^2, \tag{8.104}$$

$$\Theta_n^{(i)} = \Theta_{n-1}^{(i)} - \left( \frac{\varepsilon + \Delta_{n-1}^{(i)}}{\varepsilon + \mathbb{M}_n^{(i)}} \right)^{1/2} \left[ \frac{1}{J_n} \sum_{j=1}^{J_n} G_i(\Theta_{n-1}, X_{n,j}) \right], \tag{8.105}$$

$$and \qquad \Delta_n^{(i)} = \delta_n \Delta_{n-1}^{(i)} + (1 - \delta_n) \big| \Theta_n^{(i)} - \Theta_{n-1}^{(i)} \big|^2. \tag{8.106}$$

## 8.7 The adaptive moment estimation stochastic gradient descent optimization method (Adam stochastic gradient descent optimization method)

**Definition 8.7.1** (Adam stochastic gradient descent optimization method). *Let $d \in \mathbb{N}$, $(\gamma_n)_{n\in\mathbb{N}} \subseteq [0, \infty)$, $(J_n)_{n\in\mathbb{N}} \subseteq \mathbb{N}$, $(\alpha_n)_{n\in\mathbb{N}}$, $(\beta_n)_{n\in\mathbb{N}} \subseteq [0,1)$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $(S, \mathcal{S})$ be a measurable space, let $\xi: \Omega \to \mathbb{R}^d$ and $X_{n,j}: \Omega \to S$, $j \in \{1, 2, \dots, J_n\}$, $n \in \mathbb{N}$, be random variables, and let $F = (F(\theta, x))_{(\theta,x)\in\mathbb{R}^d\times S}: \mathbb{R}^d \times S \to \mathbb{R}$ and $G = (G_1, \dots, G_d): \mathbb{R}^d \times S \to \mathbb{R}^d$ be functions which satisfy for all $x \in S$, $\theta \in \{v \in \mathbb{R}^d: F(\cdot, x) \text{ is differentiable at } v\}$ that*

$$G(\theta, x) = (\nabla_\theta F)(\theta, x). \tag{8.107}$$

*Then we say that $\Theta$ is the Adam stochastic gradient descent process on $((\Omega, \mathcal{F}, \mathbb{P}), (S, \mathcal{S}))$ for the loss function $F$ with generalized gradient $G$, learning rates $(\gamma_n)_{n\in\mathbb{N}}$, batch sizes $(J_n)_{n\in\mathbb{N}}$, momentum decay factors $(\alpha_n)_{n\in\mathbb{N}}$, second moment decay factors $(\beta_n)_{n\in\mathbb{N}}$, initial value $\xi$, and data $(X_{n,j})_{j\in\{1,2,\dots,J_n\},n\in\mathbb{N}}$ (we say that $\Theta$ is the Adam stochastic gradient descent process for the loss function $F$ with learning rates $(\gamma_n)_{n\in\mathbb{N}}$, batch sizes $(J_n)_{n\in\mathbb{N}}$, momentum decay factors $(\alpha_n)_{n\in\mathbb{N}}$, second moment decay factors $(\beta_n)_{n\in\mathbb{N}}$, initial value $\xi$, and data $(X_{n,j})_{j\in\{1,2,\dots,J_n\},n\in\mathbb{N}}$) if and only if it holds that $\Theta = (\Theta^{(1)}, \dots, \Theta^{(d)}): \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ is the function from $\mathbb{N}_0 \times \Omega$ to $\mathbb{R}^d$ which satisfies that there exist functions $\mathbf{m} = (\mathbf{m}^{(1)}, \dots, \mathbf{m}^{(d)}), \mathbb{M} = (\mathbb{M}^{(1)}, \dots, \mathbb{M}^{(d)}): \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ such that for all $n \in \mathbb{N}$, $i \in \{1, 2, \dots, d\}$ it holds that*

$$\Theta_0 = \xi, \qquad \mathbf{m}_0 = 0, \qquad \mathbb{M}_0 = 0, \tag{8.108}$$

$$\mathbf{m}_n = \alpha_n\, \mathbf{m}_{n-1} + (1 - \alpha_n) \left[ \frac{1}{J_n} \sum_{j=1}^{J_n} G(\Theta_{n-1}, X_{n,j}) \right], \tag{8.109}$$

$$\mathbb{M}_n^{(i)} = \beta_n\, \mathbb{M}_{n-1}^{(i)} + (1 - \beta_n) \left[ \frac{1}{J_n} \sum_{j=1}^{J_n} G_i(\Theta_{n-1}, X_{n,j}) \right]^2, \tag{8.110}$$

$$\text{and} \qquad \Theta_n^{(i)} = \Theta_{n-1}^{(i)} - \gamma_n \left[ \varepsilon + \left[ \frac{\mathbb{M}_n^{(i)}}{(1 - \prod_{l=1}^n \beta_l)} \right]^{1/2} \right]^{-1} \left[ \frac{\mathbf{m}_n^{(i)}}{(1 - \prod_{l=1}^n \alpha_l)} \right]. \qquad (8.111)$$

Kingma & Ba [18] suggests the choice that for all $n \in \mathbb{N}$ it holds that that

$$\gamma_n = 0.001, \qquad \alpha_n = 0.9, \qquad \beta_n = 0.999, \qquad \text{and} \qquad \varepsilon = 10^{-8} \qquad (8.112)$$

as default values for $(\gamma_n)_{n \in \mathbb{N}}$, $(\alpha_n)_{n \in \mathbb{N}}$, $(\beta_n)_{n \in \mathbb{N}}$, and $\varepsilon$ in Definition 8.7.1.

# Chapter 9

# Generalization error

## 9.1 Concentration inequalities for random variables

This section is inspired by Duchi [8].

### 9.1.1 Markov's inequality

**Lemma 9.1.1** (Markov inequality). *Let $(\Omega, \mathcal{F}, \mu)$ be a measure space, let $X \colon \Omega \to [0, \infty)$ be an $\mathcal{F}/\mathcal{B}([0, \infty))$-measurable function, and let $\varepsilon \in (0, \infty)$. Then*

$$\mu(X \geq \varepsilon) \leq \frac{\int_\Omega X \, d\mu}{\varepsilon}. \tag{9.1}$$

*Proof of Lemma 9.1.1.* Observe that the fact that $X \geq 0$ proves that

$$\mathbb{1}_{\{X \geq \varepsilon\}} = \frac{\varepsilon \mathbb{1}_{\{X \geq \varepsilon\}}}{\varepsilon} \leq \frac{X \mathbb{1}_{\{X \geq \varepsilon\}}}{\varepsilon} \leq \frac{X}{\varepsilon}. \tag{9.2}$$

Hence, we obtain that

$$\mu(X \geq \varepsilon) = \int_\Omega \mathbb{1}_{\{X \geq \varepsilon\}} \, d\mu \leq \frac{\int_\Omega X \, d\mu}{\varepsilon}. \tag{9.3}$$

The proof of Lemma 9.1.1 is thus complete. $\qquad\square$

### 9.1.2 A first concentration inequality

#### 9.1.2.1 On the variance of bounded random variables

**Lemma 9.1.2.** *Let $x \in [0, 1]$, $y \in \mathbb{R}$. Then*

$$(x - y)^2 \leq (1 - x)y^2 + x(1 - y)^2. \tag{9.4}$$

*Proof of Lemma 9.1.2.* Observe that the assumption that $x \in [0, 1]$ assures that

$$(1 - x)y^2 + x(1 - y)^2 = y^2 - xy^2 + x - 2xy + xy^2 \geq y^2 + x^2 - 2xy = (x - y)^2. \tag{9.5}$$

This establishes (9.4). The proof of Lemma 9.1.2 is thus complete. $\qquad\square$

**Lemma 9.1.3.** *It holds that $\sup_{p \in \mathbb{R}} p(1 - p) = \frac{1}{4}$.*

*Proof of Lemma 9.1.3.* Throughout this proof let $f\colon \mathbb{R} \to \mathbb{R}$ be the function which satisfies for all $p \in \mathbb{R}$ that $f(p) = p(1-p)$. Observe that the fact that $\forall\, p \in \mathbb{R}\colon f'(p) = 1 - 2p$ implies that $\{p \in \mathbb{R}\colon f'(p) = 0\} = \{\nicefrac{1}{2}\}$. Combining this with the fact that $f$ is a strictly concave function implies that

$$\sup_{p \in \mathbb{R}} p(1-p) = \sup_{p \in \mathbb{R}} f(p) = f(\nicefrac{1}{2}) = \nicefrac{1}{4}. \tag{9.6}$$

The proof of Lemma 9.1.3 is thus complete. $\qquad\square$

**Lemma 9.1.4.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and let $X\colon \Omega \to [0,1]$ be a random variable. Then*

$$\mathrm{Var}(X) \le \nicefrac{1}{4}. \tag{9.7}$$

*Proof of Lemma 9.1.4.* Observe that Lemma 9.1.2 implies that

$$\begin{aligned}
\mathrm{Var}(X) = \mathbb{E}\big[(X - \mathbb{E}[X])^2\big] &\le \mathbb{E}\big[(1-X)(\mathbb{E}[X])^2 + X(1 - \mathbb{E}[X])^2\big] \\
&= (1 - \mathbb{E}[X])(\mathbb{E}[X])^2 + \mathbb{E}[X](1 - \mathbb{E}[X])^2 \\
&= (1 - \mathbb{E}[X])\mathbb{E}[X](\mathbb{E}[X] + (1 - \mathbb{E}[X])) \\
&= (1 - \mathbb{E}[X])\mathbb{E}[X].
\end{aligned} \tag{9.8}$$

This and Lemma 9.1.3 demonstrate that $\mathrm{Var}(X) \le \nicefrac{1}{4}$. The proof of Lemma 9.1.4 is thus complete. $\qquad\square$

**Lemma 9.1.5.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $a \in \mathbb{R}$, $b \in [a, \infty)$, and let $X\colon \Omega \to [a,b]$ be a random variable. Then*

$$\mathrm{Var}(X) \le \frac{(b-a)^2}{4}. \tag{9.9}$$

*Proof of Lemma 9.1.5.* Throughout this proof assume w.l.o.g. that $a < b$. Observe that Lemma 9.1.4 implies that

$$\begin{aligned}
\mathrm{Var}(X) = \mathbb{E}\big[(X - \mathbb{E}[X])^2\big] = (b-a)^2\, \mathbb{E}\!\left[\left(\tfrac{X - a - (\mathbb{E}[X] - a)}{b-a}\right)^2\right] \\
= (b-a)^2\, \mathbb{E}\!\left[\left(\tfrac{X-a}{b-a} - \mathbb{E}\!\left[\tfrac{X-a}{b-a}\right]\right)^2\right] \\
= (b-a)^2\, \mathrm{Var}\!\left(\tfrac{X-a}{b-a}\right) \le (b-a)^2\big(\tfrac{1}{4}\big) = \frac{(b-a)^2}{4}.
\end{aligned} \tag{9.10}$$

The proof of Lemma 9.1.5 is thus complete. $\qquad\square$

### 9.1.2.2   A concentration inequality

**Lemma 9.1.6.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $N \in \mathbb{N}$, $\varepsilon \in (0, \infty)$, $a_1, a_2, \ldots, a_N \in \mathbb{R}$, $b_1 \in [a_1, \infty)$, $b_2 \in [a_2, \infty)$, $\ldots$, $b_N \in [a_N, \infty)$, and let $X_n\colon \Omega \to [a_n, b_n]$, $n \in \{1, 2, \ldots, N\}$, be independent random variables. Then*

$$\mathbb{P}\!\left(\left|\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right| \ge \varepsilon\right) \le \frac{\sum_{n=1}^{N}(b_n - a_n)^2}{4\varepsilon^2}. \tag{9.11}$$

*Proof of Lemma 9.1.6.* Note that Lemma 9.1.1 assures that

$$\mathbb{P}\left(\left|\sum_{n=1}^{N}\big(X_n - \mathbb{E}[X_n]\big)\right| \geq \varepsilon\right) = \mathbb{P}\left(\left|\sum_{n=1}^{N}\big(X_n - \mathbb{E}[X_n]\big)\right|^2 \geq \varepsilon^2\right)$$
$$\leq \frac{\mathbb{E}\left[\left|\sum_{n=1}^{N}\big(X_n - \mathbb{E}[X_n]\big)\right|^2\right]}{\varepsilon^2}. \tag{9.12}$$

In addition, note that the assumption that $X_n\colon \Omega \to [a_n, b_n]$, $n \in \{1, 2, \ldots, N\}$, are independent variables and Lemma 9.1.5 demonstrate that

$$\mathbb{E}\left[\left|\sum_{n=1}^{N}\big(X_n - \mathbb{E}[X_n]\big)\right|^2\right] = \sum_{n,m=1}^{N} \mathbb{E}\left[\big(X_n - \mathbb{E}[X_n]\big)\big(X_m - \mathbb{E}[X_m]\big)\right]$$
$$= \sum_{n=1}^{N} \mathbb{E}\left[\big(X_n - \mathbb{E}[X_n]\big)^2\right] \leq \frac{\sum_{n=1}^{N}(b_n - a_n)^2}{4}. \tag{9.13}$$

Combining this with (9.12) establishes

$$\mathbb{P}\left(\left|\sum_{n=1}^{N}\big(X_n - \mathbb{E}[X_n]\big)\right| \geq \varepsilon\right) \leq \frac{\sum_{n=1}^{N}(b_n - a_n)^2}{4\varepsilon^2} \tag{9.14}$$

The proof of Lemma 9.1.6 is thus complete. □

### 9.1.3 Moment-generating functions

**Definition 9.1.7.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and let $X\colon \Omega \to \mathbb{R}$ be a random variable. Then we denote by $\mathbb{M}_{X,\mathbb{P}}\colon \mathbb{R} \to [0, \infty]$ (we denote by $\mathbb{M}_X\colon \mathbb{R} \to [0, \infty]$) the function which satisfies for all $t \in \mathbb{R}$ that*

$$\mathbb{M}_{X,\mathbb{P}}(t) = \mathbb{E}\big[e^{tX}\big] \tag{9.15}$$

*and we call $\mathbb{M}_{X,\mathbb{P}}$ the moment-generating function of $X$ with respect to $\mathbb{P}$ (we call $\mathbb{M}_{X,\mathbb{P}}$ the moment-generating function of $X$).*

#### 9.1.3.1 Moment-generation function for the sum of independent random variables

**Lemma 9.1.8.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $t \in \mathbb{R}$, $N \in \mathbb{N}$, and let $X_n\colon \Omega \to \mathbb{R}$, $n \in \{1, 2, \ldots, N\}$, be independent random variables. Then*

$$\mathbb{M}_{\sum_{n=1}^{N} X_n}(t) = \prod_{n=1}^{N} \mathbb{M}_{X_n}(t). \tag{9.16}$$

*Proof of Lemma 9.1.8.* Observe that Fubini's theorem ensures that for all $t \in \mathbb{R}$ it holds that

$$\mathbb{M}_{\sum_{n=1}^{N} X_n}(t) = \mathbb{E}\left[e^{t\left(\sum_{n=1}^{N} X_n\right)}\right] = \mathbb{E}\left[\prod_{n=1}^{N} e^{tX_n}\right] = \prod_{n=1}^{N} \mathbb{E}\big[e^{tX_n}\big] = \prod_{n=1}^{N} \mathbb{M}_{X_n}(t). \tag{9.17}$$

The proof of Lemma 9.1.8 is thus complete. □

### 9.1.4 Chernoff bounds

#### 9.1.4.1 Probability to cross a barrier

**Proposition 9.1.9.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $X \colon \Omega \to \mathbb{R}$ be a random variable, and let $\varepsilon \in \mathbb{R}$. Then*

$$\mathbb{P}(X \geq \varepsilon) \leq \inf_{\lambda \in [0, \infty)} \left( e^{-\lambda \varepsilon} \, \mathbb{E}\big[ e^{\lambda X} \big] \right) = \inf_{\lambda \in [0, \infty)} \left( e^{-\lambda \varepsilon} \, \mathrm{M}_X(\lambda) \right). \tag{9.18}$$

*Proof of Proposition 9.1.9.* Note that Lemma 9.1.1 ensures that for all $\lambda \in [0, \infty)$ it holds that

$$\mathbb{P}(X \geq \varepsilon) \leq \mathbb{P}(\lambda X \geq \lambda \varepsilon) = \mathbb{P}(\exp(\lambda X) \geq \exp(\lambda \varepsilon)) \leq \frac{\mathbb{E}[\exp(\lambda X)]}{\exp(\lambda \varepsilon)} = e^{-\lambda \varepsilon} \, \mathbb{E}\big[ e^{\lambda X} \big]. \tag{9.19}$$

The proof of Proposition 9.1.9 is thus complete. $\square$

**Corollary 9.1.10.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $X \colon \Omega \to \mathbb{R}$ be a random variable, and let $c, \varepsilon \in \mathbb{R}$. Then*

$$\mathbb{P}(X \geq c + \varepsilon) \leq \inf_{\lambda \in [0, \infty)} \left( e^{-\lambda \varepsilon} \, \mathrm{M}_{X-c}(\lambda) \right). \tag{9.20}$$

*Proof of Corollary 9.1.10.* Throughout this proof let $Y \colon \Omega \to \mathbb{R}$ satisfy

$$Y = X - c. \tag{9.21}$$

Observe that Proposition 9.1.9 and (9.21) ensure that

$$\mathbb{P}(X - c \geq \varepsilon) = \mathbb{P}(Y \geq \varepsilon) \leq \inf_{\lambda \in [0, \infty)} \left( e^{-\lambda \varepsilon} \, \mathrm{M}_Y(\lambda) \right) = \inf_{\lambda \in [0, \infty)} \left( e^{-\lambda \varepsilon} \, \mathrm{M}_{X-c}(\lambda) \right). \tag{9.22}$$

The proof of Corollary 9.1.10 is thus complete. $\square$

**Corollary 9.1.11.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $X \colon \Omega \to \mathbb{R}$ be a random variable with $\mathbb{E}[|X|] < \infty$, and let $\varepsilon \in \mathbb{R}$. Then*

$$\mathbb{P}(X \geq \mathbb{E}[X] + \varepsilon) \leq \inf_{\lambda \in [0, \infty)} \left( e^{-\lambda \varepsilon} \, \mathrm{M}_{X - \mathbb{E}[X]}(\lambda) \right). \tag{9.23}$$

*Proof of Corollary 9.1.11.* Observe that Corollary 9.1.10 (applied with $c \curvearrowleft \mathbb{E}[X]$ in the notation of Corollary 9.1.10) establishes (9.23). The proof of Corollary 9.1.11 is thus complete. $\square$

#### 9.1.4.2 Probability to fall below a barrier

**Corollary 9.1.12.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $X \colon \Omega \to \mathbb{R}$ be a random variable, and let $c, \varepsilon \in \mathbb{R}$. Then*

$$\mathbb{P}(X \leq c - \varepsilon) \leq \inf_{\lambda \in [0, \infty)} \left( e^{-\lambda \varepsilon} \, \mathrm{M}_{c-X}(\lambda) \right). \tag{9.24}$$

*Proof of Corollary 9.1.12.* Throughout this proof let $\mathfrak{c} \in \mathbb{R}$ satisfy $\mathfrak{c} = -c$ and let $\mathfrak{X} \colon \Omega \to \mathbb{R}$ satisfy

$$\mathfrak{X} = -X. \tag{9.25}$$

Observe that Corollary 9.1.10 and (9.25) ensure that

$$\mathbb{P}(X \leq c - \varepsilon) = \mathbb{P}(-X \geq -c + \varepsilon) = \mathbb{P}(\mathfrak{X} \geq \mathfrak{c} + \varepsilon) \leq \inf_{\lambda \in [0, \infty)} \left( e^{-\lambda \varepsilon} \, \mathrm{M}_{\mathfrak{X} - \mathfrak{c}}(\lambda) \right)$$
$$= \inf_{\lambda \in [0, \infty)} \left( e^{-\lambda \varepsilon} \, \mathrm{M}_{c-X}(\lambda) \right). \tag{9.26}$$

The proof of Corollary 9.1.12 is thus complete. $\square$

### 9.1.4.3   Sums of independent random variables

**Corollary 9.1.13.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $\varepsilon \in \mathbb{R}$, $N \in \mathbb{N}$, and let $X_n \colon \Omega \to \mathbb{R}$, $n \in \{1, 2, \ldots, N\}$, be independent random variables with $\max_{n \in \{1,2,\ldots,N\}} \mathbb{E}[|X_n|] < \infty$. Then*

$$\mathbb{P}\left(\left[\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right] \geq \varepsilon\right) \leq \inf_{\lambda \in [0,\infty)}\left(e^{-\lambda\varepsilon}\left[\prod_{n=1}^{N}\mathbb{M}_{X_n - \mathbb{E}[X_n]}(\lambda)\right]\right). \tag{9.27}$$

*Proof of Corollary 9.1.13.* Throughout this proof let $Y_n \colon \Omega \to \mathbb{R}$, $n \in \{1, 2, \ldots, N\}$, satisfy for all $n \in \{1, 2, \ldots, N\}$ that

$$Y_n = X_n - \mathbb{E}[X_n]. \tag{9.28}$$

Observe that Proposition 9.1.9, Lemma 9.1.8, and (9.28) ensure that

$$\mathbb{P}\left(\left[\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right] \geq \varepsilon\right) = \mathbb{P}\left(\left[\sum_{n=1}^{N}Y_n\right] \geq \varepsilon\right) \leq \inf_{\lambda \in [0,\infty)}\left(e^{-\lambda\varepsilon}\,\mathbb{M}_{\sum_{n=1}^{N}Y_n}(\lambda)\right)$$
$$= \inf_{\lambda \in [0,\infty)}\left(e^{-\lambda\varepsilon}\left[\prod_{n=1}^{N}\mathbb{M}_{Y_n}(\lambda)\right]\right) = \inf_{\lambda \in [0,\infty)}\left(e^{-\lambda\varepsilon}\left[\prod_{n=1}^{N}\mathbb{M}_{X_n - \mathbb{E}[X_n]}(\lambda)\right]\right). \tag{9.29}$$

The proof of Corollary 9.1.13 is thus complete. $\square$

## 9.1.5   Hoeffding's inequality

### 9.1.5.1   On the moment-generating function for bounded random variables

**Lemma 9.1.14.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $\lambda, a \in \mathbb{R}$, $b \in (a, \infty)$, $p \in [0,1]$ satisfy $p = \frac{-a}{(b-a)}$, let $X \colon \Omega \to [a,b]$ be a random variable with $\mathbb{E}[X] = 0$, and let $\phi \colon \mathbb{R} \to \mathbb{R}$ satisfy for all $x \in \mathbb{R}$ that $\phi(x) = \ln(1 - p + pe^x) - px$. Then*

$$\mathbb{E}\left[e^{\lambda X}\right] \leq e^{\phi(\lambda(b-a))}. \tag{9.30}$$

*Proof of Lemma 9.1.14.* Observe that for all $x \in \mathbb{R}$ it holds that

$$x(b-a) = bx - ax = [ab - ax] + [bx - ab] = [a(b-x)] + [b(x-a)]$$
$$= a(b-x) + b[b - a - b + x] = a(b-x) + b[(b-a) - (b-x)]. \tag{9.31}$$

Hence, we obtain that for all $x \in \mathbb{R}$ it holds that

$$x = a\left(\frac{b-x}{b-a}\right) + b\left[1 - \left(\frac{b-x}{b-a}\right)\right]. \tag{9.32}$$

This implies that for all $x \in \mathbb{R}$ it holds that

$$\lambda x = \left(\frac{b-x}{b-a}\right)\lambda a + \left[1 - \left(\frac{b-x}{b-a}\right)\right]\lambda b. \tag{9.33}$$

The fact that the function $\mathbb{R} \ni x \mapsto e^x \in \mathbb{R}$ is convex hence demonstrates that for all $x \in [a,b]$ it holds that

$$e^{\lambda x} = \exp\left(\left(\frac{b-x}{b-a}\right)\lambda a + \left[1 - \left(\frac{b-x}{b-a}\right)\right]\lambda b\right) \leq \left(\frac{b-x}{b-a}\right)e^{\lambda a} + \left[1 - \left(\frac{b-x}{b-a}\right)\right]e^{\lambda b}. \tag{9.34}$$

The assumption that $\mathbb{E}[X] = 0$ therefore assures that

$$\mathbb{E}\big[e^{\lambda X}\big] \le \left(\frac{b}{b-a}\right)e^{\lambda a} + \left[1 - \left(\frac{b}{b-a}\right)\right]e^{\lambda b}. \tag{9.35}$$

Combining this with the fact that

$$\frac{b}{(b-a)} = 1 - \left[1 - \left(\frac{b}{(b-a)}\right)\right] = 1 - \left[\left(\frac{(b-a)}{(b-a)}\right) - \left(\frac{b}{(b-a)}\right)\right] = 1 - \left[\frac{-a}{(b-a)}\right] = 1 - p \tag{9.36}$$

demonstrates that

$$\mathbb{E}\big[e^{\lambda X}\big] \le \left(\frac{b}{b-a}\right)e^{\lambda a} + \left[1 - \left(\frac{b}{b-a}\right)\right]e^{\lambda b} = (1-p)e^{\lambda a} + [1 - (1-p)]e^{\lambda b} = (1-p)e^{\lambda a} + p\, e^{\lambda b}$$
$$= \big[(1-p) + p\, e^{\lambda(b-a)}\big]e^{\lambda a}. \tag{9.37}$$

Moreover, note that the assumption that $p = \frac{-a}{(b-a)}$ shows that $p(b-a) = -a$. Hence, we obtain that $a = -p(b-a)$. This and (9.37) assure that

$$\mathbb{E}\big[e^{\lambda X}\big] \le \big[(1-p) + p\, e^{\lambda(b-a)}\big]e^{-p\lambda(b-a)} = \exp\big(\ln\big(\big[(1-p) + p\, e^{\lambda(b-a)}\big]e^{-p\lambda(b-a)}\big)\big)$$
$$= \exp\big(\ln\big((1-p) + p\, e^{\lambda(b-a)}\big) - p\lambda(b-a)\big) = \exp(\phi(\lambda(b-a))). \tag{9.38}$$

The proof of Lemma 9.1.14 is thus complete. $\qquad\square$

### 9.1.5.2 Hoeffding's lemma

**Lemma 9.1.15.** *Let $p \in [0,1]$ and let $\phi\colon \mathbb{R} \to \mathbb{R}$ satisfy for all $x \in \mathbb{R}$ that $\phi(x) = \ln(1 - p + pe^x) - px$. Then it holds for all $x \in \mathbb{R}$ that $\phi(x) \le \frac{x^2}{8}$.*

*Proof of Lemma 9.1.15.* Observe that the fundamental theorem of calculus ensures that for all $x \in \mathbb{R}$ it holds that

$$\phi(x) = \phi(0) + \int_0^x \phi'(y)\,\mathrm{d}y = \phi(0) + \phi'(0)x + \int_0^x \int_0^y \phi''(z)\,dz\,\mathrm{d}y \le \phi(0) + \phi'(0)x + \frac{x^2}{2}\left[\sup_{z \in \mathbb{R}} \phi''(z)\right]. \tag{9.39}$$

Moreover, note that for all $x \in \mathbb{R}$ it holds that

$$\phi'(x) = \left[\frac{pe^x}{1 - p + pe^x}\right] - p \qquad \text{and} \qquad \phi''(x) = \left[\frac{pe^x}{1 - p + pe^x}\right] - \left[\frac{p^2 e^{2x}}{(1 - p + pe^x)^2}\right]. \tag{9.40}$$

Hence, we obtain that

$$\phi'(0) = \left[\frac{p}{1 - p + p}\right] - p = 0. \tag{9.41}$$

In the next step we combine (9.40) and the fact that for all $a \in \mathbb{R}$ it holds that

$$a(1-a) = a - a^2 = -\left[a^2 - 2a\big[\tfrac{1}{2}\big] + \big[\tfrac{1}{2}\big]^2\right] + \big[\tfrac{1}{2}\big]^2 = \tfrac{1}{4} - \big[a - \tfrac{1}{2}\big]^2 \le \tfrac{1}{4} \tag{9.42}$$

to obtain that for all $x \in \mathbb{R}$ it holds that $\phi''(x) \le \frac{1}{4}$. This, (9.39), and (9.41) ensure that for all $x \in \mathbb{R}$ it holds that

$$\phi(x) \le \phi(0) + \phi'(0)x + \frac{x^2}{2}\left[\sup_{z \in \mathbb{R}} \phi''(z)\right] = \phi(0) + \frac{x^2}{2}\left[\sup_{z \in \mathbb{R}} \phi''(z)\right] \le \phi(0) + \frac{x^2}{8} = \frac{x^2}{8}. \tag{9.43}$$

The proof of Lemma 9.1.15 is thus complete. $\qquad\square$

**Lemma 9.1.16.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $a \in \mathbb{R}$, $b \in [a, \infty)$, $\lambda \in \mathbb{R}$, and let $X\colon \Omega \to [a,b]$ be a random variable with $\mathbb{E}[X] = 0$. Then*

$$\mathbb{E}\big[\exp(\lambda X)\big] \le \exp\Big(\tfrac{\lambda^2(b-a)^2}{8}\Big). \tag{9.44}$$

*Proof of Lemma 9.1.16.* Throughout this proof assume w.l.o.g. that $a < b$, let $p \in \mathbb{R}$ satisfy $p = \frac{-a}{(b-a)}$, and let $\phi_r\colon \mathbb{R} \to \mathbb{R}$, $r \in [0,1]$, satisfy for all $r \in [0,1]$, $x \in \mathbb{R}$ that

$$\phi_r(x) = \ln(1 - r + re^x) - rx. \tag{9.45}$$

Observe that the assumption that $\mathbb{E}[X] = 0$ and the fact that $a \le \mathbb{E}[X] \le b$ ensures that $a \le 0 \le b$. Combining this with the assumption that $a < b$ implies that

$$0 \le p = \frac{-a}{(b-a)} \le \frac{(b-a)}{(b-a)} = 1. \tag{9.46}$$

Lemma 9.1.14 and Lemma 9.1.15 hence demonstrate that

$$\mathbb{E}\big[e^{\lambda X}\big] \le e^{\phi_p(\lambda(b-a))} = \exp(\phi_p(\lambda(b-a))) \le \exp\Big(\tfrac{(\lambda(b-a))^2}{8}\Big) = \exp\Big(\tfrac{\lambda^2(b-a)^2}{8}\Big). \tag{9.47}$$

The proof of Lemma 9.1.16 is thus complete. $\qquad\square$

### 9.1.5.3   Probability to cross a barrier

**Lemma 9.1.17.** *Let $\beta \in (0, \infty)$, $\varepsilon \in [0, \infty)$ and let $f\colon [0, \infty) \to [0, \infty)$ be the function which satisfies for all $\lambda \in [0, \infty)$ that $f(\lambda) = \beta\lambda^2 - \varepsilon\lambda$. Then*

$$\inf_{\lambda \in [0,\infty)} f(\lambda) = f\big(\tfrac{\varepsilon}{2\beta}\big) = -\tfrac{\varepsilon^2}{4\beta}. \tag{9.48}$$

*Proof of Lemma 9.1.17.* Observe that for all $\lambda \in \mathbb{R}$ it holds that

$$f'(\lambda) = 2\beta\lambda - \varepsilon. \tag{9.49}$$

Moreover, note that

$$f\big(\tfrac{\varepsilon}{2\beta}\big) = \beta\Big[\tfrac{\varepsilon}{2\beta}\Big]^2 - \varepsilon\Big[\tfrac{\varepsilon}{2\beta}\Big] = \tfrac{\varepsilon^2}{4\beta} - \tfrac{\varepsilon^2}{2\beta} = -\tfrac{\varepsilon^2}{4\beta}. \tag{9.50}$$

Combining this and (9.49) establishes (9.48). The proof of Lemma 9.1.17 is thus complete. $\qquad\square$

**Corollary 9.1.18.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $N \in \mathbb{N}$, $\varepsilon \in [0, \infty)$, $a_1, a_2, \ldots, a_N \in \mathbb{R}$, $b_1 \in [a_1, \infty)$, $b_2 \in [a_2, \infty)$, ..., $b_N \in [a_N, \infty)$ satisfy $\sum_{n=1}^N (b_n - a_n)^2 \ne 0$, and let $X_n\colon \Omega \to [a_n, b_n]$, $n \in \{1, 2, \ldots, N\}$, be independent random variables. Then*

$$\mathbb{P}\left(\Big[\sum_{n=1}^N (X_n - \mathbb{E}[X_n])\Big] \ge \varepsilon\right) \le \exp\left(\frac{-2\varepsilon^2}{\sum_{n=1}^N (b_n - a_n)^2}\right). \tag{9.51}$$

*Proof of Corollary 9.1.18.* Throughout this proof let $\beta \in (0, \infty)$ satisfy

$$\beta = \frac{1}{8}\left[\sum_{n=1}^N (b_n - a_n)^2\right]. \tag{9.52}$$

Observe that Corollary 9.1.13 ensures that

$$\mathbb{P}\left(\left[\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right] \geq \varepsilon\right) \leq \inf_{\lambda \in [0,\infty)}\left(e^{-\lambda\varepsilon}\left[\prod_{n=1}^{N}\mathbb{M}_{X_n - \mathbb{E}[X_n]}(\lambda)\right]\right). \tag{9.53}$$

Moreover, note that Lemma 9.1.16 proves that for all $n \in \{1, 2, \ldots, N\}$ it holds that

$$\mathbb{M}_{X_n - \mathbb{E}[X_n]}(\lambda) \leq \exp\left(\tfrac{\lambda^2[(b_n - \mathbb{E}[X_n]) - (a_n - \mathbb{E}[X_n])]^2}{8}\right) = \exp\left(\tfrac{\lambda^2(b_n - a_n)^2}{8}\right). \tag{9.54}$$

Combining this with (9.53) and Lemma 9.1.17 ensures that

$$\mathbb{P}\left(\left[\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right] \geq \varepsilon\right) \leq \inf_{\lambda \in [0,\infty)}\left(\exp\left(\left[\sum_{n=1}^{N}\left(\tfrac{\lambda^2(b_n - a_n)^2}{8}\right)\right] - \lambda\varepsilon\right)\right)$$

$$= \inf_{\lambda \in [0,\infty)}\left[\exp\left(\lambda^2\left[\tfrac{\sum_{n=1}^{N}(b_n - a_n)^2}{8}\right] - \lambda\varepsilon\right)\right] = \exp\left(\inf_{\lambda \in [0,\infty)}\left[\beta\lambda^2 - \varepsilon\lambda\right]\right) \tag{9.55}$$

$$= \exp\left(\tfrac{-\varepsilon^2}{4\beta}\right) = \exp\left(\tfrac{-2\varepsilon^2}{\sum_{n=1}^{N}(b_n - a_n)^2}\right).$$

The proof of Corollary 9.1.18 is thus complete. $\qquad\square$

### 9.1.5.4 Probability to fall below a barrier

**Corollary 9.1.19.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $N \in \mathbb{N}$, $\varepsilon \in [0,\infty)$, $a_1, a_2, \ldots, a_N \in \mathbb{R}$, $b_1 \in [a_1, \infty)$, $b_2 \in [a_2, \infty)$, $\ldots$, $b_N \in [a_N, \infty)$ satisfy $\sum_{n=1}^{N}(b_n - a_n)^2 \neq 0$, and let $X_n \colon \Omega \to [a_n, b_n]$, $n \in \{1, 2, \ldots, N\}$, be independent random variables. Then*

$$\mathbb{P}\left(\left[\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right] \leq -\varepsilon\right) \leq \exp\left(\tfrac{-2\varepsilon^2}{\sum_{n=1}^{N}(b_n - a_n)^2}\right). \tag{9.56}$$

*Proof of Corollary 9.1.19.* Throughout this proof let $\mathfrak{X}_n \colon \Omega \to [-b_n, -a_n]$, $n \in \{1, 2, \ldots, N\}$, satisfy for all $n \in \{1, 2, \ldots, N\}$ that

$$\mathfrak{X}_n = -X_n. \tag{9.57}$$

Observe that Corollary 9.1.18 and (9.57) ensure that

$$\mathbb{P}\left(\left[\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right] \leq -\varepsilon\right) = \mathbb{P}\left(\left[\sum_{n=1}^{N}(-X_n - \mathbb{E}[-X_n])\right] \geq \varepsilon\right)$$

$$= \mathbb{P}\left(\left[\sum_{n=1}^{N}(\mathfrak{X}_n - \mathbb{E}[\mathfrak{X}_n])\right] \geq \varepsilon\right) \leq \exp\left(\tfrac{-2\varepsilon^2}{\sum_{n=1}^{N}(b_n - a_n)^2}\right). \tag{9.58}$$

The proof of Corollary 9.1.19 is thus complete. $\qquad\square$

### 9.1.5.5 Hoeffding's inequality

**Corollary 9.1.20.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $N \in \mathbb{N}$, $\varepsilon \in [0, \infty)$, $a_1, a_2, \ldots, a_N \in \mathbb{R}$, $b_1 \in [a_1, \infty)$, $b_2 \in [a_2, \infty)$, ..., $b_N \in [a_N, \infty)$ satisfy $\sum_{n=1}^{N}(b_n - a_n)^2 \neq 0$, and let $X_n \colon \Omega \to [a_n, b_n]$, $n \in \{1, 2, \ldots, N\}$, be independent random variables. Then*

$$\mathbb{P}\left(\left|\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right| \geq \varepsilon\right) \leq 2\exp\left(\frac{-2\varepsilon^2}{\sum_{n=1}^{N}(b_n - a_n)^2}\right). \tag{9.59}$$

*Proof of Corollary 9.1.20.* Observe that

$$\mathbb{P}\left(\left|\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right| \geq \varepsilon\right)$$
$$= \mathbb{P}\left(\left\{\left[\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right] \geq \varepsilon\right\} \cup \left\{\left[\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right] \leq -\varepsilon\right\}\right) \tag{9.60}$$
$$\leq \mathbb{P}\left(\left[\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right] \geq \varepsilon\right) + \mathbb{P}\left(\left[\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right] \leq -\varepsilon\right).$$

Combining this with Corollary 9.1.18 and Corollary 9.1.19 establishes (9.59). The proof of Corollary 9.1.20 is thus complete. □

**Corollary 9.1.21.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $N \in \mathbb{N}$, $\varepsilon \in [0, \infty)$, $a_1, a_2, \ldots, a_N \in \mathbb{R}$, $b_1 \in [a_1, \infty)$, $b_2 \in [a_2, \infty)$, ..., $b_N \in [a_N, \infty)$ satisfy $\sum_{n=1}^{N}(b_n - a_n)^2 \neq 0$, and let $X_n \colon \Omega \to [a_n, b_n]$, $n \in \{1, 2, \ldots, N\}$, be independent random variables. Then*

$$\mathbb{P}\left(\frac{1}{N}\left|\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right| \geq \varepsilon\right) \leq 2\exp\left(\frac{-2\varepsilon^2 N^2}{\sum_{n=1}^{N}(b_n - a_n)^2}\right). \tag{9.61}$$

*Proof of Corollary 9.1.21.* Observe that Corollary 9.1.20 ensures that

$$\mathbb{P}\left(\frac{1}{N}\left|\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right| \geq \varepsilon\right) = \mathbb{P}\left(\left|\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right| \geq \varepsilon N\right) \leq 2\exp\left(\frac{-2(\varepsilon N)^2}{\sum_{n=1}^{N}(b_n - a_n)^2}\right). \tag{9.62}$$

The proof of Corollary 9.1.21 is thus complete. □

**Exercise 9.1.1.** *Prove or disprove the following statement: For every probability space $(\Omega, \mathcal{F}, \mathbb{P})$, every $N \in \mathbb{N}$, $\varepsilon \in [0, \infty)$, and every random variable $X = (X_1, X_2, \ldots, X_N)\colon \Omega \to [-1, 1]^N$ with $\forall\, a = (a_1, a_2, \ldots, a_N) \in [-1, 1]^N \colon \mathbb{P}(\bigcap_{i=1}^{N}\{X_i \leq a_i\}) = \prod_{i=1}^{N}\frac{a_i+1}{2}$ it holds that*

$$\mathbb{P}\left(\frac{1}{N}\left|\sum_{i=1}^{N}(X_n - \mathbb{E}[X_n])\right| \geq \varepsilon\right) \leq 2\exp\left(\frac{-\varepsilon^2 N}{2}\right). \tag{9.63}$$

**Exercise 9.1.2.** *Prove or disprove the following statement: For every probability space $(\Omega, \mathcal{F}, \mathbb{P})$, every $N \in \mathbb{N}$, and every random variable $X = (X_1, X_2, \ldots, X_N)\colon \Omega \to [-1, 1]^N$ with $\forall\, a = (a_1, a_2, \ldots, a_N) \in [-1, 1]^N \colon \mathbb{P}(\bigcap_{i=1}^{N}\{X_i \leq a_i\}) = \prod_{i=1}^{N}\frac{a_i+1}{2}$ it holds that*

$$\mathbb{P}\left(\frac{1}{N}\left|\sum_{n=1}^{N}(X_n - \mathbb{E}[X_n])\right| \geq \frac{1}{2}\right) \leq 2\left[\frac{e}{4}\right]^N. \tag{9.64}$$

**Exercise 9.1.3.** *Prove or disprove the following statement: For every probability space $(\Omega, \mathcal{F}, \mathbb{P})$, every $N \in \mathbb{N}$, and every random variable $X = (X_1, X_2, \ldots, X_N) \colon \Omega \to [-1, 1]^N$ with $\forall\, a = (a_1, a_2, \ldots, a_N) \in [-1, 1]^N \colon \mathbb{P}(\bigcap_{i=1}^{N} \{X_i \leq a_i\}) = \prod_{i=1}^{N} \frac{a_i + 1}{2}$ it holds that*

$$\mathbb{P}\left( \frac{1}{N} \left| \sum_{n=1}^{N} (X_n - \mathbb{E}[X_n]) \right| \geq \frac{1}{2} \right) \leq 2 \left[ \frac{e - e^{-3}}{4} \right]^N. \tag{9.65}$$

**Exercise 9.1.4.** *Prove or disprove the following statement: For every probability space $(\Omega, \mathcal{F}, \mathbb{P})$, every $N \in \mathbb{N}$, $\varepsilon \in [0, \infty)$, and every standard normal random variable $X = (X_1, X_2, \ldots, X_N) \colon \Omega \to \mathbb{R}^N$ it holds that*

$$\mathbb{P}\left( \frac{1}{N} \left| \sum_{n=1}^{N} (X_n - \mathbb{E}[X_n]) \right| \geq \varepsilon \right) \leq 2 \exp\left( \frac{-\varepsilon^2 N}{2} \right). \tag{9.66}$$

### 9.1.6   A strengthened Hoeffding's inequality

**Lemma 9.1.22.** *Let $f, g \colon (0, \infty) \to \mathbb{R}$ satisfy for all $x \in (0, \infty)$ that $f(x) = 2 \exp(-2x)$ and $g(x) = \frac{1}{4x}$. Then*

(i) *it holds that $\lim_{x \to \infty} \frac{f(x)}{g(x)} = \lim_{x \searrow 0} \frac{f(x)}{g(x)} = 0$ and*

(ii) *it holds that $g(\frac{1}{2}) = \frac{1}{2} < \frac{2}{3} < \frac{2}{e} = f(\frac{1}{2})$.*

*Proof of Lemma 9.1.22.* Note that the fact that $\lim_{x \to \infty} \frac{\exp(-x)}{x^{-1}} = \lim_{x \searrow 0} \frac{\exp(-x)}{x^{-1}} = 0$ establishes item (i). Moreover, observe that the fact that $e < 3$ implies item (ii). The proof of Lemma 9.1.22 is thus complete. $\qquad \square$

**Corollary 9.1.23.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $N \in \mathbb{N}$, $\varepsilon \in (0, \infty)$, $a_1, a_2, \ldots, a_N \in \mathbb{R}$, $b_1 \in [a_1, \infty)$, $b_2 \in [a_2, \infty)$, $\ldots$, $b_N \in [a_N, \infty)$ satisfy $\sum_{n=1}^{N} (b_n - a_n)^2 \neq 0$, and let $X_n \colon \Omega \to [a_n, b_n]$, $n \in \{1, 2, \ldots, N\}$, be independent random variables. Then*

$$\mathbb{P}\left( \left| \sum_{n=1}^{N} (X_n - \mathbb{E}[X_n]) \right| \geq \varepsilon \right) \leq \min\left\{ 1, 2 \exp\left( \frac{-2\varepsilon^2}{\sum_{n=1}^{N} (b_n - a_n)^2} \right), \frac{\sum_{n=1}^{N} (b_n - a_n)^2}{4\varepsilon^2} \right\}. \tag{9.67}$$

*Proof of Corollary 9.1.23.* Observe that Lemma 9.1.6, Corollary 9.1.20, and the fact that for all $B \in \mathcal{F}$ it holds that $\mathbb{P}(B) \leq 1$ establish (9.67). The proof of Corollary 9.1.23 is thus complete. $\qquad \square$

## 9.2   Covering number estimates

This section is inspired by Section 6 in Chapter I in Cucker & Smale [6] and Section 1.1 in Carl & Stephani [4].

## 9.2.1 Entropy quantities

### 9.2.1.1 Covering radii (Outer entropy numbers)

**Definition 9.2.1.** *Let $(X, d)$ be a metric space and let $n \in \mathbb{N}$. Then we denote by $\mathcal{C}_{(X,d),n} \in [0, \infty]$ (we denote by $\mathcal{C}_{X,n} \in [0, \infty]$) the extended real number given by*

$$\mathcal{C}_{(X,d),n} = \inf\left(\left\{ r \in [0, \infty] \colon \left(\exists\, A \subseteq X \colon \left[(|A| \leq n) \wedge (\forall\, x \in X \colon \exists\, a \in A \colon d(a, x) \leq r)\right]\right)\right\}\right) \tag{9.68}$$

*and we call $\mathcal{C}_{(X,d),n}$ the n-covering radius of $(X, d)$ (we call $\mathcal{C}_{X,r}$ the n-covering radius of $X$).*

**Lemma 9.2.2.** *Let $(X, d)$ be a metric space, let $n \in \mathbb{N}$, $r \in [0, \infty]$, assume $X \neq \emptyset$, and let $A \subseteq X$ satisfy $|A| \leq n$ and $\forall\, x \in X \colon \exists\, a \in A \colon d(a, x) \leq r$. Then there exist $x_1, x_2, \ldots, x_n \in X$ such that*

$$X \subseteq \left[\bigcup_{i=1}^{n} \{v \in X \colon d(x_i, v) \leq r\}\right]. \tag{9.69}$$

*Proof of Lemma 9.2.2.* Note that the assumption that $X \neq \emptyset$ and the assumption that $|A| \leq n$ imply that there exist $x_1, x_2, \ldots, x_n \in X$ which satisfy $A \subseteq \{x_1, x_2, \ldots, x_n\}$. This and the assumption that $\forall\, x \in X \colon \exists\, a \in A \colon d(a, x) \leq r$ ensure that

$$X \subseteq \left[\bigcup_{a \in A} \{v \in X \colon d(a, v) \leq r\}\right] \subseteq \left[\bigcup_{i=1}^{n} \{v \in X \colon d(x_i, v) \leq r\}\right]. \tag{9.70}$$

The proof of Lemma 9.2.2 is thus complete. $\qquad\square$

**Lemma 9.2.3.** *Let $(X, d)$ be a metric space, let $n \in \mathbb{N}$, $r \in [0, \infty]$, $x_1, x_2, \ldots, x_n \in X$ satisfy $X \subseteq \left[\bigcup_{i=1}^{n} \{v \in X \colon d(x_i, v) \leq r\}\right]$. Then there exists $A \subseteq X$ such that $|A| \leq n$ and*

$$\forall\, x \in X \colon \exists\, a \in A \colon d(a, x) \leq r. \tag{9.71}$$

*Proof of Lemma 9.2.3.* Throughout this proof let $A = \{x_1, x_2, \ldots, x_n\}$. Note that the assumption that $X \subseteq \left[\bigcup_{i=1}^{n} \{v \in X \colon d(x_i, v) \leq r\}\right]$ implies that for all $v \in X$ there exists $i \in \{1, 2, \ldots, n\}$ such that $d(x_i, v) \leq r$. Hence, we obtain that

$$\forall\, x \in X \colon \exists\, a \in A \colon d(a, x) \leq r. \tag{9.72}$$

The proof of Lemma 9.2.3 is thus complete. $\qquad\square$

**Lemma 9.2.4.** *Let $(X, d)$ be a metric space, let $n \in \mathbb{N}$, $r \in [0, \infty]$, and assume $X \neq \emptyset$. Then the following two statements are equivalent:*

(i) *There exists $A \subseteq X$ such that $|A| \leq n$ and $\forall\, x \in X \colon \exists\, a \in A \colon d(a, x) \leq r$.*

(ii) *There exist $x_1, x_2, \ldots, x_n \in X$ such that $X \subseteq \left[\bigcup_{i=1}^{n} \{v \in X \colon d(x_i, v) \leq r\}\right]$.*

*Proof of Lemma 9.2.4.* Note that Lemma 9.2.2 and Lemma 9.2.3 prove that ((i)⇔(ii)). The proof of Lemma 9.2.4 is thus complete. $\qquad\square$

**Lemma 9.2.5.** *Let $(X, d)$ be a metric space and let $n \in \mathbb{N}$. Then*

$$
\mathcal{C}_{(X,d),n} = \begin{cases} 0 & : X = \emptyset \\ \inf\left(\left\{r \in [0,\infty) \colon \left(\exists\, x_1, x_2, \ldots, x_n \in X \colon \right.\right.\right. \\ \qquad\qquad \left.\left.\left. X \subseteq \left[\bigcup_{m=1}^{n} \{v \in X \colon d(x_m, v) \le r\}\right]\right)\right\} \cup \{\infty\}\right) & : X \ne \emptyset. \end{cases}
$$

$$(9.73)$$

*Proof of Lemma 9.2.5.* Throughout this proof assume w.l.o.g. that $X \ne \emptyset$ and let $a \in X$. Note that the assumption that $d$ is a metric implies that for all $x \in X$ it holds that $d(a, x) \le \infty$. Combining this with Lemma 9.2.4 proves (9.73). This completes the proof of Lemma 9.2.5. $\qquad\square$

**Exercise 9.2.1.** *Prove or disprove the following statement: For every metric space $(X, d)$ and every $n, m \in \mathbb{N}$ it holds that $\mathcal{C}_{(X,d),n} < \infty$ if and only if $\mathcal{C}_{(X,d),m} < \infty$.*

**Exercise 9.2.2.** *Prove or disprove the following statement: For every metric space $(X, d)$ and every $n \in \mathbb{N}$ it holds that $(X, d)$ is bounded if and only if $\mathcal{C}_{(X,d),n} < \infty$.*

**Exercise 9.2.3.** *Prove or disprove the following statement: For every $n \in \mathbb{N}$ and every metric space $(X, d)$ with $X \ne \emptyset$ it holds that*

$$
\mathcal{C}_{(X,d),n} = \inf_{x_1, x_2, \ldots, x_n \in X} \sup_{v \in X} \min_{i \in \{1, 2, \ldots, n\}} d(x_i, v) = \inf_{x_1, x_2, \ldots, x_n \in X} \sup_{x_{n+1} \in X} \min_{i \in \{1, 2, \ldots, n\}} d(x_i, x_{n+1})
$$

$$(9.74)$$

### 9.2.1.2 Covering numbers

**Definition 9.2.6.** *Let $(X, d)$ be a metric space and let $r \in [0, \infty]$. Then we denote by $\mathcal{C}^{(X,d),r} \in [0, \infty]$ (we denote by $\mathcal{C}^{X,r} \in [0, \infty]$) the extended real number given by*

$$
\mathcal{C}^{(X,d),r} = \inf\left(\left\{n \in \mathbb{N}_0 \colon \left(\exists\, A \subseteq X \colon \left[(|A| \le n) \wedge (\forall\, x \in X \colon \exists\, a \in A \colon d(a, x) \le r)\right]\right)\right\} \cup \{\infty\}\right)
$$

$$(9.75)$$

*and we call $\mathcal{C}^{(X,d),r}$ the $r$-covering number of $(X, d)$ (we call $\mathcal{C}^{X,r}$ the $r$-covering number of $X$).*

**Lemma 9.2.7.** *Let $(X, d)$ be a metric space and let $r \in [0, \infty]$. Then*

$$
\mathcal{C}^{(X,d),r} = \begin{cases} 0 & : X = \emptyset \\ \inf\left(\left\{n \in \mathbb{N} \colon \left(\exists\, x_1, x_2, \ldots, x_n \in X \colon \right.\right.\right. \\ \qquad\qquad \left.\left.\left. X \subseteq \left[\bigcup_{m=1}^{n} \{v \in X \colon d(x_m, v) \le r\}\right]\right)\right\} \cup \{\infty\}\right) & : X \ne \emptyset. \end{cases}
$$

$$(9.76)$$

*Proof of Lemma 9.2.7.* Throughout this proof assume w.l.o.g. that $X \ne \emptyset$. Observe that Lemma 9.2.4 establishes (9.76). The proof of Lemma 9.2.7 is thus complete. $\qquad\square$

**Exercise 9.2.4.** *Prove or disprove the following statement: For every $r \in [0, \infty]$, every metric space $(X, d)$, and every $Y \subseteq X$ it holds that*

$$
\mathcal{C}^{(Y, d|_{Y \times Y}), r} \le \mathcal{C}^{(X,d),r}.
$$

$$(9.77)$$

### 9.2.2.3 Upper bounds for packing radii based on upper bounds for covering radii

**Lemma 9.2.12.** *Let $(X, d)$ be a metric space and let $n \in \mathbb{N}$. Then $\mathcal{P}_{(X,d),n} \leq \mathcal{C}_{(X,d),n}$.*

*Proof of Lemma 9.2.12.* Throughout this proof assume w.l.o.g. that $\mathcal{C}_{(X,d),n} < \infty$ and $\mathcal{P}_{(X,d),n} > 0$, let $r \in [0, \infty)$, $x_1, x_2, \ldots, x_n \in X$ satisfy

$$X \subseteq \left[ \bigcup_{m=1}^{n} \{v \in X : d(x_m, v) \leq r\} \right], \tag{9.82}$$

let $\mathbf{r} \in [0, \infty)$, $\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_{n+1} \in X$ satisfy

$$\left[ \min_{i,j \in \{1,2,\ldots,n+1\}, i \neq j} d(\mathbf{x}_i, \mathbf{x}_j) \right] > 2\mathbf{r}, \tag{9.83}$$

and let $\varphi \colon X \to \{1, 2, \ldots, n\}$ satisfy for all $v \in X$ that

$$\varphi(v) = \min\{m \in \{1, 2, \ldots, n\} : v \in \{w \in X : d(x_m, w) \leq r\}\} \tag{9.84}$$

(cf. Lemma 9.2.5). Observe that (9.84) shows that for all $v \in X$ it holds that

$$v \in \{w \in X : d(x_{\varphi(v)}, w) \leq r\}. \tag{9.85}$$

Hence, we obtain that for all $v \in X$ it holds that

$$d(v, x_{\varphi(v)}) \leq r \tag{9.86}$$

Moreover, note that the fact that $\varphi(\mathbf{x}_1), \varphi(\mathbf{x}_2), \ldots, \varphi(\mathbf{x}_{n+1}) \in \{1, 2, \ldots, n\}$ ensures that there exist $i, j \in \{1, 2, \ldots, n+1\}$ which satisfy

$$i \neq j \qquad \text{and} \qquad \varphi(\mathbf{x}_i) = \varphi(\mathbf{x}_j). \tag{9.87}$$

The triangle inequality, (9.83), and (9.86) hence show that

$$2\mathbf{r} < d(\mathbf{x}_i, \mathbf{x}_j) \leq d(\mathbf{x}_i, x_{\varphi(\mathbf{x}_i)}) + d(x_{\varphi(\mathbf{x}_i)}, \mathbf{x}_j) = d(\mathbf{x}_i, x_{\varphi(\mathbf{x}_i)}) + d(\mathbf{x}_j, x_{\varphi(\mathbf{x}_j)}) \leq 2r. \tag{9.88}$$

This implies that $\mathbf{r} < r$. The proof of Lemma 9.2.12 is thus complete. $\qquad \square$

### 9.2.2.4 Upper bounds for packing radii in balls of metric spaces

**Lemma 9.2.13.** *Let $(X, d)$ be a metric space, let $n \in \mathbb{N}$, $x \in X$, $r \in (0, \infty]$, and let $S = \{v \in X : d(x, v) \leq r\}$. Then $\mathcal{P}_{(S,d|_{S \times S}),n} \leq r$.*

*Proof of Lemma 9.2.13.* Throughout this proof assume w.l.o.g. that $\mathcal{P}_{(S,d|_{S \times S}),n} > 0$. Observe that for all $\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_{n+1} \in S$, $i, j \in \{1, 2, \ldots, n+1\}$ it holds that

$$d(\mathbf{x}_i, \mathbf{x}_j) \leq d(\mathbf{x}_i, x) + d(x, \mathbf{x}_j) \leq 2r. \tag{9.89}$$

Hence, we obtain that for all $\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_{n+1} \in S$ it holds that

$$\min_{i,j \in \{1,2,\ldots,n+1\}, i \neq j} d(\mathbf{x}_i, \mathbf{x}_j) \leq 2r. \tag{9.90}$$

Moreover, note that (9.78) ensures that for all $\rho \in [0, \mathcal{P}_{(S,d|_{S \times S}),n})$ there exist $\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_{n+1} \in S$ such that

$$\min_{i,j \in \{1,2,\ldots,n+1\}, i \neq j} d(\mathbf{x}_i, \mathbf{x}_j) > 2\rho. \tag{9.91}$$

This and (9.90) demonstrate that for all $\rho \in [0, \mathcal{P}_{(S,d|_{S \times S}),n})$ it holds that $2\rho < 2r$. The proof of Lemma 9.2.13 is thus complete. $\qquad \square$

### 9.2.3   Inequalities for covering entropy quantities in metric spaces

#### 9.2.3.1   Upper bounds for covering numbers based on upper bounds for covering radii

**Lemma 9.2.14.** *Let $(X, d)$ be a metric space and let $r \in [0, \infty]$, $n \in \mathbb{N}$ satisfy $\mathcal{C}_{(X,d),n} < r$. Then $\mathcal{C}^{(X,d),r} \leq n$.*

*Proof of Lemma 9.2.14.* Observe that the assumption that $\mathcal{C}_{(X,d),n} < r$ ensures that there exists $A \subseteq X$ such that $|A| \leq n$ and

$$X \subseteq \left[ \bigcup_{a \in A} \{v \in X : d(a, v) \leq r\} \right]. \tag{9.92}$$

This establishes that $\mathcal{C}^{(X,d),r} \leq n$. The proof of Lemma 9.2.14 is thus complete. $\square$

**Lemma 9.2.15.** *Let $(X, d)$ be a compact metric space and let $r \in [0, \infty]$, $n \in \mathbb{N}$, satisfy $\mathcal{C}_{(X,d),n} \leq r$. Then $\mathcal{C}^{(X,d),r} \leq n$.*

*Proof of Lemma 9.2.15.* Throughout this proof assume w.l.o.g. that $X \neq \emptyset$ and let $x_{k,m} \in X$, $m \in \{1, 2, \ldots, n\}$, $k \in \mathbb{N}$, satisfy for all $k \in \mathbb{N}$ that

$$X \subseteq \left[ \bigcup_{m=1}^{n} \{v \in X : d(x_{k,m}, v) \leq r + \tfrac{1}{k}\} \right] \tag{9.93}$$

(cf. Lemma 9.2.4). Note that the assumption that $(X, d)$ is a compact metric space demonstrates that there exist $\mathfrak{x} = (\mathfrak{x}_m)_{m \in \{1,2,\ldots,n\}} \colon \{1, 2, \ldots, n\} \to X$ and $k = (k_l)_{l \in \mathbb{N}} \colon \mathbb{N} \to \mathbb{N}$ which satisfy that

$$\limsup_{l \to \infty} \max_{m \in \{1,2,\ldots,n\}} d(\mathfrak{x}_m, x_{k_l,m}) = 0 \qquad \text{and} \qquad \limsup_{l \to \infty} k_l = \infty. \tag{9.94}$$

Next observe that the assumption that $d$ is a metric ensures that for all $v \in X$, $m \in \{1, 2, \ldots, n\}$, $l \in \mathbb{N}$ it holds that

$$d(v, \mathfrak{x}_m) \leq d(v, x_{k_l,m}) + d(x_{k_l,m}, \mathfrak{x}_m). \tag{9.95}$$

This and (9.93) prove that for all $v \in X$, $l \in \mathbb{N}$ it holds that

$$
\begin{aligned}
\min_{m \in \{1,2,\ldots,n\}} d(v, \mathfrak{x}_m) &\leq \min_{m \in \{1,2,\ldots,n\}} [d(v, x_{k_l,m}) + d(x_{k_l,m}, \mathfrak{x}_m)] \\
&\leq \left[ \min_{m \in \{1,2,\ldots,n\}} d(v, x_{k_l,m}) \right] + \left[ \max_{m \in \{1,2,\ldots,n\}} d(x_{k_l,m}, \mathfrak{x}_m) \right] \\
&\leq \left[ r + \tfrac{1}{k_l} \right] + \left[ \max_{m \in \{1,2,\ldots,n\}} d(x_{k_l,m}, \mathfrak{x}_m) \right].
\end{aligned}
\tag{9.96}
$$

Hence, we obtain for all $v \in X$ that

$$\min_{m \in \{1,2,\ldots,n\}} d(v, \mathfrak{x}_m) \leq \limsup_{l \to \infty} \left( \left[ r + \tfrac{1}{k_l} \right] + \left[ \max_{m \in \{1,2,\ldots,n\}} d(x_{k_l,m}, \mathfrak{x}_m) \right] \right) = r. \tag{9.97}$$

This establishes that $\mathcal{C}^{(X,d),r} \leq n$. The proof of Lemma 9.2.15 is thus complete. $\square$

### 9.2.3.2   Upper bounds for covering radii based on upper bounds for covering numbers

**Lemma 9.2.16.** *Let $(X, d)$ be a metric space and let $r \in [0, \infty]$, $n \in \mathbb{N}$ satisfy $\mathcal{C}^{(X,d),r} \leq n$. Then $\mathcal{C}_{(X,d),n} \leq r$.*

*Proof of Lemma 9.2.16.* Observe that the assumption that $\mathcal{C}^{(X,d),r} \leq n$ ensures that there exists $A \subseteq X$ such that $|A| \leq n$ and

$$X \subseteq \left[ \bigcup_{a \in A} \{v \in X \colon d(a, v) \leq r\} \right]. \tag{9.98}$$

This establishes that $\mathcal{C}_{(X,d),n} \leq r$. The proof of Lemma 9.2.16 is thus complete. □

### 9.2.3.3   Upper bounds for covering radii based on upper bounds for packing radii

**Lemma 9.2.17.** *Let $(X, d)$ be a metric space and let $n \in \mathbb{N}$. Then $\mathcal{C}_{(X,d),n} \leq 2\mathcal{P}_{(X,d),n}$.*

*Proof of Lemma 9.2.17.* Throughout this proof assume w.l.o.g. that $X \neq \emptyset$, assume w.l.o.g. that $\mathcal{P}_{(X,d),n} < \infty$, let $r \in [0, \infty]$ satisfy $r > \mathcal{P}_{(X,d),n}$, and let $N \in \mathbb{N}_0 \cup \{\infty\}$ satisfy $N = \mathcal{P}^{(X,d),r}$. Observe that Lemma 9.2.11 ensures that

$$N = \mathcal{P}^{(X,d),r} < n. \tag{9.99}$$

Moreover, note that the fact that $N = \mathcal{P}^{(X,d),r}$ and (9.80) demonstrate that for all $x_1, x_2, \ldots, x_{N+1}, x_{N+2} \in X$ it holds that

$$\min_{i,j \in \{1,2,\ldots,N+2\}, i \neq j} d(x_i, x_j) \leq 2r. \tag{9.100}$$

In addition, observe that the fact that $N = \mathcal{P}^{(X,d),r}$ and (9.80) imply that there exist $x_1, x_2, \ldots, x_{N+1} \in X$ which satisfy that

$$\min\big(\{d(x_i, x_j) \colon i, j \in \{1, 2, \ldots, N+1\}, \, i \neq j\} \cup \{\infty\}\big) > 2r. \tag{9.101}$$

Combining this with (9.100) establishes that for all $v \in X$ it holds that

$$\min_{i \in \{1,2,\ldots,N\}} d(x_i, v) \leq 2r. \tag{9.102}$$

Hence, we obtain that for all $w \in X$ it holds that

$$w \in \left[ \bigcup_{m=1}^{n} \{v \in X \colon d(x_i, v) \leq 2r\} \right]. \tag{9.103}$$

Therefore, we obtain that

$$X \subseteq \left[ \bigcup_{m=1}^{n} \{v \in X \colon d(x_i, v) \leq 2r\} \right]. \tag{9.104}$$

Combining this and Lemma 9.2.5 shows that $\mathcal{C}_{(X,d),n} \leq 2r$. The proof of Lemma 9.2.17 is thus complete. □

#### 9.2.3.4   Equivalence of covering and packing radii

**Corollary 9.2.18.** *Let $(X, d)$ be a metric space and let $n \in \mathbb{N}$. Then $\mathcal{P}_{(X,d),n} \leq \mathcal{C}_{(X,d),n} \leq 2\mathcal{P}_{(X,d),n}$.*

*Proof of Corollary 9.2.18.* Observe that Lemma 9.2.12 and Lemma 9.2.17 establish that $\mathcal{P}_{(X,d),n} \leq \mathcal{C}_{(X,d),n} \leq 2\mathcal{P}_{(X,d),n}$. The proof of Corollary 9.2.18 is thus complete. □

### 9.2.4   Inequalities for entropy quantities in finite dimensional vector spaces

#### 9.2.4.1   Measures induced by Lebesgue-Borel measures

**Lemma 9.2.19.** *Let $(V, \|\!|\cdot|\!\|)$ be a normed vector space, let $N \in \mathbb{N}$, let $b_1, b_2, \ldots, b_N \in V$ be a Hamel-basis of $V$, let $\lambda \colon \mathcal{B}(\mathbb{R}^N) \to [0, \infty]$ be the Lebesgue-Borel measure on $\mathbb{R}^N$, let $\Phi \colon \mathbb{R}^N \to V$ satisfy for all $r = (r_1, r_2, \ldots, r_N) \in \mathbb{R}^N$ that $\Phi(r) = r_1 b_1 + r_2 b_2 + \ldots + r_N b_N$, and let $\nu \colon \mathcal{B}(V) \to [0, \infty]$ satisfy for all $A \in \mathcal{B}(V)$ that*

$$\nu(A) = \lambda(\Phi^{-1}(A)). \tag{9.105}$$

*Then*

*(i) it holds that $\Phi$ is linear,*

*(ii) it holds for all $r = (r_1, r_2, \ldots, r_N) \in \mathbb{R}^N$ that $\|\!|\Phi(r)|\!\| \leq \left[\sum_{n=1}^N \|\!|b_n|\!\|^2\right]^{1/2}\left[\sum_{n=1}^N |r_n|^2\right]^{1/2}$,*

*(iii) it holds that $\Phi \in C(\mathbb{R}^N, V)$,*

*(iv) it holds that $\Phi$ is bijective,*

*(v) it holds that $(V, \mathcal{B}(V), \nu)$ is a measure space,*

*(vi) it holds for all $r \in (0, \infty)$, $v \in V$, $A \in \mathcal{B}(V)$ that $\nu(\{(ra + v) \in V \colon a \in A\}) = r^N \nu(A)$,*

*(vii) it holds for all $r \in (0, \infty)$ that $\nu(\{v \in V \colon \|\!|v|\!\| \leq r\}) = r^N \nu(\{v \in V \colon \|\!|v|\!\| \leq 1\})$, and*

*(viii) it holds that $\nu(\{v \in V \colon \|\!|v|\!\| \leq 1\}) > 0$.*

*Proof of Lemma 9.2.19.* Note that for all $r = (r_1, r_2, \ldots, r_N)$, $s = (s_1, s_2, \ldots, s_N) \in \mathbb{R}^N$, $\rho \in \mathbb{R}$ it holds that

$$\Phi(\rho r + s) = (\rho r_1 + s_1)b_1 + (\rho r_2 + s_2)b_2 + \cdots + (\rho r_N + s_N)b_N = \rho\Phi(r) + \Phi(s). \tag{9.106}$$

This establishes item (i). Next observe that Hölder's inequality shows that for all $r = (r_1, r_2, \ldots, r_N) \in \mathbb{R}^N$ it holds that

$$\|\!|\Phi(r)|\!\| = \|\!|r_1 b_1 + r_2 b_2 + \cdots + r_N b_N|\!\| \leq \sum_{n=1}^N |r_n|\|\!|b_n|\!\| \leq \left[\sum_{n=1}^N \|\!|b_n|\!\|^2\right]^{1/2}\left[\sum_{n=1}^N |r_n|^2\right]^{1/2}. \tag{9.107}$$

This establishes item (ii). Moreover, note that item (ii) proves item (iii). Furthermore, observe that the assumption that $b_1, b_2, \ldots, b_N \in V$ is a Hamel-basis of $V$ establishes

item (iv). Next note that (9.105) and item (iii) prove item (v). In addition, observe that the integral transformation theorem shows that for all $r \in (0, \infty)$, $v \in \mathbb{R}^N$, $A \in \mathcal{B}(\mathbb{R}^N)$ it holds that

$$
\lambda(\{(ra + v) \in \mathbb{R}^N : a \in A\}) = \lambda(\{ra \in \mathbb{R}^N : a \in A\}) = \int_{\mathbb{R}^N} \mathbb{1}_{\{ra \in \mathbb{R}^N : a \in A\}}(x) \, dx
$$
$$
= \int_{\mathbb{R}^N} \mathbb{1}_A(\tfrac{x}{r}) \, dx = r^N \int_{\mathbb{R}^N} \mathbb{1}_A(x) \, dx = r^N \lambda(A).
$$
(9.108)

Combining item (i) and item (iv) hence demonstrates that for all $r \in (0, \infty)$, $v \in V$, $A \in \mathcal{B}(V)$ it holds that

$$
\nu(\{(ra + v) \in V : a \in A\}) = \lambda\big(\Phi^{-1}(\{(ra + v) \in V : a \in A\})\big) = \lambda\big(\{\Phi^{-1}(ra + v) \in \mathbb{R}^N : a \in A\}\big)
$$
$$
= \lambda\big(\{[r\Phi^{-1}(a) + \Phi^{-1}(v)] \in \mathbb{R}^N : a \in A\}\big)
$$
$$
= \lambda\big(\{[ra + \Phi^{-1}(v)] \in \mathbb{R}^N : a \in \Phi^{-1}(A)\}\big) = r^N \lambda(\Phi^{-1}(A)) = r^N \nu(A).
$$
(9.109)

This establishes item (vi). Hence, we obtain that for all $r \in (0, \infty)$ it holds that

$$
\nu(\{v \in V : \|v\| \le r\}) = \nu(\{rv \in V : \|v\| \le 1\}) = r^N \nu(\{v \in V : \|v\| \le 1\}) = r^N \nu(X).
$$
(9.110)

This establishes item (vii). Furthermore, observe that (9.110) demonstrates that

$$
\infty = \lambda(\mathbb{R}^N) = \nu(V) = \limsup_{r \to \infty}\big[\nu(\{v \in V : \|v\| \le r\})\big] = \limsup_{r \to \infty}\big[r^N \nu(\{v \in V : \|v\| \le 1\})\big].
$$
(9.111)

Hence, we obtain that $\nu(\{v \in V : \|v\| \le 1\}) \ne 0$. This establishes item (viii). The proof of Lemma 9.2.19 is thus complete. $\qquad\square$

### 9.2.4.2 Upper bounds for packing radii

**Lemma 9.2.20.** *Let $(V, \|\cdot\|)$ be a normed vector space, let $X = \{v \in V : \|v\| \le 1\}$, let $d\colon X \times X \to [0, \infty)$ satisfy for all $v, w \in X$ that $d(v, w) = \|v - w\|$, and let $n, N \in \mathbb{N}$ satisfy $N = \dim(V)$. Then*

$$
\mathcal{P}_{(X,d),n} \le 2\,(n + 1)^{-1/N}.
$$
(9.112)

*Proof of Lemma 9.2.20.* Throughout this proof assume w.l.o.g. that $\mathcal{P}_{(X,d),n} > 0$, let $\rho \in [0, \mathcal{P}_{(X,d),n})$, let $\lambda\colon \mathcal{B}(\mathbb{R}^N) \to [0, \infty]$ be the Lebesgue-Borel measure on $\mathbb{R}^N$, let $b_1, b_2, \ldots, b_N \in V$ be a Hamel-basis of $V$, let $\Phi\colon \mathbb{R}^N \to V$ satisfy for all $r = (r_1, r_2, \ldots, r_N) \in \mathbb{R}^N$ that

$$
\Phi(r) = r_1 b_1 + r_2 b_2 + \ldots + r_N b_N,
$$
(9.113)

and let $\nu\colon \mathcal{B}(V) \to [0, \infty]$ satisfy for all $A \in \mathcal{B}(V)$ that

$$
\nu(A) = \lambda(\Phi^{-1}(A)).
$$
(9.114)

Observe that Lemma 9.2.13 ensures that $\rho < \mathcal{P}_{(X,d),n} \le 1$. Moreover, note that (9.78) shows that there exist $x_1, x_2, \ldots, x_{n+1} \in X$ which satisfy

$$
\min_{i,j \in \{1,2,\ldots,n+1\}, i \ne j} \|x_i - x_j\| = \min_{i,j \in \{1,2,\ldots,n+1\}, i \ne j} d(x_i, x_j) > 2\rho.
$$
(9.115)

Observe that (9.115) ensures that for all $i, j \in \{1, 2, \ldots, n+1\}$ with $i \neq j$ it holds that

$$\{v \in V \colon \|\!|x_i - v|\!\| \leq \rho\} \cap \{v \in V \colon \|\!|x_j - v|\!\| \leq \rho\} = \emptyset. \tag{9.116}$$

Moreover, note that (9.115) and the fact that $\rho < 1$ show that for all $j \in \{1, 2, \ldots, n+1\}$, $w \in \{v \in X \colon d(x_j, v) \leq \rho\}$ it holds that

$$\|\!|w|\!\| \leq \|\!|w - x_j|\!\| + \|\!|x_j|\!\| \leq \rho + 1 \leq 2. \tag{9.117}$$

Therefore, we obtain that for all $j \in \{1, 2, \ldots, n+1\}$ it holds that

$$\{v \in V \colon \|\!|v - x_j|\!\| \leq \rho\} \subseteq \{v \in V \colon \|\!|v|\!\| \leq 2\}. \tag{9.118}$$

Next observe that Lemma 9.2.19 ensures that $(V, \mathcal{B}(V), \nu)$ is a measure space. Combining this and (9.116) with (9.118) proves that

$$\sum_{j=1}^{n+1} \nu(\{v \in V \colon \|\!|v - x_j|\!\| \leq \rho\}) = \nu\!\left(\bigcup_{j=1}^{n+1}\{v \in V \colon \|\!|v - x_j|\!\| \leq \rho\}\right) \leq \nu(\{v \in V \colon \|\!|v|\!\| \leq 2\}). \tag{9.119}$$

Lemma 9.2.19 hence shows that

$$(n+1)\rho^N \nu(X) = \sum_{j=1}^{n+1}\left[\rho^N \nu(\{v \in V \colon \|\!|v|\!\| \leq 1\})\right] = \sum_{j=1}^{n+1} \nu(\{v \in V \colon \|\!|v|\!\| \leq \rho\})$$

$$= \sum_{j=1}^{n+1} \nu(\{v \in V \colon \|\!|v - x_j|\!\| \leq \rho\}) \leq \nu(\{v \in V \colon \|\!|v|\!\| \leq 2\}) = 2^N \nu(\{v \in V \colon \|\!|v|\!\| \leq 1\}) = 2^N \nu(X). \tag{9.120}$$

Next observe that Lemma 9.2.19 demonstrates that $\nu(X) > 0$. Combining this with (9.120) assures that $(n+1)\rho^N \leq 2^N$. Therefore, we obtain that $\rho^N \leq (n+1)^{-1} 2^N$. Hence, we obtain that $\rho \leq 2(n+1)^{-1/N}$. The proof of Lemma 9.2.20 is thus complete. □

### 9.2.4.3   Upper bounds for covering radii

**Corollary 9.2.21.** *Let $(V, \|\!|\cdot|\!\|)$ be a normed vector space, let $X = \{v \in V \colon \|\!|v|\!\| \leq 1\}$, let $d \colon X \times X \to [0, \infty)$ satisfy for all $v, w \in X$ that $d(v, w) = \|\!|v - w|\!\|$, and let $n, N \in \mathbb{N}$ satisfy $N = \dim(V)$. Then*

$$\mathcal{C}_{(X,d),n} \leq 4\,(n+1)^{-1/N}. \tag{9.121}$$

*Proof of Corollary 9.2.21.* Observe that Corollary 9.2.18 and Lemma 9.2.20 establish (9.121). The proof of Corollary 9.2.21 is thus complete. □

### 9.2.4.4   Lower bounds for covering radii

**Lemma 9.2.22.** *Let $(V, \|\!|\cdot|\!\|)$ be a normed vector space, let $X = \{v \in V \colon \|\!|v|\!\| \leq 1\}$, let $d \colon X \times X \to [0, \infty)$ satisfy for all $v, w \in X$ that $d(v, w) = \|\!|v - w|\!\|$, and let $n, N \in \mathbb{N}$ satisfy $N = \dim(V)$. Then*

$$n^{-1/N} \leq \mathcal{C}_{(X,d),n}. \tag{9.122}$$

*Proof of Lemma 9.2.22.* Throughout this proof assume w.l.o.g. that $\mathcal{C}_{(X,d),n} < \infty$, let $\rho \in (\mathcal{C}_{(X,d),n}, \infty)$, let $\lambda\colon \mathcal{B}(\mathbb{R}^N) \to [0,\infty]$ be the Lebesgue-Borel measure on $\mathbb{R}^N$, let $b_1, b_2, \ldots, b_N \in V$ be a Hamel-basis of $V$, let $\Phi\colon \mathbb{R}^N \to V$ satisfy for all $r = (r_1, r_2, \ldots, r_N) \in \mathbb{R}^N$ that

$$\Phi(r) = r_1 b_1 + r_2 b_2 + \ldots + r_N b_N, \tag{9.123}$$

and let $\nu\colon \mathcal{B}(V) \to [0,\infty]$ satisfy for all $A \in \mathcal{B}(V)$ that

$$\nu(A) = \lambda(\Phi^{-1}(A)). \tag{9.124}$$

The fact that $\rho > \mathcal{C}_{(X,d),n}$ demonstrates that there exist $x_1, x_2, \ldots, x_n \in X$ which satisfy

$$X \subseteq \left[\bigcup_{m=1}^{n} \{v \in X\colon d(x_m, v) \le \rho\}\right]. \tag{9.125}$$

Lemma 9.2.19 hence shows that

$$
\begin{aligned}
\nu(X) &\le \nu\left(\bigcup_{m=1}^{n} \{v \in X\colon d(x_m, v) \le \rho\}\right) \le \sum_{m=1}^{n} \nu(\{v \in X\colon d(x_m, v) \le \rho\}) \\
&= \sum_{m=1}^{n} \left[\rho^N \nu(\{v \in X\colon d(x_m, v) \le 1\})\right] \le n\rho^N \nu(X).
\end{aligned}
\tag{9.126}
$$

This and Lemma 9.2.19 demonstrate that $1 \le n\rho^N$. Hence, we obtain that $\rho^N \ge n^{-1}$. This ensures that $\rho \ge n^{-1/N}$. The proof of Lemma 9.2.22 is thus complete. $\qquad\square$

#### 9.2.4.5 Lower and upper bounds for covering radii

**Corollary 9.2.23.** *Let $(V, \|\!|\cdot|\!\|)$ be a normed vector space, let $X = \{v \in V\colon \|\!|v|\!\| \le 1\}$, let $d\colon X \times X \to [0,\infty)$ satisfy for all $v, w \in X$ that $d(v,w) = \|\!|v - w|\!\|$, and let $n, N \in \mathbb{N}$ satisfy $N = \dim(V)$. Then*

$$n^{-1/N} \le \mathcal{C}_{(X,d),n} \le 4\,(n+1)^{-1/N}. \tag{9.127}$$

*Proof of Corollary 9.2.23.* Observe that Corollary 9.2.21 and Lemma 9.2.22 establish (9.127). The proof of Corollary 9.2.23 is thus complete. $\qquad\square$

#### 9.2.4.6 Scaling property for covering radii

**Lemma 9.2.24.** *Let $(V, \|\!|\cdot|\!\|)$ be a normed vector space, let $d\colon V \times V \to [0,\infty)$ satisfy for all $v, w \in V$ that $d(v,w) = \|\!|v - w|\!\|$, let $n \in \mathbb{N}$, $r \in (0,\infty)$, and let $X \subseteq V$ and $\mathfrak{X} \subseteq V$ satisfy $\mathfrak{X} = \{rv \in V\colon v \in X\}$. Then*

$$\mathcal{C}_{(\mathfrak{X}, d|_{\mathfrak{X} \times \mathfrak{X}}),n} = r\,\mathcal{C}_{(X, d|_{X \times X}),n}. \tag{9.128}$$

*Proof of Lemma 9.2.24.* Throughout this proof let $\Phi\colon V \to V$ satisfy for all $v \in V$ that $\Phi(v) = rv$. Observe that Exercise 9.2.3 shows that

$$
\begin{aligned}
r\,\mathcal{C}_{(X,d),n} &= r\left[\inf_{x_1,x_2,\ldots,x_n \in X} \sup_{v \in X} \min_{i \in \{1,2,\ldots,n\}} d(x_i, v)\right] \\
&= \inf_{x_1,x_2,\ldots,x_n \in X} \sup_{v \in X} \min_{i \in \{1,2,\ldots,n\}} \|\!|rx_i - rv|\!\| \\
&= \inf_{x_1,x_2,\ldots,x_n \in X} \sup_{v \in X} \min_{i \in \{1,2,\ldots,n\}} \|\!|\Phi(x_i) - \Phi(v)|\!\| \\
&= \inf_{x_1,x_2,\ldots,x_n \in X} \sup_{v \in X} \min_{i \in \{1,2,\ldots,n\}} d(\Phi(x_i), \Phi(v)) \\
&= \inf_{x_1,x_2,\ldots,x_n \in X} \sup_{v \in \mathfrak{X}} \min_{i \in \{1,2,\ldots,n\}} d(\Phi(x_i), v) \\
&= \inf_{x_1,x_2,\ldots,x_n \in \mathfrak{X}} \sup_{v \in \mathfrak{X}} \min_{i \in \{1,2,\ldots,n\}} d(x_i, v) = \mathcal{C}_{(\mathfrak{X}, d|_{\mathfrak{X} \times \mathfrak{X}}),n}.
\end{aligned}
\tag{9.129}
$$

This establishes (9.128). The proof of Lemma 9.2.24 is thus complete. $\qquad\square$

### 9.2.4.7 Upper bounds for covering numbers

**Proposition 9.2.25.** *Let $(V, \|\cdot\|)$ be a normed vector space with $\dim(V) < \infty$, let $r, R \in (0, \infty)$, $X = \{v \in V \colon \|v\| \leq R\}$, and let $d \colon X \times X \to [0, \infty)$ satisfy for all $v, w \in X$ that $d(v, w) = \|v - w\|$. Then*

$$\mathcal{C}^{(X,d),r} \leq \begin{cases} 1 & : r \geq R \\ \left[\frac{4R}{r}\right]^{\dim(V)} & : r < R \end{cases} \tag{9.130}$$

*(cf. Definition 9.2.6).*

*Proof of Proposition 9.2.25.* Throughout this proof assume w.l.o.g. that $\dim(V) > 0$, assume w.l.o.g. that $r < R$, let $\lceil \cdot \rceil \colon [0, \infty) \to [0, \infty)$ satisfy for all $x \in [0, \infty)$ that

$$\lceil x \rceil = \inf([x, \infty) \cap \mathbb{N}), \tag{9.131}$$

let $N \in \mathbb{N}$ satisfy $N = \dim(V)$, let $n \in \mathbb{N}$ satisfy

$$n = \left\lceil \left[\frac{4R}{r}\right]^N - 1 \right\rceil, \tag{9.132}$$

let $\mathfrak{X} = \{v \in V \colon \|v\| \leq 1\}$, and let $\mathfrak{d} \colon \mathfrak{X} \times \mathfrak{X} \to [0, \infty)$ satisfy for all $v, w \in \mathfrak{X}$ that

$$\mathfrak{d}(v, w) = \|v - w\|. \tag{9.133}$$

Observe that Corollary 9.2.21 proves that

$$\mathcal{C}_{(\mathfrak{X},\mathfrak{d}),n} \leq 4\,(n+1)^{-1/N}. \tag{9.134}$$

The fact that

$$n + 1 = \left\lceil \left[\frac{4R}{r}\right]^N - 1 \right\rceil + 1 \geq \left[\left[\frac{4R}{r}\right]^N - 1\right] + 1 = \left[\frac{4R}{r}\right]^N \tag{9.135}$$

therefore ensures that

$$\mathcal{C}_{(\mathfrak{X},\mathfrak{d}),n} \leq 4\,(n+1)^{-1/N} \leq 4\left[\left[\frac{4R}{r}\right]^N\right]^{-1/N} = 4\left[\frac{4R}{r}\right]^{-1} = \frac{r}{R}. \tag{9.136}$$

This and Lemma 9.2.24 demonstrate that

$$\mathcal{C}_{(X,d),n} = R\,\mathcal{C}_{(\mathfrak{X},\mathfrak{d}),n} \leq R\left[\frac{r}{R}\right] = r. \tag{9.137}$$

Lemma 9.2.15 hence ensures that

$$\mathcal{C}^{(X,d),r} \leq n \leq \left[\frac{4R}{r}\right]^N = \left[\frac{4R}{r}\right]^{\dim(V)}. \tag{9.138}$$

The proof of Proposition 9.2.25 is thus complete. $\square$

# 9.3 Risk minimization

## 9.3.1 Bias-variance decomposition

**Lemma 9.3.1** (Bias-variance decomposition). *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $(S, \mathcal{S})$ be a measurable space, let $X \colon \Omega \to S$ and $Y \colon \Omega \to \mathbb{R}$ be random variables with $\mathbb{E}[|Y|^2] < \infty$, and let $\mathcal{E} \colon \mathcal{L}^2(\mathbb{P}_X; \mathbb{R}) \to [0, \infty)$ satisfy for all $f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ that $\mathcal{E}(f) = \mathbb{E}[|f(X) - Y|^2]$. Then*

*(i) it holds for all $f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ that*

$$\mathcal{E}(f) = \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] + \mathbb{E}\big[|Y - \mathbb{E}[Y|X]|^2\big], \tag{9.139}$$

*(ii) it holds for all $f, g \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ that*

$$\mathcal{E}(f) - \mathcal{E}(g) = \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] - \mathbb{E}\big[|g(X) - \mathbb{E}[Y|X]|^2\big], \tag{9.140}$$

*and*

*(iii) it holds for all $f, g \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ that*

$$\mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] = \mathbb{E}\big[|g(X) - \mathbb{E}[Y|X]|^2\big] + \big(\mathcal{E}(f) - \mathcal{E}(g)\big). \tag{9.141}$$

*Proof of Lemma 9.3.1.* First, observe that the assumption that for all $f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ it holds that $\mathcal{E}(f) = \mathbb{E}[|f(X) - Y|^2]$ shows that for all $f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ it holds that

$$
\begin{aligned}
\mathcal{E}(f) &= \mathbb{E}\big[|f(X) - Y|^2\big] = \mathbb{E}\big[|(f(X) - \mathbb{E}[Y|X]) + (\mathbb{E}[Y|X] - Y)|^2\big] \\
&= \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] + 2\,\mathbb{E}\big[\big(f(X) - \mathbb{E}[Y|X]\big)\big(\mathbb{E}[Y|X] - Y\big)\big] + \mathbb{E}\big[|\mathbb{E}[Y|X] - Y|^2\big] \\
&= \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] + 2\,\mathbb{E}\Big[\mathbb{E}\big[\big(f(X) - \mathbb{E}[Y|X]\big)\big(\mathbb{E}[Y|X] - Y\big)|X\big]\Big] + \mathbb{E}\big[|\mathbb{E}[Y|X] - Y|^2\big] \\
&= \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] + 2\,\mathbb{E}\Big[\big(f(X) - \mathbb{E}[Y|X]\big)\mathbb{E}\big[\big(\mathbb{E}[Y|X] - Y\big)|X\big]\Big] + \mathbb{E}\big[|\mathbb{E}[Y|X] - Y|^2\big] \\
&= \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] + 2\,\mathbb{E}\big[\big(f(X) - \mathbb{E}[Y|X]\big)\big(\mathbb{E}[Y|X] - \mathbb{E}[Y|X]\big)\big] + \mathbb{E}\big[|\mathbb{E}[Y|X] - Y|^2\big] \\
&= \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] + \mathbb{E}\big[|\mathbb{E}[Y|X] - Y|^2\big].
\end{aligned}
\tag{9.142}
$$

This implies that for all $f, g \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ it holds that

$$\mathcal{E}(f) - \mathcal{E}(g) = \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] - \mathbb{E}\big[|g(X) - \mathbb{E}[Y|X]|^2\big]. \tag{9.143}$$

Hence, we obtain that for all $f, g \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ it holds that

$$\mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] = \mathbb{E}\big[|g(X) - \mathbb{E}[Y|X]|^2\big] + \mathcal{E}(f) - \mathcal{E}(g). \tag{9.144}$$

Combining this with (9.142) and (9.143) establishes items (i), (ii), and (iii). The proof of Lemma 9.3.1 is thus complete. □

## 9.3.2 Risk minimization for measurable functions

**Proposition 9.3.2.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $(S, \mathcal{S})$ be a measurable space, let $X\colon \Omega \to S$ and $Y\colon \Omega \to \mathbb{R}$ be random variables, assume $\mathbb{E}[|Y|^2] < \infty$, let $\mathcal{E}\colon \mathcal{L}^2(\mathbb{P}_X; \mathbb{R}) \to [0, \infty)$ satisfy for all $f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ that $\mathcal{E}(f) = \mathbb{E}[|f(X) - Y|^2]$. Then it holds that*

$$
\begin{aligned}
\big\{ f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R}) \colon \mathcal{E}(f) = \inf\nolimits_{g \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})} \mathcal{E}(g) \big\} &= \big\{ f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R}) \colon \mathcal{E}(f) = \mathbb{E}\big[|\mathbb{E}[Y|X] - Y|^2\big] \big\} \\
&= \{ f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R}) \colon f(X) = \mathbb{E}[Y|X] \ \mathbb{P}\text{-}a.s. \}.
\end{aligned}
\tag{9.145}
$$

*Proof of Proposition 9.3.2.* Note that Lemma 9.3.1 shows that for all $g \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ it holds that

$$
\mathcal{E}(g) = \mathbb{E}\big[|g(X) - \mathbb{E}[Y|X]|^2\big] + \mathbb{E}\big[|\mathbb{E}[Y|X] - Y|^2\big].
\tag{9.146}
$$

Hence, we obtain that for all $g \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R})$ it holds that

$$
\mathcal{E}(g) \geq \mathbb{E}\big[|\mathbb{E}[Y|X] - Y|^2\big].
\tag{9.147}
$$

Furthermore, note that (9.146) shows that

$$
\begin{aligned}
\big\{ f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R}) \colon \mathcal{E}(f) = \mathbb{E}\big[|\mathbb{E}[Y|X] - Y|^2\big] \big\} &= \big\{ f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R}) \colon \mathbb{E}\big[|f(X) - \mathbb{E}[Y|X]|^2\big] = 0 \big\} \\
&= \{ f \in \mathcal{L}^2(\mathbb{P}_X; \mathbb{R}) \colon f(X) = \mathbb{E}[Y|X] \ \mathbb{P}\text{-a.s.} \}.
\end{aligned}
\tag{9.148}
$$

Combining this with (9.147) establishes (9.145). The proof of Proposition 9.3.2 is thus complete. $\qquad \square$

**Proposition 9.3.3.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $(S, \mathcal{S})$ be a measurable space, let $X\colon \Omega \to S$ be a random variable, let $\mathcal{M} = \{ (f\colon S \to \mathbb{R}) \colon f \text{ is } \mathcal{S}/\mathcal{B}(\mathbb{R})\text{-measurable} \}$, let $\varphi \in \mathcal{M}$, let $\mathcal{E}\colon \mathcal{M} \to [0, \infty)$ satisfy for all $f \in \mathcal{M}$ that $\mathcal{E}(f) = \mathbb{E}\big[|f(X) - \varphi(X)|^2\big]$. Then it holds that*

$$
\{ f \in \mathcal{M} \colon \mathcal{E}(f) = \inf\nolimits_{g \in \mathcal{M}} \mathcal{E}(g) \} = \{ f \in \mathcal{M} \colon \mathcal{E}(f) = 0 \} = \{ f \in \mathcal{M} \colon \mathbb{P}(f(X) = \varphi(X)) = 1 \}.
\tag{9.149}
$$

*Proof of Proposition 9.3.3.* Note that the assumption that for all $f \in \mathcal{M}$ it holds that $\mathcal{E}(f) = \mathbb{E}[|f(X) - \varphi(X)|^2]$ implies that $\mathcal{E}(\varphi) = 0$. Hence, we obtain that

$$
\inf_{g \in \mathcal{M}} \mathcal{E}(g) = 0.
\tag{9.150}
$$

Furthermore, observe that

$$
\begin{aligned}
\{ f \in \mathcal{M} \colon \mathcal{E}(f) = 0 \} &= \big\{ f \in \mathcal{M} \colon \mathbb{E}\big[|f(X) - \varphi(X)|^2\big] = 0 \big\} \\
&= \big\{ f \in \mathcal{M} \colon \mathbb{P}\big(\{\omega \in \Omega \colon f(X(\omega)) \neq \varphi(X(\omega))\}\big) = 0 \big\} \\
&= \big\{ f \in \mathcal{M} \colon \mathbb{P}\big(X^{-1}(\{x \in S \colon f(x) \neq \varphi(x)\})\big) = 0 \big\} \\
&= \{ f \in \mathcal{M} \colon \mathbb{P}_X(\{x \in S \colon f(x) \neq \varphi(x)\}) = 0 \}.
\end{aligned}
\tag{9.151}
$$

The proof of Proposition 9.3.3 is thus complete. $\qquad \square$

### 9.3.3 Risk minimization for continuous functions

**Proposition 9.3.4.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $(E, \delta)$ be a metric space, let $X\colon \Omega \to E$ be a random variable, assume for all $x \in E$, $r \in (0,\infty)$ that $\mathbb{P}_X(\{y \in E\colon d(x,y) < r\}) > 0$, let $\varphi \in C(E,\mathbb{R})$, let $\mathcal{E}\colon C(E,\mathbb{R}) \to [0,\infty]$ satisfy for all $f \in C(E,\mathbb{R})$ that $\mathcal{E}(f) = \mathbb{E}\big[|f(X) - \varphi(X)|^2\big]$. Then it holds that*

$$\{f \in C(E,\mathbb{R})\colon \mathcal{E}(f) = \inf_{g \in C(E,\mathbb{R})} \mathcal{E}(g)\} = \{f \in C(E,\mathbb{R})\colon \mathcal{E}(f) = 0\} = \{\varphi\}. \quad (9.152)$$

*Proof of Proposition 9.3.4.* Note that the assumption that for all $f \in C(E,\mathbb{R})$ it holds that $\mathcal{E}(f) = \mathbb{E}[|f(X) - \varphi(X)|^2]$ implies that $\mathcal{E}(\varphi) = 0$. Furthermore, note that the fact that $\varphi \in C(E,\mathbb{R})$ implies that for all $f \in C(E,\mathbb{R})$, $x \in E$ with $f(x) \neq \varphi(x)$ there exists $r \in (0,\infty)$ such that

$$\{y \in E\colon d(x,y) < r\} \subseteq \{y \in E\colon f(y) \neq \varphi(y)\}. \quad (9.153)$$

Combining this with the assumption that for all $x \in E$, $r \in (0,\infty)$ it holds that $\mathbb{P}_X(\{y \in E\colon d(x,y) < r\}) > 0$ shows that for all $f \in C(E,\mathbb{R})$ with $f \neq \varphi$ it holds that

$$\mathbb{P}_X(\{y \in E\colon f(y) \neq \varphi(y)\}) > 0. \quad (9.154)$$

This implies that for all $f \in C(E,\mathbb{R})$ with $f \neq \varphi$ it holds that

$$\mathcal{E}(f) = \mathbb{E}\big[|f(X) - \varphi(X)|^2\big] = \int_S |f(x) - \varphi(x)|^2 \, \mathbb{P}_X(\mathrm{d}x) > 0. \quad (9.155)$$

The proof of Proposition 9.3.4 is thus complete. $\square$

## 9.4 Empirical risk minimization

### 9.4.1 Measurability properties for suprema

**Lemma 9.4.1.** *Let $(E, \mathscr{E})$ be a topological space, assume $E \neq \emptyset$, let $\mathbf{E} \subseteq E$ be an at most countable set, assume that $\mathbf{E}$ is dense in $E$, let $(\Omega, \mathcal{F})$ be a measurable space, let $f_x\colon \Omega \to \mathbb{R}$, $x \in E$, be $\mathcal{F}/\mathcal{B}(\mathbb{R})$-measurable functions, assume for all $\omega \in \Omega$ that $E \ni x \mapsto f_x(\omega) \in \mathbb{R}$ is a continuous function, and let $F\colon \Omega \to \mathbb{R} \cup \{\infty\}$ satisfy for all $\omega \in \Omega$ that $F(\omega) = \sup_{x \in E} f_x(\omega)$. Then*

*(i) it holds for all $\omega \in \Omega$ that $F(\omega) = \sup_{x \in \mathbf{E}} f_x(\omega)$ and*

*(ii) it holds that $F$ is an $\mathcal{F}/\mathcal{B}(\mathbb{R} \cup \{\infty\})$-measurable function.*

*Proof of Lemma 9.4.1.* Note that the assumption that $\mathbf{E}$ is dense in $E$ implies that for all $g \in C(E,\mathbb{R})$ it holds that

$$\sup_{x \in E} g(x) = \sup_{x \in \mathbf{E}} g(x). \quad (9.156)$$

This and the assumption that for all $\omega \in \Omega$ it holds that $E \ni x \mapsto f_x(\omega) \in \mathbb{R}$ is a continuous function show that for all $\omega \in \Omega$ it holds that

$$F(\omega) = \sup_{x \in E} f_x(\omega) = \sup_{x \in \mathbf{E}} f_x(\omega). \quad (9.157)$$

This establishes item (i). Next note that item (i) and the assumption that for all $x \in E$ it holds that $f_x\colon \Omega \to \mathbb{R}$ is an $\mathcal{F}/\mathcal{B}(\mathbb{R})$-measurable function demonstrate item (ii). The proof of Lemma 9.4.1 is thus complete. $\square$

**Lemma 9.4.2.** *Let $(E, \delta)$ be a separable metric space, assume $E \neq \emptyset$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $L \in \mathbb{R}$, and let $Z_x \colon \Omega \to \mathbb{R}$, $x \in E$, be random variables which satisfy for all $x, y \in E$ that $\mathbb{E}[|Z_x|] < \infty$ and $|Z_x - Z_y| \leq L\delta(x, y)$. Then*

(i) *it holds for all $x, y \in E$, $\eta \in \Omega$ that $|(Z_x(\eta) - \mathbb{E}[Z_x]) - (Z_y(\eta) - \mathbb{E}[Z_y])| \leq 2L\delta(x, y)$ and*

(ii) *it holds that $\Omega \ni \eta \mapsto \sup_{x \in E} |Z_x(\eta) - \mathbb{E}[Z_x]| \in [0, \infty]$ is an $\mathcal{F}/\mathcal{B}([0, \infty])$-measurable function.*

*Proof of Lemma 9.4.2.* Note that the assumption that for all $x, y \in E$ it holds that $|Z_x - Z_y| \leq L\delta(x, y)$ shows that for all $x, y \in E$, $\eta \in \Omega$ it holds that

$$
\begin{aligned}
|(Z_x(\eta) - \mathbb{E}[Z_x]) - (Z_y(\eta) - \mathbb{E}[Z_y])| &= |(Z_x(\eta) - Z_y(\eta)) + (\mathbb{E}[Z_y] - \mathbb{E}[Z_x])| \\
&\leq |Z_x(\eta) - Z_y(\eta)| + |\mathbb{E}[Z_x] - \mathbb{E}[Z_y]| \leq L\delta(x, y) + |\mathbb{E}[Z_x] - \mathbb{E}[Z_y]| \\
&= L\delta(x, y) + |\mathbb{E}[Z_x - Z_y]| \leq L\delta(x, y) + \mathbb{E}[|Z_x - Z_y|] \leq L\delta(x, y) + L\delta(x, y) = 2L\delta(x, y).
\end{aligned}
$$
(9.158)

This proves item (i). Next observe that item (i) implies that for all $\eta \in \Omega$ it holds that $E \ni x \mapsto |Z_x(\eta) - \mathbb{E}[Z_x]| \in \mathbb{R}$ is a continuous function. Combining this and the assumption that $E$ is separable with Lemma 9.4.1 establishes item (ii). The proof of Lemma 9.4.2 is thus complete. $\qquad\square$

## 9.4.2   Concentration inequalities for random fields

**Lemma 9.4.3.** *Let $(E, d)$ be a separable metric space and let $F \subseteq E$ be a set. Then*

$$
(F, d|_{F \times F})
$$
(9.159)

*is a separable metric space.*

*Proof of Lemma 9.4.3.* Throughout this proof assume w.l.o.g. that $F \neq \emptyset$, let $e = (e_n)_{n \in \mathbb{N}} \colon \mathbb{N} \to E$ be a sequence of elements in $E$ such that $\{e_n \in E \colon n \in \mathbb{N}\}$ is dense in $E$, and let $f = (f_n)_{n \in \mathbb{N}} \colon \mathbb{N} \to F$ be a sequence of elements in $F$ such that for all $n \in \mathbb{N}$ it holds that

$$
d(f_n, e_n) \leq \begin{cases} 0 & : e_n \in F \\ \left[\inf_{x \in F} d(x, e_n)\right] + \frac{1}{2^n} & : e_n \notin F. \end{cases}
$$
(9.160)

Observe that for all $v \in F \setminus \{e_m \in E \colon m \in \mathbb{N}\}$, $n \in \mathbb{N}$ it holds that

$$
\begin{aligned}
\inf_{m \in \mathbb{N}} d(v, f_m) &\leq \inf_{m \in \mathbb{N} \cap [n, \infty)} d(v, f_m) \\
&\leq \inf_{m \in \mathbb{N} \cap [n, \infty)} [d(v, e_m) + d(e_m, f_m)] \\
&\leq \inf_{m \in \mathbb{N} \cap [n, \infty)} \left[ d(v, e_m) + \left[\inf_{x \in F} d(x, e_m)\right] + \frac{1}{2^m} \right] \\
&\leq \inf_{m \in \mathbb{N} \cap [n, \infty)} \left[ 2\, d(v, e_m) + \frac{1}{2^m} \right] \\
&\leq 2\left[ \inf_{m \in \mathbb{N} \cap [n, \infty)} d(v, e_m) \right] + \frac{1}{2^n} = \frac{1}{2^n}.
\end{aligned}
$$
(9.161)

Combining this with the fact that for all $v \in F \cap \{e_m \in E \colon m \in \mathbb{N}\}$ it holds that $\inf_{m \in \mathbb{N}} d(v, f_m) = 0$ ensures that the set $\{f_n \in F \colon n \in \mathbb{N}\}$ is dense in $F$. The proof of Lemma 9.4.3 is thus complete. $\qquad\square$

**Lemma 9.4.4.** *Let $(E, \delta)$ be a separable metric space, let $\varepsilon, L \in \mathbb{R}$, $N \in \mathbb{N}$, $z_1, z_2, \ldots, z_N \in E$ satisfy $E \subseteq \bigcup_{i=1}^{N} \{x \in E \colon 2L\delta(x, z_i) \leq \varepsilon\}$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, and let $Z_x \colon \Omega \to \mathbb{R}$, $x \in E$, be random variables which satisfy for all $x, y \in E$ that $|Z_x - Z_y| \leq L\delta(x, y)$. Then*

$$\mathbb{P}(\sup_{x \in E} |Z_x| \geq \varepsilon) \leq \sum_{i=1}^{N} \mathbb{P}(|Z_{z_i}| \geq \tfrac{\varepsilon}{2}) \tag{9.162}$$

*(cf. Lemma 9.4.1).*

*Proof of Lemma 9.4.4.* Throughout this proof let $B_1, B_2, \ldots, B_N \subseteq E$ satisfy for all $i \in \{1, 2, \ldots, N\}$ that $B_i = \{x \in E \colon 2L\delta(x, z_i) \leq \varepsilon\}$. Observe that the triangle inequality and the assumption that for all $x, y \in E$ it holds that $|Z_x - Z_y| \leq L\delta(x, y)$ show that for all $i \in \{1, 2, \ldots, N\}$, $x \in B_i$ it holds that

$$|Z_x| = |Z_x - Z_{z_i} + Z_{z_i}| \leq |Z_x - Z_{z_i}| + |Z_{z_i}| \leq L\delta(x, z_i) + |Z_{z_i}| \leq \tfrac{\varepsilon}{2} + |Z_{z_i}|. \tag{9.163}$$

Combining this with Lemma 9.4.1 and Lemma 9.4.3 proves that for all $i \in \{1, 2, \ldots, N\}$ it holds that

$$\mathbb{P}(\sup_{x \in B_i} |Z_x| \geq \varepsilon) \leq \mathbb{P}(\tfrac{\varepsilon}{2} + |Z_{z_i}| \geq \varepsilon) = \mathbb{P}(|Z_{z_i}| \geq \tfrac{\varepsilon}{2}). \tag{9.164}$$

This, Lemma 9.4.1, and Lemma 9.4.3 establish that

$$\begin{aligned}
\mathbb{P}(\sup_{x \in E} |Z_x| \geq \varepsilon) &= \mathbb{P}\left(\sup_{x \in \left(\bigcup_{i=1}^{N} B_i\right)} |Z_x| \geq \varepsilon\right) = \mathbb{P}\left(\bigcup_{i=1}^{N} \{\sup_{x \in B_i} |Z_x| \geq \varepsilon\}\right) \\
&\leq \sum_{i=1}^{N} \mathbb{P}(\sup_{x \in B_i} |Z_x| \geq \varepsilon) \leq \sum_{i=1}^{N} \mathbb{P}(|Z_{z_i}| \geq \tfrac{\varepsilon}{2}).
\end{aligned} \tag{9.165}$$

This completes the proof of Lemma 9.4.4. $\qquad\square$

**Lemma 9.4.5.** *Let $(E, \delta)$ be a separable metric space, assume $E \neq \emptyset$, let $\varepsilon, L \in (0, \infty)$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, and let $Z_x \colon \Omega \to \mathbb{R}$, $x \in E$, be random variables which satisfy for all $x, y \in E$ that $|Z_x - Z_y| \leq L\delta(x, y)$. Then*

$$\left[\mathcal{C}^{(E,\delta), \frac{\varepsilon}{2L}}\right]^{-1} \mathbb{P}(\sup_{x \in E} |Z_x| \geq \varepsilon) \leq \sup_{x \in E} \mathbb{P}(|Z_x| \geq \tfrac{\varepsilon}{2}). \tag{9.166}$$

*(cf. Definition 9.2.6 and Lemma 9.4.1).*

*Proof of Lemma 9.4.5.* Throughout this proof let $N \in \mathbb{N} \cup \{\infty\}$ satisfy $N = \mathcal{C}^{(E,\delta), \frac{\varepsilon}{2L}}$, assume without loss of generality that $N < \infty$, and let $z_1, z_2, \ldots, z_N \in E$ satisfy $E \subseteq \bigcup_{i=1}^{N} \{x \in E \colon \delta(x, z_i) \leq \tfrac{\varepsilon}{2L}\}$ (cf. Definition 9.2.6). Observe that Lemma 9.4.1 and Lemma 9.4.4 establish that

$$\mathbb{P}(\sup_{x \in E} |Z_x| \geq \varepsilon) \leq \sum_{i=1}^{N} \mathbb{P}(|Z_{z_i}| \geq \tfrac{\varepsilon}{2}) \leq N\left[\sup_{x \in E} \mathbb{P}(|Z_x| \geq \tfrac{\varepsilon}{2})\right]. \tag{9.167}$$

This completes the proof of Lemma 9.4.5. $\qquad\square$

**Lemma 9.4.6.** *Let $(E, \delta)$ be a separable metric space, assume $E \neq \emptyset$, let $\varepsilon, L \in (0, \infty)$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, and let $Z_x \colon \Omega \to \mathbb{R}$, $x \in E$, be random variables which satisfy for all $x, y \in E$ that $\mathbb{E}[|Z_x|] < \infty$ and $|Z_x - Z_y| \leq L\delta(x, y)$. Then*

$$\left[\mathcal{C}^{(E,\delta), \frac{\varepsilon}{4L}}\right]^{-1} \mathbb{P}(\sup_{x \in E} |Z_x - \mathbb{E}[Z_x]| \geq \varepsilon) \leq \sup_{x \in E} \mathbb{P}(|Z_x - \mathbb{E}[Z_x]| \geq \tfrac{\varepsilon}{2}). \tag{9.168}$$

*(cf. Definition 9.2.6 and Lemma 9.4.2).*

*Proof of Lemma 9.4.6.* Throughout this proof let $Y_x \colon \Omega \to \mathbb{R}$, $x \in E$, satisfy for all $x \in E$, $\eta \in \Omega$ that $Y_x(\eta) = Z_x(\eta) - \mathbb{E}[Z_x]$. Observe that Lemma 9.4.2 ensures that for all $x, y \in E$ it holds that

$$|Y_x - Y_y| \le 2L\delta(x, y). \tag{9.169}$$

This and Lemma 9.4.5 (applied with $(E, \delta) \curvearrowright (E, \delta)$, $\varepsilon \curvearrowright \varepsilon$, $L \curvearrowright 2L$, $(\Omega, \mathcal{F}, \mathbb{P}) \curvearrowright (\Omega, \mathcal{F}, \mathbb{P})$, $(Z_x)_{x \in E} \curvearrowright (Y_x)_{x \in E}$ in the notation of Lemma 9.4.5) establish (9.168). The proof of Lemma 9.4.6 is thus complete. $\qquad\square$

**Lemma 9.4.7.** *Let $(E, \delta)$ be a separable metric space, assume $E \ne \emptyset$, let $M \in \mathbb{N}$, $\varepsilon, L, D \in (0, \infty)$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, for every $x \in E$ let $Y_{x,1}, Y_{x,2}, \ldots, Y_{x,M} \colon \Omega \to [0, D]$ be independent random variables, assume for all $x, y \in E$, $m \in \{1, 2, \ldots, M\}$ that $|Y_{x,m} - Y_{y,m}| \le L\delta(x, y)$, and let $Z_x \colon \Omega \to [0, \infty)$, $x \in E$, satisfy for all $x \in E$ that*

$$Z_x = \frac{1}{M} \left[ \sum_{m=1}^{M} Y_{x,m} \right]. \tag{9.170}$$

*Then*

*(i) it holds for all $x \in E$ that $\mathbb{E}[|Z_x|] \le D < \infty$,*

*(ii) it holds that $\Omega \ni \eta \mapsto \sup_{x \in E} |Z_x(\eta) - \mathbb{E}[Z_x]| \in [0, \infty]$ is an $\mathcal{F}/\mathcal{B}([0, \infty])$-measurable function, and*

*(iii) it holds that*

$$\mathbb{P}(\sup_{x \in E} |Z_x - \mathbb{E}[Z_x]| \ge \varepsilon) \le 2\mathcal{C}^{(E,\delta),\frac{\varepsilon}{4L}} \exp\left( \frac{-\varepsilon^2 M}{2D^2} \right) \tag{9.171}$$

*(cf. Definition 9.2.6).*

*Proof of Lemma 9.4.7.* First, observe that the triangle inequality and the assumption that for all $x, y \in E$, $m \in \{1, 2, \ldots, M\}$ it holds that $|Y_{x,m} - Y_{y,m}| \le L\delta(x, y)$ imply that for all $x, y \in E$ it holds that

$$
\begin{aligned}
|Z_x - Z_y| &= \left| \frac{1}{M} \left[ \sum_{m=1}^{M} Y_{x,m} \right] - \frac{1}{M} \left[ \sum_{m=1}^{M} Y_{y,m} \right] \right| = \frac{1}{M} \left| \sum_{m=1}^{M} (Y_{x,m} - Y_{y,m}) \right| \\
&\le \frac{1}{M} \left[ \sum_{m=1}^{M} |Y_{x,m} - Y_{y,m}| \right] \le L\delta(x, y).
\end{aligned}
\tag{9.172}
$$

Next note that the assumption that for all $x \in E$, $m \in \{1, 2, \ldots, M\}$, $\omega \in \Omega$ it holds that $|Y_{x,m}(\omega)| \in [0, D]$ ensures that for all $x \in E$ it holds that

$$\mathbb{E}[|Z_x|] = \mathbb{E}\left[ \frac{1}{M} \left[ \sum_{m=1}^{M} Y_{x,m} \right] \right] = \frac{1}{M} \left[ \sum_{m=1}^{M} \mathbb{E}[Y_{x,m}] \right] \le D < \infty. \tag{9.173}$$

This proves item (i). Furthermore, note that item (i), (9.172), and Lemma 9.4.2 establish item (ii). Next observe that (9.170) shows that for all $x \in E$ it holds that

$$|Z_x - \mathbb{E}[Z_x]| = \left| \frac{1}{M} \left[ \sum_{m=1}^{M} Y_{x,m} \right] - \mathbb{E}\left[ \frac{1}{M} \left[ \sum_{m=1}^{M} Y_{x,m} \right] \right] \right| = \frac{1}{M} \left| \sum_{m=1}^{M} (Y_{x,m} - \mathbb{E}[Y_{x,m}]) \right|. \tag{9.174}$$

Combining this with Corollary 9.1.21 (applied with $(\Omega, \mathcal{F}, \mathbb{P}) \curvearrowright (\Omega, \mathcal{F}, \mathbb{P})$, $N \curvearrowright M$, $\varepsilon \curvearrowright \frac{\varepsilon}{2}$, $(a_1, a_2, \ldots, a_N) \curvearrowright (0, 0, \ldots, 0)$, $(b_1, b_2, \ldots, b_N) \curvearrowright (D, D, \ldots, D)$, $(X_n)_{n \in \{1,2,\ldots,N\}} \curvearrowright (Y_{x,m})_{m \in \{1,2,\ldots,M\}}$ for $x \in E$ in the notation of Corollary 9.1.21) ensures that for all $x \in E$ it holds that

$$\mathbb{P}\big(|Z_x - \mathbb{E}[Z_x]| \geq \tfrac{\varepsilon}{2}\big) \leq 2 \exp\left(\frac{-2\left[\frac{\varepsilon}{2}\right]^2 M^2}{MD^2}\right) = 2 \exp\left(\frac{-\varepsilon^2 M}{2D^2}\right). \tag{9.175}$$

Combining this, (9.172), and (9.173) with Lemma 9.4.6 establishes item (iii). The proof of Lemma 9.4.7 is thus complete. $\qquad\square$

### 9.4.2.1 Uniform estimates for the statistical learning error

**Lemma 9.4.8.** *Let $(E, \delta)$ be a separable metric space, assume $E \neq \emptyset$, let $M \in \mathbb{N}$, $\varepsilon, L, D \in (0, \infty)$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $X_{x,m} \colon \Omega \to \mathbb{R}$, $x \in E$, $m \in \{1, 2, \ldots, M\}$, and $Y_m \colon \Omega \to \mathbb{R}$, $m \in \{1, 2, \ldots, M\}$, be functions, assume for all $x \in E$ that $(X_{x,m}, Y_m)$, $m \in \{1, 2, \ldots, M\}$, are i.i.d. random variables, assume for all $x, y \in E$, $m \in \{1, 2, \ldots, M\}$ that $|X_{x,m} - X_{y,m}| \leq L\delta(x, y)$ and $|X_{x,m} - Y_m| \leq D$, let $\mathfrak{E}_x \colon \Omega \to [0, \infty)$, $x \in E$, satisfy for all $x \in E$ that*

$$\mathfrak{E}_x = \frac{1}{M}\left[\sum_{m=1}^{M} |X_{x,m} - Y_m|^2\right], \tag{9.176}$$

*and let $\mathcal{E}_x \in [0, \infty)$, $x \in E$, satisfy for all $x \in E$ that $\mathcal{E}_x = \mathbb{E}[|X_{x,1} - Y_1|^2]$. Then $\Omega \ni \omega \mapsto \sup_{x \in E}|\mathfrak{E}_x(\omega) - \mathcal{E}_x| \in [0, \infty]$ is an $\mathcal{F}/\mathcal{B}([0, \infty])$-measurable function and*

$$\mathbb{P}(\sup_{x \in E}|\mathfrak{E}_x - \mathcal{E}_x| \geq \varepsilon) \leq 2\mathcal{C}^{(E,\delta), \frac{\varepsilon}{8LD}} \exp\left(\frac{-\varepsilon^2 M}{2D^4}\right) \tag{9.177}$$

*(cf. Definition 9.2.6).*

*Proof of Lemma 9.4.8.* Throughout this proof let $\mathscr{E}_{x,m} \colon \Omega \to [0, D^2]$, $x \in E$, $m \in \{1, 2, \ldots, M\}$, satisfy for all $x \in E$, $m \in \{1, 2, \ldots, M\}$ that

$$\mathscr{E}_{x,m} = |X_{x,m} - Y_m|^2. \tag{9.178}$$

Observe that the fact that for all $x_1, x_2, y \in \mathbb{R}$ it holds that $(x_1 - y)^2 - (x_2 - y)^2 = (x_1 - x_2)((x_1 - y) + (x_2 - y))$, the assumption that for all $x \in E$, $m \in \{1, 2, \ldots, M\}$ it holds that $|X_{x,m} - Y_m| \leq D$, and the assumption that for all $x, y \in E$, $m \in \{1, 2, \ldots, M\}$ it holds that $|X_{x,m} - X_{y,m}| \leq L\delta(x, y)$ imply that for all $x, y \in E$, $m \in \{1, 2, \ldots, M\}$ it holds that

$$|\mathscr{E}_{x,m} - \mathscr{E}_{y,m}| = \big|(X_{x,m} - Y_m)^2 - (X_{y,m} - Y_m)^2\big| = |X_{x,m} - X_{y,m}|\big|(X_{x,m} - Y_m) + (X_{y,m} - Y_m)\big|$$
$$\leq |X_{x,m} - X_{y,m}|\big(|X_{x,m} - Y_m| + |X_{y,m} - Y_m|\big) \leq 2D|X_{x,m} - X_{y,m}| \leq 2LD\delta(x, y). \tag{9.179}$$

In addition, note that (9.176) and the assumption that for all $x \in E$ it holds that $(X_{x,m}, Y_m)$, $m \in \{1, 2, \ldots, M\}$, are i.i.d. random variables show that for all $x \in E$ it holds that

$$\mathbb{E}[\mathfrak{E}_x] = \frac{1}{M}\left[\sum_{m=1}^{M} \mathbb{E}[|X_{x,m} - Y_m|^2]\right] = \frac{1}{M}\left[\sum_{m=1}^{M} \mathbb{E}[|X_{x,1} - Y_1|^2]\right] = \frac{1}{M}\left[\sum_{m=1}^{M} \mathcal{E}_x\right] = \mathcal{E}_x. \tag{9.180}$$

Furthermore, observe that the assumption that for all $x \in E$ it holds that $(X_{x,m}, Y_m)$, $m \in \{1, 2, \ldots, M\}$, are i.i.d. random variables ensures that for all $x \in E$ it holds that $\mathscr{E}_{x,m}$, $m \in \{1, 2, \ldots, M\}$, are i.i.d. random variables. Combining this, (9.179), and (9.180) with Lemma 9.4.7 (applied with $(E, \delta) \curvearrowleft (E, \delta)$, $M \curvearrowleft M$, $\varepsilon \curvearrowleft \varepsilon$, $L \curvearrowleft 2LD$, $D \curvearrowleft D^2$, $(\Omega, \mathcal{F}, \mathbb{P}) \curvearrowleft (\Omega, \mathcal{F}, \mathbb{P})$, $(Y_{x,m})_{x \in E, m \in \{1,2,\ldots,M\}} \curvearrowleft (\mathscr{E}_{x,m})_{x \in E, m \in \{1,2,\ldots,M\}}$, $(Z_x)_{x \in E} = (\mathfrak{E}_x)_{x \in E}$ in the notation of Lemma 9.4.7) establishes (9.177). The proof of Lemma 9.4.8 is thus complete. $\qquad\square$

**Lemma 9.4.9.** *Let $d, \mathfrak{d}, M \in \mathbb{N}$, $R, L, \mathcal{R}, \varepsilon \in (0, \infty)$, let $D \subseteq \mathbb{R}^d$ be a compact set, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $X_m \colon \Omega \to D$, $m \in \{1, 2, \ldots, M\}$, and $Y_m \colon \Omega \to \mathbb{R}$, $m \in \{1, 2, \ldots, M\}$, be functions, assume that $(X_m, Y_m)$, $m \in \{1, 2, \ldots, M\}$, are i.i.d. random variables, let $H = (H_\theta)_{\theta \in [-R,R]^{\mathfrak{d}}} \colon [-R, R]^{\mathfrak{d}} \to C(D, \mathbb{R})$ satisfy for all $\theta, \vartheta \in [-R, R]^{\mathfrak{d}}$, $x \in D$ that $|H_\theta(x) - H_\vartheta(x)| \leq L\|\theta - \vartheta\|_\infty$, assume for all $\theta \in [-R, R]^{\mathfrak{d}}$, $m \in \{1, 2, \ldots, M\}$ that $|H_\theta(X_m) - Y_m| \leq \mathcal{R}$ and $\mathbb{E}[|Y_1|^2] < \infty$, let $\mathcal{E} \colon C(D, \mathbb{R}) \to [0, \infty)$ satisfy for all $f \in C(D, \mathbb{R})$ that $\mathcal{E}(f) = \mathbb{E}[|f(X_1) - Y_1|^2]$, and let $\mathfrak{E} \colon [-R, R]^{\mathfrak{d}} \times \Omega \to [0, \infty)$ satisfy for all $\theta \in [-R, R]^{\mathfrak{d}}$, $\omega \in \Omega$ that*

$$\mathfrak{E}(\theta, \omega) = \frac{1}{M}\left[\sum_{m=1}^{M} |H_\theta(X_m(\omega)) - Y_m(\omega)|^2\right] \tag{9.181}$$

*(cf. Definition 3.1.16). Then $\Omega \ni \omega \mapsto \sup_{\theta \in [-R,R]^{\mathfrak{d}}}|\mathfrak{E}(\theta, \omega) - \mathcal{E}(H_\theta)| \in [0, \infty]$ is an $\mathcal{F}/\mathcal{B}([0, \infty])$-measurable function and*

$$\mathbb{P}\big(\sup\nolimits_{\theta \in [-R,R]^{\mathfrak{d}}}|\mathfrak{E}(\theta) - \mathcal{E}(H_\theta)| \geq \varepsilon\big) \leq 2\max\left\{1, \left[\frac{32LR\mathcal{R}}{\varepsilon}\right]^{\mathfrak{d}}\right\}\exp\left(\frac{-\varepsilon^2 M}{2\mathcal{R}^4}\right). \tag{9.182}$$

*Proof of Lemma 9.4.9.* Throughout this proof let $B \subseteq \mathbb{R}^{\mathfrak{d}}$ satisfy $B = [-R, R]^{\mathfrak{d}} = \{\theta \in \mathbb{R}^{\mathfrak{d}} \colon \|\theta\|_\infty \leq R\}$ and let $\delta \colon B \times B \to [0, \infty)$ satisfy for all $\theta, \vartheta \in B$ that

$$\delta(\theta, \vartheta) = \|\theta - \vartheta\|_\infty. \tag{9.183}$$

Observe that the assumption that $(X_m, Y_m)$, $m \in \{1, 2, \ldots, M\}$, are i.i.d. random variables and the assumption that for all $\theta \in [-R, R]^{\mathfrak{d}}$ it holds that $H_\theta$ is a continuous function imply that for all $\theta \in B$ it holds that $(H_\theta(X_m), Y_m)$, $m \in \{1, 2, \ldots, M\}$, are i.i.d. random variables. Combining this, the assumption that for all $\theta, \vartheta \in B$, $x \in D$ it holds that $|H_\theta(x) - H_\vartheta(x)| \leq L\|\theta - \vartheta\|_\infty$, and the assumption that for all $\theta \in B$, $m \in \{1, 2, \ldots, M\}$ it holds that $|H_\theta(X_m) - Y_m| \leq \mathcal{R}$ with Lemma 9.4.8 (applied with $(E, \delta) \curvearrowleft (B, \delta)$, $M \curvearrowleft M$, $\varepsilon \curvearrowleft \varepsilon$, $L \curvearrowleft L$, $D \curvearrowleft \mathcal{R}$, $(\Omega, \mathcal{F}, \mathbb{P}) \curvearrowleft (\Omega, \mathcal{F}, \mathbb{P})$, $(X_{x,m})_{x \in E, m \in \{1,2,\ldots,M\}} \curvearrowleft (H_\theta(X_m))_{\theta \in B, m \in \{1,2,\ldots,M\}}$, $(Y_m)_{m \in \{1,2,\ldots,M\}} \curvearrowleft (Y_m)_{m \in \{1,2,\ldots,M\}}$, $(\mathfrak{E}_x)_{x \in E} \curvearrowleft ((\Omega \ni \omega \mapsto \mathfrak{E}(\theta, \omega) \in [0, \infty)))_{\theta \in B}$, $(\mathcal{E}_x)_{x \in E} \curvearrowleft (\mathcal{E}(H_\theta))_{\theta \in B}$ in the notation of Lemma 9.4.8) establishes that $\Omega \ni \omega \mapsto \sup_{\theta \in B}|\mathfrak{E}(\theta, \omega) - \mathcal{E}(H_\theta)| \in [0, \infty]$ is an $\mathcal{F}/\mathcal{B}([0, \infty])$-measurable function and

$$\mathbb{P}\big(\sup\nolimits_{\theta \in B}|\mathfrak{E}(\theta) - \mathcal{E}(H_\theta)| \geq \varepsilon\big) \leq 2\mathcal{C}^{(B,\delta), \frac{\varepsilon}{8LR}}\exp\left(\frac{-\varepsilon^2 M}{2\mathcal{R}^4}\right) \tag{9.184}$$

(cf. Definition 9.2.6). Moreover, note that Proposition 9.2.25 (applied with $V \curvearrowleft \mathbb{R}^{\mathfrak{d}}$, $\|\cdot\| \curvearrowleft (\mathbb{R}^{\mathfrak{d}} \ni x \mapsto \|x\|_\infty \in [0, \infty))$, $r \curvearrowleft \frac{\varepsilon}{8LR}$, $R \curvearrowleft R$, $X \curvearrowleft B$, $d \curvearrowleft \delta$ in the notation of Proposition 9.2.25) demonstrates that

$$\mathcal{C}^{(B,\delta), \frac{\varepsilon}{8LR}} \leq \max\left\{1, \left[\frac{32LR\mathcal{R}}{\varepsilon}\right]^{\mathfrak{d}}\right\}. \tag{9.185}$$

This and (9.184) prove (9.182). The proof of Lemma 9.4.9 is thus complete. $\qquad\square$

**Lemma 9.4.10.** *Let $\mathfrak{d}, M, L \in \mathbb{N}$, $u \in \mathbb{R}$, $v \in (u,\infty)$, $R \in [1,\infty)$, $\varepsilon, b \in (0,\infty)$, $l = (l_0, l_1, \ldots, l_L) \in \mathbb{N}^{L+1}$ satisfy $l_L = 1$ and $\sum_{k=1}^{L} l_k(l_{k-1}+1) \leq \mathfrak{d}$, let $D \subseteq [-b,b]^{l_0}$ be a compact set, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $X_m \colon \Omega \to D$, $m \in \{1, 2, \ldots, M\}$, and $Y_m \colon \Omega \to [u,v]$, $m \in \{1, 2, \ldots, M\}$, be functions, assume that $(X_m, Y_m)$, $m \in \{1, 2, \ldots, M\}$, are i.i.d. random variables, let $\mathcal{E} \colon C(D, \mathbb{R}) \to [0, \infty)$ satisfy for all $f \in C(D, \mathbb{R})$ that $\mathcal{E}(f) = \mathbb{E}[|f(X_1) - Y_1|^2]$, and let $\mathfrak{E} \colon [-R, R]^{\mathfrak{d}} \times \Omega \to [0, \infty)$ satisfy for all $\theta \in [-R, R]^{\mathfrak{d}}$, $\omega \in \Omega$ that*

$$\mathfrak{E}(\theta, \omega) = \frac{1}{M}\left[\sum_{m=1}^{M} |\mathcal{N}_{u,v}^{\theta,l}(X_m(\omega)) - Y_m(\omega)|^2\right] \tag{9.186}$$

*(cf. Definitions 2.1.27 and 3.1.16). Then $\Omega \ni \omega \mapsto \sup_{\theta \in [-R,R]^{\mathfrak{d}}} |\mathfrak{E}(\theta, \omega) - \mathcal{E}(\mathcal{N}_{u,v}^{\theta,l}|_D)| \in [0,\infty]$ is an $\mathcal{F}/\mathcal{B}([0,\infty])$-measurable function and*

$$\mathbb{P}\big(\sup_{\theta \in [-R,R]^{\mathfrak{d}}}|\mathfrak{E}(\theta) - \mathcal{E}(\mathcal{N}_{u,v}^{\theta,l}|_D)| \geq \varepsilon\big)$$
$$\leq 2\max\left\{1, \left[\frac{32L\max\{1,b\}(\|l\|_\infty + 1)^L R^L(v-u)}{\varepsilon}\right]^{\mathfrak{d}}\right\}\exp\left(\frac{-\varepsilon^2 M}{2(v-u)^4}\right). \tag{9.187}$$

*Proof of Lemma 9.4.10.* Throughout this proof let $\mathfrak{L} \in (0, \infty)$ satisfy

$$\mathfrak{L} = L\max\{1,b\}(\|l\|_\infty + 1)^L R^{L-1}. \tag{9.188}$$

Observe that Corollary 5.3.7 (applied with $a \curvearrowleft -b$, $b \curvearrowleft b$, $u \curvearrowleft u$, $v \curvearrowleft v$, $d \curvearrowleft \mathfrak{d}$, $L \curvearrowleft L$, $l \curvearrowleft l$ in the notation of Corollary 5.3.7) and the assumption that $D \subseteq [-b,b]^{l_0}$ show that for all $\theta, \vartheta \in [-R, R]^{\mathfrak{d}}$ it holds that

$$\sup_{x \in D} |\mathcal{N}_{u,v}^{\theta,l}(x) - \mathcal{N}_{u,v}^{\vartheta,l}(x)| \leq \sup_{x \in [-b,b]^{l_0}} |\mathcal{N}_{u,v}^{\theta,l}(x) - \mathcal{N}_{u,v}^{\vartheta,l}(x)|$$
$$\leq L\max\{1,b\}(\|l\|_\infty + 1)^L (\max\{1, \|\theta\|_\infty, \|\vartheta\|_\infty\})^{L-1}\|\theta - \vartheta\|_\infty$$
$$\leq L\max\{1,b\}(\|l\|_\infty + 1)^L R^{L-1}\|\theta - \vartheta\|_\infty = \mathfrak{L}\|\theta - \vartheta\|_\infty. \tag{9.189}$$

Furthermore, observe that the fact that for all $\theta \in \mathbb{R}^{\mathfrak{d}}$, $x \in \mathbb{R}^{l_0}$ it holds that $\mathcal{N}_{u,v}^{\theta,l}(x) \in [u,v]$ and the assumption that for all $m \in \{1, 2, \ldots, M\}$, $\omega \in \Omega$ it holds that $Y_m(\omega) \in [u,v]$ demonstrate that for all $\theta \in [-R, R]^{\mathfrak{d}}$, $m \in \{1, 2, \ldots, M\}$ it holds that

$$|\mathcal{N}_{u,v}^{\theta,l}(X_m) - Y_m| \leq v - u. \tag{9.190}$$

Combining this and (9.189) with Lemma 9.4.9 (applied with $d \curvearrowleft l_0$, $\mathfrak{d} \curvearrowleft \mathfrak{d}$, $M \curvearrowleft M$, $R \curvearrowleft R$, $L \curvearrowleft \mathfrak{L}$, $\mathcal{R} \curvearrowleft v - u$, $\varepsilon \curvearrowleft \varepsilon$, $D \curvearrowleft D$, $(\Omega, \mathcal{F}, \mathbb{P}) \curvearrowleft (\Omega, \mathcal{F}, \mathbb{P})$, $(X_m)_{m \in \{1,2,\ldots,M\}} \curvearrowleft (X_m)_{m \in \{1,2,\ldots,M\}}$, $(Y_m)_{m \in \{1,2,\ldots,M\}} \curvearrowleft ((\Omega \ni \omega \mapsto Y_m(\omega) \in \mathbb{R}))_{m \in \{1,2,\ldots,M\}}$, $H \curvearrowleft ([-R,R]^{\mathfrak{d}} \ni \theta \mapsto \mathcal{N}_{u,v}^{\theta,l}|_D \in C(D, \mathbb{R}))$, $\mathcal{E} \curvearrowleft \mathcal{E}$, $\mathfrak{E} \curvearrowleft \mathfrak{E}$ in the notation of Lemma 9.4.9) establishes that $\Omega \ni \omega \mapsto \sup_{\theta \in [-R,R]^{\mathfrak{d}}} |\mathfrak{E}(\theta, \omega) - \mathcal{E}(\mathcal{N}_{u,v}^{\theta,l}|_D)| \in [0,\infty]$ is an $\mathcal{F}/\mathcal{B}([0,\infty])$-measurable function and

$$\mathbb{P}\big(\sup_{\theta \in [-R,R]^{\mathfrak{d}}}|\mathfrak{E}(\theta) - \mathcal{E}(\mathcal{N}_{u,v}^{\theta,l}|_D)| \geq \varepsilon\big) \leq 2\max\left\{1, \left[\frac{32\mathfrak{L}R(v-u)}{\varepsilon}\right]^{\mathfrak{d}}\right\}\exp\left(\frac{-\varepsilon^2 M}{2(v-u)^4}\right). \tag{9.191}$$

The proof of Lemma 9.4.10 is thus complete. $\square$

# Chapter 10

# Analysis of the optimization error

## 10.1 Convergence rates for the minimum Monte Carlo method

**Lemma 10.1.1.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $\mathfrak{d}, N \in \mathbb{N}$, let $\|\|\cdot\|\| \colon \mathbb{R}^{\mathfrak{d}} \to [0, \infty)$ be a norm, let $\mathfrak{H} \subseteq \mathbb{R}^{\mathfrak{d}}$ be a set, let $\vartheta \in \mathfrak{H}$, $L, \varepsilon \in (0, \infty)$, let $\mathfrak{E} \colon \mathfrak{H} \times \Omega \to \mathbb{R}$ be a $(\mathcal{B}(\mathfrak{H}) \otimes \mathcal{F})/\mathcal{B}(\mathbb{R})$-measurable function, assume for all $x, y \in \mathfrak{H}$, $\omega \in \Omega$ that $|\mathfrak{E}(x, \omega) - \mathfrak{E}(y, \omega)| \leq L\|\|x - y\|\|$, and let $\Theta_n \colon \Omega \to \mathfrak{H}$, $n \in \{1, 2, \ldots, N\}$, be i.i.d. random variables. Then*

$$\mathbb{P}\big(\big[\min_{n \in \{1,2,\ldots,N\}} \mathfrak{E}(\Theta_n)\big] - \mathfrak{E}(\vartheta) > \varepsilon\big) \leq \big[\mathbb{P}\big(\|\|\Theta_1 - \vartheta\|\| > \tfrac{\varepsilon}{L}\big)\big]^N \leq \exp\big(-N\,\mathbb{P}\big(\|\|\Theta_1 - \vartheta\|\| \leq \tfrac{\varepsilon}{L}\big)\big). \tag{10.1}$$

*Proof of Lemma 10.1.1.* Note that the assumption that for all $x, y \in \mathfrak{H}$, $\omega \in \Omega$ it holds that $|\mathfrak{E}(x, \omega) - \mathfrak{E}(y, \omega)| \leq L\|\|x - y\|\|$ implies that

$$\begin{aligned}
\big[\min_{n \in \{1,2,\ldots,N\}} \mathfrak{E}(\Theta_n)\big] - \mathfrak{E}(\vartheta) &= \min_{n \in \{1,2,\ldots,N\}}[\mathfrak{E}(\Theta_n) - \mathfrak{E}(\vartheta)] \\
&\leq \min_{n \in \{1,2,\ldots,N\}}|\mathfrak{E}(\Theta_n) - \mathfrak{E}(\vartheta)| \leq \min_{n \in \{1,2,\ldots,N\}}\big[L\|\|\Theta_n - \vartheta\|\|\big] \\
&= L\big[\min_{n \in \{1,2,\ldots,N\}}\|\|\Theta_n - \vartheta\|\|\big].
\end{aligned} \tag{10.2}$$

The assumption that $\Theta_n$, $n \in \{1, 2, \ldots, N\}$, are i.i.d. random variables and the fact that $\forall\, x \in \mathbb{R} \colon 1 - x \leq e^{-x}$ hence show that

$$\begin{aligned}
\mathbb{P}\Big(\big[\min_{n \in \{1,2,\ldots,N\}} \mathfrak{E}(\Theta_n)\big] - \mathfrak{E}(\vartheta) > \varepsilon\Big) &\leq \mathbb{P}\Big(L\big[\min_{n \in \{1,2,\ldots,N\}}\|\|\Theta_n - \vartheta\|\|\big] > \varepsilon\Big) \\
&= \mathbb{P}\big(\min_{n \in \{1,2,\ldots,N\}}\|\|\Theta_n - \vartheta\|\| > \tfrac{\varepsilon}{L}\big) = \big[\mathbb{P}\big(\|\|\Theta_1 - \vartheta\|\| > \tfrac{\varepsilon}{L}\big)\big]^N \\
&= \big[1 - \mathbb{P}\big(\|\|\Theta_1 - \vartheta\|\| \leq \tfrac{\varepsilon}{L}\big)\big]^N \leq \exp\big(-N\,\mathbb{P}\big(\|\|\Theta_1 - \vartheta\|\| \leq \tfrac{\varepsilon}{L}\big)\big).
\end{aligned} \tag{10.3}$$

The proof of Lemma 10.1.1 is thus complete. $\qquad\square$

## 10.2 Continuous uniformly distributed samples

**Lemma 10.2.1.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $\mathfrak{d}, N \in \mathbb{N}$, $a \in \mathbb{R}$, $b \in (a, \infty)$, $\vartheta \in [a, b]^{\mathfrak{d}}$, $L, \varepsilon \in (0, \infty)$, let $\mathfrak{E} \colon [a, b]^{\mathfrak{d}} \times \Omega \to \mathbb{R}$ be a $(\mathcal{B}([a, b]^{\mathfrak{d}}) \otimes \mathcal{F})/\mathcal{B}(\mathbb{R})$-measurable function, assume for all $x, y \in [a, b]^{\mathfrak{d}}$, $\omega \in \Omega$ that $|\mathfrak{E}(x, \omega) - \mathfrak{E}(y, \omega)| \leq L\|x - y\|_{\infty}$, let*

$\Theta_n \colon \Omega \to [a, b]^{\mathfrak{d}}$, $n \in \{1, 2, \ldots, N\}$, *be i.i.d. random variables, and assume that* $\Theta_1$ *is continuous uniformly distributed on* $[a, b]^{\mathfrak{d}}$ *(cf. Definition 3.1.16). Then*

$$\mathbb{P}\Big( \big[\min_{n \in \{1, 2, \ldots, N\}} \mathfrak{E}(\Theta_n)\big] - \mathfrak{E}(\vartheta) > \varepsilon \Big) \leq \exp\Big( -N \min\Big\{ 1, \frac{\varepsilon^{\mathfrak{d}}}{L^{\mathfrak{d}}(b-a)^{\mathfrak{d}}} \Big\} \Big). \qquad (10.4)$$

*Proof of Lemma 10.2.1.* Note that the assumption that $\Theta_1$ is continuous uniformly distributed on $[a, b]^{\mathfrak{d}}$ ensures that

$$\mathbb{P}\big( \|\Theta_1 - \vartheta\|_\infty \leq \tfrac{\varepsilon}{L} \big) \geq \mathbb{P}\big( \|\Theta_1 - (a, a, \ldots, a)\|_\infty \leq \tfrac{\varepsilon}{L} \big) = \mathbb{P}\big( \|\Theta_1 - (a, a, \ldots, a)\|_\infty \leq \min\{\tfrac{\varepsilon}{L}, b-a\} \big)$$

$$= \left[ \frac{\min\{\tfrac{\varepsilon}{L}, b-a\}}{(b-a)} \right]^{\mathfrak{d}} = \min\left\{ 1, \left[ \frac{\varepsilon}{L(b-a)} \right]^{\mathfrak{d}} \right\}.$$

$$(10.5)$$

Combining this with Lemma 10.1.1 proves (10.4). The proof of Lemma 10.2.1 is thus complete. $\qquad \square$

# Chapter 11

# Full error analysis for training algorithms for DNNs

## 11.1 Overall error decomposition

**Lemma 11.1.1.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $d, M \in \mathbb{N}$, let $D \subseteq \mathbb{R}^d$ be a compact set, let $X_m \colon \Omega \to D$, $m \in \{1, 2, \ldots, M\}$, and $Y_m \colon \Omega \to \mathbb{R}$, $m \in \{1, 2, \ldots, M\}$, be functions, assume that $(X_m, Y_m)$, $m \in \{1, 2, \ldots, M\}$, are i.i.d. random variables, assume $\mathbb{E}[|Y_1|^2] < \infty$, let $\mathcal{E} \colon C(D, \mathbb{R}) \to [0, \infty)$ satisfy for all $f \in C(D, \mathbb{R})$ that $\mathcal{E}(f) = \mathbb{E}[|f(X_1) - Y_1|^2]$, and let $\mathfrak{E} \colon C(D, \mathbb{R}) \times \Omega \to [0, \infty)$ satisfy for all $f \in C(D, \mathbb{R})$, $\omega \in \Omega$ that*

$$\mathfrak{E}(f, \omega) = \frac{1}{M}\left[\sum_{m=1}^{M} |f(X_m(\omega)) - Y_m(\omega)|^2\right]. \tag{11.1}$$

*Then it holds for all $f, \phi \in C(D, \mathbb{R})$ that*

$$
\begin{aligned}
\mathbb{E}\big[|f(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] &= \mathbb{E}\big[|\phi(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] + \mathcal{E}(f) - \mathcal{E}(\phi) \\
&\leq \mathbb{E}\big[|\phi(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] + \big[\mathfrak{E}(f) - \mathfrak{E}(\phi)\big] + 2\Big[\max_{v \in \{f, \phi\}} |\mathfrak{E}(v) - \mathcal{E}(v)|\Big].
\end{aligned}
\tag{11.2}
$$

*Proof of Lemma 11.1.1.* Note that Lemma 9.3.1 ensures that for all $f, \phi \in C(D, \mathbb{R})$ it holds that

$$
\begin{aligned}
&\mathbb{E}\big[|f(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] \\
&= \mathbb{E}\big[|\phi(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] + \mathcal{E}(f) - \mathcal{E}(\phi) \\
&= \mathbb{E}\big[|\phi(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] + \mathcal{E}(f) - \mathfrak{E}(f) + \mathfrak{E}(f) - \mathfrak{E}(\phi) + \mathfrak{E}(\phi) - \mathcal{E}(\phi) \\
&= \mathbb{E}\big[|\phi(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] + \big[\big(\mathcal{E}(f) - \mathfrak{E}(f)\big) + \big(\mathfrak{E}(\phi) - \mathcal{E}(\phi)\big)\big] + \big[\mathfrak{E}(f) - \mathfrak{E}(\phi)\big] \\
&\leq \mathbb{E}\big[|\phi(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] + \Big[\sum_{v \in \{f, \phi\}} |\mathfrak{E}(v) - \mathcal{E}(v)|\Big] + \big[\mathfrak{E}(f) - \mathfrak{E}(\phi)\big] \\
&\leq \mathbb{E}\big[|\phi(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] + 2\Big[\max_{v \in \{f, \phi\}} |\mathfrak{E}(v) - \mathcal{E}(v)|\Big] + \big[\mathfrak{E}(f) - \mathfrak{E}(\phi)\big].
\end{aligned}
\tag{11.3}
$$

The proof of Lemma 11.1.1 is thus complete. $\qquad\square$

**Lemma 11.1.2.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $d, \mathfrak{d}, M \in \mathbb{N}$, let $D \subseteq \mathbb{R}^d$ be a compact set, let $B \subseteq \mathbb{R}^{\mathfrak{d}}$ be a set, let $H = (H_\theta)_{\theta \in B} \colon B \to C(D, \mathbb{R})$ be a function, let*

$X_m \colon \Omega \to D$, $m \in \{1, 2, \ldots, M\}$, and $Y_m \colon \Omega \to \mathbb{R}$, $m \in \{1, 2, \ldots, M\}$, be functions, assume that $(X_m, Y_m)$, $m \in \{1, 2, \ldots, M\}$, are i.i.d. random variables, assume $\mathbb{E}[|Y_1|^2] < \infty$, let $\varphi \colon D \to \mathbb{R}$ be a $\mathcal{B}(D)/\mathcal{B}(\mathbb{R})$-measurable function, assume it holds $\mathbb{P}$-a.s. that $\varphi(X_1) = \mathbb{E}[Y_1|X_1]$, let $\mathcal{E} \colon C(D, \mathbb{R}) \to [0, \infty)$ satisfy for all $f \in C(D, \mathbb{R})$ that $\mathcal{E}(f) = \mathbb{E}[|f(X_1) - Y_1|^2]$, and let $\mathfrak{E} \colon B \times \Omega \to [0, \infty)$ satisfy for all $\theta \in B$, $\omega \in \Omega$ that

$$\mathfrak{E}(\theta, \omega) = \frac{1}{M} \left[ \sum_{m=1}^{M} |H_\theta(X_m(\omega)) - Y_m(\omega)|^2 \right]. \tag{11.4}$$

*Then it holds for all* $\theta, \vartheta \in B$ *that*

$$\int_D |H_\theta(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) = \int_D |H_\vartheta(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) + \mathcal{E}(H_\theta) - \mathcal{E}(H_\vartheta)$$
$$\leq \int_D |H_\vartheta(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) + \left[ \mathfrak{E}(\theta) - \mathfrak{E}(\vartheta) \right] + 2 \left[ \sup_{\eta \in B} |\mathfrak{E}(\eta) - \mathcal{E}(H_\eta)| \right]. \tag{11.5}$$

*Proof of Lemma 11.1.2.* First, observe that Lemma 11.1.1 (applied with $(\Omega, \mathcal{F}, \mathbb{P}) \curvearrowright (\Omega, \mathcal{F}, \mathbb{P})$, $d \curvearrowright d$, $M \curvearrowright M$, $D \curvearrowright D$, $(X_m)_{m \in \{1,2,\ldots,M\}} \curvearrowright (X_m)_{m \in \{1,2,\ldots,M\}}$, $(Y_m)_{m \in \{1,2,\ldots,M\}} \curvearrowright (Y_m)_{m \in \{1,2,\ldots,M\}}$, $\mathcal{E} \curvearrowright \mathcal{E}$, $\mathfrak{E} \curvearrowright \left( C(D, \mathbb{R}) \times \Omega \ni (f, \omega) \mapsto \frac{1}{M} \left[ \sum_{m=1}^{M} |f(X_m(\omega)) - Y_m(\omega)|^2 \right] \in [0, \infty) \right)$ in the notation of Lemma 11.1.1) shows that for all $\theta, \vartheta \in B$ it holds that

$$\mathbb{E}\big[|H_\theta(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] = \mathbb{E}\big[|H_\vartheta(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] + \mathcal{E}(H_\theta) - \mathcal{E}(H_\vartheta)$$
$$\leq \mathbb{E}\big[|H_\vartheta(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] + \left[ \mathfrak{E}(\theta) - \mathfrak{E}(\vartheta) \right] + 2 \left[ \max_{\eta \in \{\theta, \vartheta\}} |\mathfrak{E}(\eta) - \mathcal{E}(H_\eta)| \right] \tag{11.6}$$
$$\leq \mathbb{E}\big[|H_\vartheta(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] + \left[ \mathfrak{E}(\theta) - \mathfrak{E}(\vartheta) \right] + 2 \left[ \sup_{\eta \in B} |\mathfrak{E}(\eta) - \mathcal{E}(H_\eta)| \right].$$

In addition, note that the assumption that it holds $\mathbb{P}$-a.s. that $\varphi(X_1) = \mathbb{E}[Y_1|X_1]$ ensures that for all $\eta \in B$ it holds that

$$\mathbb{E}\big[|H_\eta(X_1) - \mathbb{E}[Y_1|X_1]|^2\big] = \mathbb{E}\big[|H_\eta(X_1) - \varphi(X_1)|^2\big] = \int_D |H_\eta(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x). \tag{11.7}$$

Combining this with (11.6) establishes (11.5). The proof of Lemma 11.1.2 is thus complete. $\qquad \square$

## 11.2 Analysis of the convergence speed

### 11.2.1 Convergence rates for convergence in probability

**Lemma 11.2.1.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $u \in \mathbb{R}$, $v \in (u, \infty)$, $\mathfrak{d}, L \in \mathbb{N}$, let $l = (l_0, l_1, \ldots, l_L) \in \mathbb{N}^{L+1}$ satisfy $l_L = 1$ and $\sum_{i=1}^{L} l_i(l_{i-1} + 1) \leq \mathfrak{d}$, let $B \subseteq \mathbb{R}^{\mathfrak{d}}$ be a non-empty compact set, and let $X \colon \Omega \to \mathbb{R}^{l_0}$ and $Y \colon \Omega \to [u, v]$ be random variables. Then*

*(i) it holds for all $\theta \in B$, $\omega \in \Omega$ that $|\mathscr{N}_{u,v}^{\theta,l}(X(\omega)) - Y(\omega)|^2 \in [0, (v - u)^2]$,*

*(ii) it holds that $B \ni \theta \mapsto \mathbb{E}\big[|\mathscr{N}_{u,v}^{\theta,l}(X) - Y|^2\big] \in [0, \infty)$ is continuous, and*

*(iii) there exists $\vartheta \in B$ such that $\mathbb{E}\big[|\mathcal{N}_{u,v}^{\vartheta,l}(X) - Y|^2\big] = \inf_{\theta \in B} \mathbb{E}\big[|\mathcal{N}_{u,v}^{\theta,l}(X) - Y|^2\big]$*

*(cf. Definition 2.1.27).*

*Proof of Lemma 11.2.1.* First, note that the fact that for all $\theta \in \mathbb{R}^{\mathfrak{d}}$, $x \in \mathbb{R}^{l_0}$ it holds that $\mathcal{N}_{u,v}^{\theta,l}(x) \in [u,v]$ and the assumption that for all $\omega \in \Omega$ it holds that $Y(\omega) \in [u,v]$ demonstrate item (i). Next observe that Corollary 5.3.7 ensures that for all $\omega \in \Omega$ it holds that $B \ni \theta \mapsto |\mathcal{N}_{u,v}^{\theta,l}(X(\omega)) - Y(\omega)|^2 \in [0,\infty)$ is a continuous function. Combining this and item (i) with Lebesgue's dominated convergence theorem establishes item (ii). Furthermore, note that item (ii) and the assumption that $B \subseteq \mathbb{R}^{\mathfrak{d}}$ is a non-empty compact set prove item (iii). The proof of Lemma 11.2.1 is thus complete. □

**Theorem 11.2.2.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $d, \mathfrak{d}, K, M \in \mathbb{N}$, $\varepsilon \in (0,\infty)$, $L, u \in \mathbb{R}$, $v \in (u,\infty)$, let $D \subseteq \mathbb{R}^d$ be a compact set, assume $|D| \geq 2$, let $X_m \colon \Omega \to D$, $m \in \{1,2,\ldots,M\}$, and $Y_m \colon \Omega \to [u,v]$, $m \in \{1,2,\ldots,M\}$, be functions, assume that $(X_m, Y_m)$, $m \in \{1,2,\ldots,M\}$, are i.i.d. random variables, let $\delta \colon D \times D \to [0,\infty)$ satisfy for all $x = (x_1, x_2, \ldots, x_d), y = (y_1, y_2, \ldots, y_d) \in D$ that $\delta(x,y) = \sum_{i=1}^{d} |x_i - y_i|$, let $\varphi \colon D \to [u,v]$ satisfy $\mathbb{P}$-a.s. that $\varphi(X_1) = \mathbb{E}[Y_1|X_1]$, assume for all $x, y \in D$ that $|\varphi(x) - \varphi(y)| \leq L\delta(x,y)$, let $N \in \mathbb{N} \cap [\max\{2, \mathcal{C}^{(D,\delta),\frac{\varepsilon}{4L}}\}, \infty)$, let $l \in \mathbb{N} \cap (N, \infty)$, let $\mathfrak{l} = (\mathfrak{l}_0, \mathfrak{l}_1, \ldots, \mathfrak{l}_l) \in \mathbb{N}^{l+1}$ satisfy for all $i \in \mathbb{N} \cap [2, N]$, $j \in \mathbb{N} \cap [N, l]$ that $\mathfrak{l}_0 = d$, $\mathfrak{l}_1 \geq 2dN$, $\mathfrak{l}_i \geq 2N - 2i + 3$, $\mathfrak{l}_j \geq 2$, $\mathfrak{l}_l = 1$, and $\sum_{k=1}^{l} \mathfrak{l}_k(\mathfrak{l}_{k-1} + 1) \leq \mathfrak{d}$, let $R \in [\max\{1, L, \sup_{z \in D} \|z\|_\infty, 2[\sup_{z \in D} |\varphi(z)|]\}, \infty)$, let $B \subseteq \mathbb{R}^{\mathfrak{d}}$ satisfy $B = [-R, R]^{\mathfrak{d}}$, let $\mathfrak{E} \colon B \times \Omega \to [0,\infty)$ satisfy for all $\theta \in B$, $\omega \in \Omega$ that*

$$\mathfrak{E}(\theta, \omega) = \frac{1}{M}\left[\sum_{m=1}^{M} |\mathcal{N}_{u,v}^{\theta,l}(X_m(\omega)) - Y_m(\omega)|^2\right], \tag{11.8}$$

*let $\Theta_k \colon \Omega \to B$, $k \in \{1,2,\ldots,K\}$, be i.i.d. random variables, assume that $\Theta_1$ is continuous uniformly distributed on $B$, and let $\Xi \colon \Omega \to B$ satisfy $\Xi = \Theta_{\min\{k \in \{1,2,\ldots,K\} \colon \mathfrak{E}(\Theta_k) = \min_{l \in \{1,2,\ldots,K\}} \mathfrak{E}(\Theta_l)\}}$ (cf. Definitions 2.1.27, 3.1.16, and 9.2.6). Then*

$$\mathbb{P}\left(\int_D |\mathcal{N}_{u,v}^{\Xi,l}(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(dx) > \varepsilon^2\right) \leq \exp\left(-K\min\left\{1, \frac{\varepsilon^{2\mathfrak{d}}}{(16(v-u)l(\|\mathfrak{l}\|_\infty + 1)^l R^{l+1})^{\mathfrak{d}}}\right\}\right)$$
$$+ 2\exp\left(\mathfrak{d}\ln\left(\max\left\{1, \frac{128l(\|\mathfrak{l}\|_\infty + 1)^l R^{l+1}(v-u)}{\varepsilon^2}\right\}\right) - \frac{\varepsilon^4 M}{32(v-u)^4}\right). \tag{11.9}$$

*Proof of Theorem 11.2.2.* Throughout this proof let $\mathcal{M} \subseteq D$ satisfy $|\mathcal{M}| = \max\{2, \mathcal{C}^{(D,\delta),\frac{\varepsilon}{4L}}\}$ and

$$4L\left[\sup_{x \in D}\left(\inf_{y \in \mathcal{M}} \delta(x,y)\right)\right] \leq \varepsilon, \tag{11.10}$$

let $b \in [0,\infty)$ satisfy $b = \sup_{z \in D} \|z\|_\infty$, let $\mathcal{E} \colon C(D, \mathbb{R}) \to [0,\infty)$ satisfy for all $f \in C(D, \mathbb{R})$ that $\mathcal{E}(f) = \mathbb{E}[|f(X_1) - Y_1|^2]$, and let $\vartheta \in B$ satisfy $\mathcal{E}(\mathcal{N}_{u,v}^{\vartheta,l}|_D) = \inf_{\theta \in B} \mathcal{E}(\mathcal{N}_{u,v}^{\theta,l}|_D)$ (cf. Lemma 11.2.1). Observe that the assumption that for all $x, y \in D$ it holds that $|\varphi(x) - \varphi(y)| \leq L\delta(x,y)$ implies that $\varphi$ is a $\mathcal{B}(D)/\mathcal{B}([u,v])$-measurable function. Lemma 11.1.2 (applied with $(\Omega, \mathcal{F}, \mathbb{P}) \curvearrowleft (\Omega, \mathcal{F}, \mathbb{P})$, $d \curvearrowleft d$, $\mathfrak{d} \curvearrowleft \mathfrak{d}$, $M \curvearrowleft M$, $D \curvearrowleft D$, $B \curvearrowleft B$, $H \curvearrowleft (B \ni \theta \mapsto \mathcal{N}_{u,v}^{\theta,l}|_D \in C(D, \mathbb{R}))$, $(X_m)_{m \in \{1,2,\ldots,M\}} \curvearrowleft (X_m)_{m \in \{1,2,\ldots,M\}}$, $(Y_m)_{m \in \{1,2,\ldots,M\}} \curvearrowleft ((\Omega \ni \omega \mapsto Y_m(\omega) \in \mathbb{R}))_{m \in \{1,2,\ldots,M\}}$, $\varphi \curvearrowleft (D \ni x \mapsto \varphi(x) \in \mathbb{R})$, $\mathcal{E} \curvearrowleft \mathcal{E}$, $\mathfrak{E} \curvearrowleft \mathfrak{E}$ in the

notation of Lemma 11.1.2) therefore ensures that for all $\omega \in \Omega$ it holds that

$$
\int_D |\mathscr{N}_{u,v}^{\Xi(\omega),\mathfrak{l}}(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x)
$$
$$
\leq \underbrace{\int_D |\mathscr{N}_{u,v}^{\vartheta,\mathfrak{l}}(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x)}_{\text{Approximation error}} + \underbrace{\left[\mathfrak{E}(\Xi(\omega),\omega) - \mathfrak{E}(\vartheta,\omega)\right]}_{\text{Optimization error}} + \underbrace{2\left[\sup_{\theta \in B} |\mathfrak{E}(\theta,\omega) - \mathcal{E}(\mathscr{N}_{u,v}^{\theta,\mathfrak{l}}|_D)|\right]}_{\text{Generalization error}}.
$$
$$(11.11)$$

Next observe that the assumption that $N \geq \max\{2, \mathcal{C}^{(D,\delta),\frac{\varepsilon}{4L}}\} = |\mathcal{M}|$ shows that for all $i \in \mathbb{N} \cap [2, N]$ it holds that $l \geq |\mathcal{M}| + 1$, $\mathfrak{l}_1 \geq 2d|\mathcal{M}|$ and $\mathfrak{l}_i \geq 2|\mathcal{M}| - 2i + 3$. The assumption that for all $x, y \in D$ it holds that $|\varphi(x) - \varphi(y)| \leq L\delta(x,y)$, the assumption that $R \geq \max\{1, L, \sup_{z \in D} \|z\|_\infty, 2[\sup_{z \in D} |\varphi(z)|]\}$, **??** (applied with $d \curvearrowleft d$, $\mathfrak{d} \curvearrowleft \mathfrak{d}$, $\mathfrak{L} \curvearrowleft l$, $L \curvearrowleft L$, $u \curvearrowleft u$, $v \curvearrowleft v$, $D \curvearrowleft D$, $f \curvearrowleft \varphi$, $\mathcal{M} \curvearrowleft \mathcal{M}$, $l \curvearrowleft \mathfrak{l}$ in the notation of **??**), and (11.10) hence ensure that there exists $\eta \in B$ which satisfies

$$
\sup_{x \in D} |\mathscr{N}_{u,v}^{\eta,\mathfrak{l}}(x) - \varphi(x)| \leq 2L \left[ \sup_{x=(x_1,x_2,\ldots,x_d) \in D} \left( \inf_{y=(y_1,y_2,\ldots,y_d) \in \mathcal{M}} \sum_{i=1}^{d} |x_i - y_i| \right) \right]
$$
$$
= 2L \left[ \sup_{x \in D} \left( \inf_{y \in \mathcal{M}} \delta(x,y) \right) \right] \leq \frac{\varepsilon}{2}.
$$
$$(11.12)$$

Lemma 11.1.2 (applied with $(\Omega, \mathcal{F}, \mathbb{P}) \curvearrowleft (\Omega, \mathcal{F}, \mathbb{P})$, $d \curvearrowleft d$, $\mathfrak{d} \curvearrowleft \mathfrak{d}$, $M \curvearrowleft M$, $D \curvearrowleft D$, $B \curvearrowleft B$, $H \curvearrowleft (B \ni \theta \mapsto \mathscr{N}_{u,v}^{\theta,\mathfrak{l}}|_D \in C(D, \mathbb{R}))$, $(X_m)_{m \in \{1,2,\ldots,M\}} \curvearrowleft (X_m)_{m \in \{1,2,\ldots,M\}}$, $(Y_m)_{m \in \{1,2,\ldots,M\}} \curvearrowleft ((\Omega \ni \omega \mapsto Y_m(\omega) \in \mathbb{R}))_{m \in \{1,2,\ldots,M\}}$, $\varphi \curvearrowleft (D \ni x \mapsto \varphi(x) \in \mathbb{R})$, $\mathcal{E} \curvearrowleft \mathcal{E}$, $\mathfrak{E} \curvearrowleft \mathfrak{E}$ in the notation of Lemma 11.1.2) and the assumption that $\mathcal{E}(\mathscr{N}_{u,v}^{\vartheta,\mathfrak{l}}|_D) = \inf_{\theta \in B} \mathcal{E}(\mathscr{N}_{u,v}^{\theta,\mathfrak{l}}|_D)$ therefore prove that

$$
\int_D |\mathscr{N}_{u,v}^{\vartheta,\mathfrak{l}}(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) = \int_D |\mathscr{N}_{u,v}^{\eta,\mathfrak{l}}(x) - \varphi|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) + \underbrace{\mathcal{E}(\mathscr{N}_{u,v}^{\vartheta,\mathfrak{l}}|_D) - \mathcal{E}(\mathscr{N}_{u,v}^{\eta,\mathfrak{l}}|_D)}_{\leq 0}
$$
$$
\leq \int_D |\mathscr{N}_{u,v}^{\eta,\mathfrak{l}}(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) \leq \sup_{x \in D} |\mathscr{N}_{u,v}^{\eta,\mathfrak{l}}(x) - \varphi(x)|^2 \leq \frac{\varepsilon^2}{4}.
$$
$$(11.13)$$

Combining this with (11.11) shows that for all $\omega \in \Omega$ it holds that

$$
\int_D |\mathscr{N}_{u,v}^{\Xi(\omega),\mathfrak{l}}(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) \leq \frac{\varepsilon^2}{4} + \left[\mathfrak{E}(\Xi(\omega),\omega) - \mathfrak{E}(\vartheta,\omega)\right] + 2\left[\sup_{\theta \in B} |\mathfrak{E}(\theta,\omega) - \mathcal{E}(\mathscr{N}_{u,v}^{\theta,\mathfrak{l}}|_D)|\right].
$$
$$(11.14)$$

Hence, we obtain that

$$
\mathbb{P}\left( \int_D |\mathscr{N}_{u,v}^{\Xi,\mathfrak{l}}(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) > \varepsilon^2 \right) \leq \mathbb{P}\left( \left[\mathfrak{E}(\Xi) - \mathfrak{E}(\vartheta)\right] + 2\left[\sup_{\theta \in B} |\mathfrak{E}(\theta) - \mathcal{E}(\mathscr{N}_{u,v}^{\theta,\mathfrak{l}}|_D)|\right] > \frac{3\varepsilon^2}{4} \right)
$$
$$
\leq \mathbb{P}\left( \mathfrak{E}(\Xi) - \mathfrak{E}(\vartheta) > \frac{\varepsilon^2}{4} \right) + \mathbb{P}\left( \sup_{\theta \in B} |\mathfrak{E}(\theta) - \mathcal{E}(\mathscr{N}_{u,v}^{\theta,\mathfrak{l}}|_D)| > \frac{\varepsilon^2}{4} \right).
$$
$$(11.15)$$

Next observe that Corollary 5.3.7 (applied with $a \curvearrowleft -b$, $b \curvearrowleft b$, $u \curvearrowleft u$, $v \curvearrowleft v$, $d \curvearrowleft \mathfrak{d}$, $L \curvearrowleft l$, $l \curvearrowleft \mathfrak{l}$ in the notation of Corollary 5.3.7) demonstrates that for all $\theta, \xi \in B$ it holds

that

$$\sup_{x\in D}|\mathcal{N}_{u,v}^{\theta,\mathfrak{l}}(x) - \mathcal{N}_{u,v}^{\xi,\mathfrak{l}}(x)| \leq \sup_{x\in[-b,b]^d}|\mathcal{N}_{u,v}^{\theta,\mathfrak{l}}(x) - \mathcal{N}_{u,v}^{\xi,\mathfrak{l}}(x)|$$

$$\leq l\max\{1,b\}(\|\mathfrak{l}\|_\infty + 1)^l(\max\{1,\|\theta\|_\infty,\|\xi\|_\infty\})^{l-1}\|\theta - \xi\|_\infty$$

$$\leq lR(\|\mathfrak{l}\|_\infty + 1)^l R^{l-1}\|\theta - \xi\|_\infty = l(\|\mathfrak{l}\|_\infty + 1)^l R^l\|\theta - \xi\|_\infty. \tag{11.16}$$

Combining this with the fact that for all $\theta \in \mathbb{R}^\mathfrak{d}$, $x \in D$ it holds that $\mathcal{N}_{u,v}^{\theta,\mathfrak{l}}(x) \in [u,v]$, the assumption that for all $m \in \{1,2,\ldots,M\}$ $\omega \in \Omega$ it holds that $Y_m(\omega) \in [u,v]$, the fact that for all $x_1,x_2,y \in \mathbb{R}$ it holds that $(x_1-y)^2 - (x_2-y)^2 = (x_1-x_2)((x_1-y)+(x_2-y))$, and (11.8) ensures that for all $\theta,\xi \in B$, $\omega \in \Omega$ it holds that

$$|\mathfrak{E}(\theta,\omega) - \mathfrak{E}(\xi,\omega)|$$

$$= \left|\frac{1}{M}\left[\sum_{m=1}^{M}|\mathcal{N}_{u,v}^{\theta,\mathfrak{l}}(X_m(\omega)) - Y_m(\omega)|^2\right] - \frac{1}{M}\left[\sum_{m=1}^{M}|\mathcal{N}_{u,v}^{\xi,\mathfrak{l}}(X_m(\omega)) - Y_m(\omega)|^2\right]\right| \tag{11.17}$$

$$= \frac{1}{M}\left|\sum_{m=1}^{M}\left(\left(\mathcal{N}_{u,v}^{\theta,\mathfrak{l}}(X_m(\omega)) - \mathcal{N}_{u,v}^{\xi,\mathfrak{l}}(X_m(\omega))\right)\left[\left(\mathcal{N}_{u,v}^{\theta,\mathfrak{l}}(X_m(\omega)) - Y_m(\omega)\right) + \left(\mathcal{N}_{u,v}^{\xi,\mathfrak{l}}(X_m(\omega)) - Y_m(\omega)\right)\right]\right)\right|$$

$$\leq \frac{1}{M}\left[\sum_{m=1}^{M}\left(|\mathcal{N}_{u,v}^{\theta,\mathfrak{l}}(X_m(\omega)) - \mathcal{N}_{u,v}^{\xi,\mathfrak{l}}(X_m(\omega))|\underbrace{\left[|\mathcal{N}_{u,v}^{\theta,\mathfrak{l}}(X_m(\omega)) - Y_m(\omega)| + |\mathcal{N}_{u,v}^{\xi,\mathfrak{l}}(X_m(\omega)) - Y_m(\omega)|\right]}_{\leq 2(v-u)}\right)\right]$$

$$\leq 2(v-u)l(\|\mathfrak{l}\|_\infty + 1)^l R^l\|\theta - \xi\|_\infty.$$

Lemma 10.2.1 (applied with $(\Omega,\mathcal{F},\mathbb{P}) \curvearrowright (\Omega,\mathcal{F},\mathbb{P})$, $\mathfrak{d} \curvearrowright \mathfrak{d}$, $N \curvearrowright K$, $a \curvearrowright -R$, $b \curvearrowright R$, $\vartheta \curvearrowright \vartheta$, $L \curvearrowright 2(v-u)l(\|\mathfrak{l}\|_\infty + 1)^l R^l$, $\varepsilon \curvearrowright \frac{\varepsilon^2}{4}$, $\mathfrak{E} \curvearrowright \mathfrak{E}$, $(\Theta_n)_{n\in\{1,2,\ldots,N\}} \curvearrowright (\Theta_k)_{k\in\{1,2,\ldots,K\}}$ in the notation of Lemma 10.2.1) therefore shows that

$$\mathbb{P}\left(\mathfrak{E}(\Xi) - \mathfrak{E}(\vartheta) > \frac{\varepsilon^2}{4}\right) = \mathbb{P}\left(\left[\min_{k\in\{1,2,\ldots,K\}}\mathfrak{E}(\Theta_k)\right] - \mathfrak{E}(\vartheta) > \frac{\varepsilon^2}{4}\right)$$

$$\leq \exp\left(-K\min\left\{1, \frac{\left(\frac{\varepsilon^2}{4}\right)^\mathfrak{d}}{[2(v-u)l(\|\mathfrak{l}\|_\infty + 1)^l R^l]^\mathfrak{d}(2R)^\mathfrak{d}}\right\}\right) \tag{11.18}$$

$$= \exp\left(-K\min\left\{1, \frac{\varepsilon^{2\mathfrak{d}}}{(16(v-u)l(\|\mathfrak{l}\|_\infty + 1)^l R^{l+1})^\mathfrak{d}}\right\}\right).$$

Moreover, note that Lemma 9.4.10 (applied with $\mathfrak{d} \curvearrowright \mathfrak{d}$, $M \curvearrowright M$, $L \curvearrowright l$, $u \curvearrowright u$, $v \curvearrowright v$, $R \curvearrowright R$, $\varepsilon \curvearrowright \frac{\varepsilon^2}{4}$, $b \curvearrowright b$, $l \curvearrowright \mathfrak{l}$, $D \curvearrowright D$, $(\Omega,\mathcal{F},\mathbb{P}) \curvearrowright (\Omega,\mathcal{F},\mathbb{P})$, $(X_m)_{m\in\{1,2,\ldots,M\}} \curvearrowright (X_m)_{m\in\{1,2,\ldots,M\}}$, $(Y_m)_{m\in\{1,2,\ldots,M\}} \curvearrowright (Y_m)_{m\in\{1,2,\ldots,M\}}$, $\mathcal{E} \curvearrowright \mathcal{E}$, $\mathfrak{E} \curvearrowright \mathfrak{E}$ in the notation of Lemma 9.4.10) establishes that

$$\mathbb{P}\left(\sup_{\theta\in B}|\mathfrak{E}(\theta) - \mathcal{E}(\mathcal{N}_{u,v}^{\theta,\mathfrak{l}}|_D)| \geq \frac{\varepsilon^2}{4}\right)$$

$$\leq 2\max\left\{1, \left[\frac{128 l\max\{1,b\}(\|\mathfrak{l}\|_\infty + 1)^l R^l(v-u)}{\varepsilon^2}\right]^\mathfrak{d}\right\}\exp\left(\frac{-\varepsilon^4 M}{32(v-u)^4}\right)$$

$$\leq 2\max\left\{1, \left[\frac{128 l(\|\mathfrak{l}\|_\infty + 1)^l R^{l+1}(v-u)}{\varepsilon^2}\right]^\mathfrak{d}\right\}\exp\left(\frac{-\varepsilon^4 M}{32(v-u)^4}\right) \tag{11.19}$$

$$= 2\exp\left(\mathfrak{d}\ln\left(\max\left\{1, \frac{128 l(\|\mathfrak{l}\|_\infty + 1)^l R^{l+1}(v-u)}{\varepsilon^2}\right\}\right) - \frac{\varepsilon^4 M}{32(v-u)^4}\right).$$

Combining this and (11.18) with (11.15) proves that

$$\mathbb{P}\left(\int_D |\mathscr{N}_{u,v}^{\Xi,\mathfrak{l}}(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x) > \varepsilon^2\right) \leq \exp\left(-K\min\left\{1, \frac{\varepsilon^{2\mathfrak{d}}}{(16(v-u)l(\|\mathfrak{l}\|_\infty + 1)^l R^{l+1})^{\mathfrak{d}}}\right\}\right)$$
$$+ 2\exp\left(\mathfrak{d}\ln\left(\max\left\{1, \frac{128l(\|\mathfrak{l}\|_\infty + 1)^l R^{l+1}(v-u)}{\varepsilon^2}\right\}\right) - \frac{\varepsilon^4 M}{32(v-u)^4}\right). \quad (11.20)$$

The proof of Theorem 11.2.2 is thus complete. $\qquad\square$

**Corollary 11.2.3.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $d, \mathfrak{d}, K, M, \tau \in \mathbb{N}$, $\varepsilon \in (0, \infty)$, $L, a, u \in \mathbb{R}$, $b \in (a, \infty)$, $v \in (u, \infty)$, $R \in [\max\{1, L, |a|, |b|, 2|u|, 2|v|\}, \infty)$, let $X_m \colon \Omega \to [a,b]^d$, $m \in \{1, 2, \ldots, M\}$, be i.i.d. random variables, let $\varphi \colon [a,b]^d \to [u,v]$ satisfy for all $x, y \in [a,b]^d$ that $|\varphi(x) - \varphi(y)| \leq L\|x - y\|_2$, assume $\tau \geq 2d(2dL(b-a)\varepsilon^{-1} + 2)^d$ and $\mathfrak{d} \geq \tau(d+1) + (\tau - 3)\tau(\tau + 1) + \tau + 1$, let $\mathfrak{l} \in \mathbb{N}^\tau$ satisfy $\mathfrak{l} = (d, \tau, \tau, \ldots, \tau, 1)$, let $B \subseteq \mathbb{R}^{\mathfrak{d}}$ satisfy $B = [-R, R]^{\mathfrak{d}}$, let $\mathfrak{E} \colon B \times \Omega \to [0, \infty)$ satisfy for all $\theta \in B$, $\omega \in \Omega$ that*

$$\mathfrak{E}(\theta, \omega) = \frac{1}{M}\left[\sum_{m=1}^M |\mathscr{N}_{u,v}^{\theta,\mathfrak{l}}(X_m(\omega)) - \varphi(X_m(\omega))|^2\right], \quad (11.21)$$

*let $\Theta_k \colon \Omega \to B$, $k \in \{1, 2, \ldots, K\}$, be i.i.d. random variables, assume that $\Theta_1$ is continuous uniformly distributed on $B$, and let $\Xi \colon \Omega \to B$ satisfy $\Xi = \Theta_{\min\{k \in \{1,2,\ldots,K\} \colon \mathfrak{E}(\Theta_k) = \min_{l \in \{1,2,\ldots,K\}} \mathfrak{E}(\Theta_l)\}}$ (cf. Definition 2.1.27). Then*

$$\mathbb{P}\left(\left[\int_{[a,b]^d} |\mathscr{N}_{u,v}^{\Xi,\mathfrak{l}}(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x)\right]^{1/2} > \varepsilon\right) \leq \exp\left(-K\min\left\{1, \frac{\varepsilon^{2\mathfrak{d}}}{(16(v-u)(\tau+1)^\tau R^\tau)^{\mathfrak{d}}}\right\}\right)$$
$$+ 2\exp\left(\mathfrak{d}\ln\left(\max\left\{1, \frac{128(\tau+1)^\tau R^\tau(v-u)}{\varepsilon^2}\right\}\right) - \frac{\varepsilon^4 M}{32(v-u)^4}\right). \quad (11.22)$$

*Proof of Corollary 11.2.3.* Throughout this proof let $N \in \mathbb{N}$ satisfy

$$N = \min\left\{k \in \mathbb{N} \colon k \geq \frac{2dL(b-a)}{\varepsilon}\right\}, \quad (11.23)$$

let $\mathcal{M} \subseteq [a,b]^d$ satisfy $\mathcal{M} = \{a, a + \frac{b-a}{N}, \ldots, a + \frac{(N-1)(b-a)}{N}, b\}^d$, let $\delta \colon [a,b]^d \times [a,b]^d \to [0,\infty)$ satisfy for all $x = (x_1, x_2, \ldots, x_d), y = (y_1, y_2, \ldots, y_d) \in [a,b]^d$ that $\delta(x,y) = \sum_{i=1}^d |x_i - y_i|$, and let $l_0, l_1, \ldots, l_{\tau-1} \in \mathbb{N}$ satisfy $\mathfrak{l} = (l_0, l_1, \ldots, l_{\tau-1})$. Observe that for all $x \in [a,b]$ there exists $y \in \{a, a + \frac{b-a}{N}, \ldots, a + \frac{(N-1)(b-a)}{N}, b\}$ such that $|x - y| \leq \frac{b-a}{2N}$. This demonstrates that

$$4L\left[\sup_{x=(x_1,x_2,\ldots,x_d)\in[a,b]^d}\left(\inf_{y=(y_1,y_2,\ldots,y_d)\in\mathcal{M}} \sum_{i=1}^d |x_i - y_i|\right)\right] \leq \frac{2Ld(b-a)}{N} \leq \varepsilon. \quad (11.24)$$

Hence, we obtain that

$$\mathcal{C}^{([a,b]^d,\delta),\frac{\varepsilon}{4L}} \leq |\mathcal{M}| = (N+1)^d. \quad (11.25)$$

Next note that (11.23) implies that $N < 2dL(b-a)\varepsilon^{-1} + 1$. The assumption that $\tau \geq 2d(2dL(b-a)\varepsilon^{-1} + 2)^d$ therefore ensures that

$$\tau > 2d(N+1)^d \geq (N+1)^d + 2. \quad (11.26)$$

Hence, we obtain that for all $i \in \{2, 3, \ldots, (N+1)^d\}$, $j \in \{(N+1)^d + 1, (N+1)^d + 2, \ldots, \tau - 2\}$ it holds that

$$l_0 = d, \quad l_1 = \tau \geq 2d(N+1)^d, \quad l_{\tau-1} = 1, \quad l_i = \tau \geq 2(N+1)^d - 2i + 3, \quad \text{and} \quad l_j = \tau \geq 2. \tag{11.27}$$

Furthermore, observe that the assumption that for all $x, y \in [a,b]^d$ it holds that $|\varphi(x) - \varphi(y)| \leq L\|x - y\|_2$ implies that for all $x, y \in [a,b]^d$ it holds that $|\varphi(x) - \varphi(y)| \leq L\delta(x, y)$. Combining this, (11.25), (11.26), (11.27), and the assumption that $\mathfrak{d} \geq \tau(d+1) + (\tau - 3)\tau(\tau + 1) + \tau + 1 = \sum_{i=1}^{\tau-1} l_i(l_{i-1} + 1)$ with Theorem 11.2.2 (applied with $(\Omega, \mathcal{F}, \mathbb{P}) \curvearrowright (\Omega, \mathcal{F}, \mathbb{P})$, $d \curvearrowright d$, $\mathfrak{d} \curvearrowright \mathfrak{d}$, $K \curvearrowright K$, $M \curvearrowright M$, $\varepsilon \curvearrowright \varepsilon$, $L \curvearrowright L$, $u \curvearrowright u$, $v \curvearrowright v$, $D \curvearrowright [a,b]^d$, $(X_m)_{m \in \{1,2,\ldots,M\}} \curvearrowright (X_m)_{m \in \{1,2,\ldots,M\}}$, $(Y_m)_{m \in \{1,2,\ldots,M\}} \curvearrowright (\varphi(X_m))_{m \in \{1,2,\ldots,M\}}$, $\delta \curvearrowright \delta$, $\varphi \curvearrowright \varphi$, $N \curvearrowright (N+1)^d$, $l \curvearrowright \tau - 1$, $\mathfrak{l} \curvearrowright \mathfrak{l}$, $R \curvearrowright R$, $B \curvearrowright B$, $\mathfrak{E} \curvearrowright \mathfrak{E}$, $(\Theta_k)_{k \in \{1,2,\ldots,K\}} \curvearrowright (\Theta_k)_{k \in \{1,2,\ldots,K\}}$, $\Xi \curvearrowright \Xi$ in the notation of Theorem 11.2.2) establishes that

$$\mathbb{P}\left(\left[\int_{[a,b]^d} |\mathcal{N}_{u,v}^{\Xi,\mathfrak{l}}(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x)\right]^{1/2} > \varepsilon\right)$$
$$\leq \exp\left(-K \min\left\{1, \frac{\varepsilon^{2\mathfrak{d}}}{(16(v-u)(\tau-1)(\tau+1)^{\tau-1}R^\tau)^{\mathfrak{d}}}\right\}\right)$$
$$+ 2\exp\left(\mathfrak{d}\ln\left(\max\left\{1, \frac{128(\tau-1)(\tau+1)^{\tau-1}R^\tau(v-u)}{\varepsilon^2}\right\}\right) - \frac{\varepsilon^4 M}{32(v-u)^4}\right)$$
$$\leq \exp\left(-K \min\left\{1, \frac{\varepsilon^{2\mathfrak{d}}}{(16(v-u)(\tau+1)^\tau R^\tau)^{\mathfrak{d}}}\right\}\right)$$
$$+ 2\exp\left(\mathfrak{d}\ln\left(\max\left\{1, \frac{128(\tau+1)^\tau R^\tau(v-u)}{\varepsilon^2}\right\}\right) - \frac{\varepsilon^4 M}{32(v-u)^4}\right). \tag{11.28}$$

The proof of Corollary 11.2.3 is thus complete. $\qquad\square$

**Corollary 11.2.4.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $d \in \mathbb{N}$, $L, a, u \in \mathbb{R}$, $b \in (a, \infty)$, $v \in (u, \infty)$, $R \in [\max\{1, L, |a|, |b|, 2|u|, 2|v|\}, \infty)$, let $X_m \colon \Omega \to [a,b]^d$, $m \in \mathbb{N}$, be i.i.d. random variables, let $\varphi \colon [a,b]^d \to [u,v]$ satisfy for all $x, y \in [a,b]^d$ that $|\varphi(x) - \varphi(y)| \leq L\|x - y\|_2$, let $\mathfrak{l}_\tau \in \mathbb{N}^\tau$, $\tau \in \mathbb{N}$, satisfy for all $\tau \in \mathbb{N} \cap [3, \infty)$ that $\mathfrak{l}_\tau = (d, \tau, \tau, \ldots, \tau, 1)$, let $\mathfrak{E}_{\mathfrak{d},M,\tau} \colon [-R,R]^{\mathfrak{d}} \times \Omega \to [0, \infty)$, $\mathfrak{d}, M, \tau \in \mathbb{N}$, satisfy for all $\mathfrak{d}, M \in \mathbb{N}$, $\tau \in \mathbb{N} \cap [3, \infty)$, $\theta \in [-R,R]^{\mathfrak{d}}$, $\omega \in \Omega$ with $\mathfrak{d} \geq \tau(d+1) + (\tau-3)\tau(\tau+1) + \tau + 1$ that*

$$\mathfrak{E}_{\mathfrak{d},M,\tau}(\theta, \omega) = \frac{1}{M}\left[\sum_{m=1}^M |\mathcal{N}_{u,v}^{\theta,\mathfrak{l}_\tau}(X_m(\omega)) - \varphi(X_m(\omega))|^2\right], \tag{11.29}$$

*for every $\mathfrak{d} \in \mathbb{N}$ let $\Theta_{\mathfrak{d},k} \colon \Omega \to [-R,R]^{\mathfrak{d}}$, $k \in \mathbb{N}$, be i.i.d. random variables, assume for all $\mathfrak{d} \in \mathbb{N}$ that $\Theta_{\mathfrak{d},1}$ is continuous uniformly distributed on $[-R,R]^{\mathfrak{d}}$, and let $\Xi_{\mathfrak{d},K,M,\tau} \colon \Omega \to [-R,R]^{\mathfrak{d}}$, $\mathfrak{d}, K, M, \tau \in \mathbb{N}$, satisfy for all $\mathfrak{d}, K, M, \tau \in \mathbb{N}$ that $\Xi_{\mathfrak{d},K,M,\tau} = \Theta_{\mathfrak{d},\min\{k \in \{1,2,\ldots,K\}: \mathfrak{E}_{\mathfrak{d},M,\tau}(\Theta_{\mathfrak{d},k}) = \min_{l \in \{}\}}$ (cf. Definition 2.1.27). Then there exists $c \in (0, \infty)$ such that for all $\mathfrak{d}, K, M, \tau \in \mathbb{N}$, $\varepsilon \in (0, \sqrt{v - u}]$ with $\tau \geq 2d(2dL(b-a)\varepsilon^{-1} + 2)^d$ and $\mathfrak{d} \geq \tau(d+1) + (\tau-3)\tau(\tau+1) + \tau + 1$ it holds that*

$$\mathbb{P}\left(\left[\int_{[a,b]^d} |\mathcal{N}_{u,v}^{\Xi_{\mathfrak{d},K,M,\tau},\mathfrak{l}_\tau}(x) - \varphi(x)|^2 \, \mathbb{P}_{X_1}(\mathrm{d}x)\right]^{1/2} > \varepsilon\right)$$
$$\leq \exp\left(-K(c\tau)^{-\tau\mathfrak{d}}\varepsilon^{2\mathfrak{d}}\right) + 2\exp\left(\mathfrak{d}\ln\left((c\tau)^\tau \varepsilon^{-2}\right) - c^{-1}\varepsilon^4 M\right). \tag{11.30}$$

*Proof of Corollary 11.2.4.* Throughout this proof let $c \in (0, \infty)$ satisfy

$$c = \max\{32(v-u)^4, 256(v-u+1)R\}. \tag{11.31}$$

Note that Corollary 11.2.3 establishes that for all $\mathfrak{d}, K, M, \tau \in \mathbb{N}$, $\varepsilon \in (0, \infty)$ with $\tau \geq 2d(2dL(b-a)\varepsilon^{-1}+2)^d$ and $\mathfrak{d} \geq \tau(d+1)+(\tau-3)\tau(\tau+1)+\tau+1$ it holds that

$$\mathbb{P}\left(\left[\int_{[a,b]^d}|\mathcal{N}_{u,v}^{\Xi_{\mathfrak{d}},K,M,\tau,\mathfrak{l}_\tau}(x)-\varphi(x)|^2\,\mathbb{P}_{X_1}(\mathrm{d}x)\right]^{1/2}>\varepsilon\right) \leq \exp\left(-K\min\left\{1,\frac{\varepsilon^{2\mathfrak{d}}}{(16(v-u)(\tau+1)^\tau R^\tau)^{\mathfrak{d}}}\right\}\right)$$
$$+2\exp\left(\mathfrak{d}\ln\left(\max\left\{1,\frac{128(\tau+1)^\tau R^\tau(v-u)}{\varepsilon^2}\right\}\right)-\frac{\varepsilon^4 M}{32(v-u)^4}\right). \tag{11.32}$$

Next observe that (11.31) ensures that for all $\tau \in \mathbb{N}$ it holds that

$$16(v-u)(\tau+1)^\tau R^\tau \leq (16(v-u+1)(\tau+1)R)^\tau \leq (32(v-u+1)R\tau)^\tau \leq (c\tau)^\tau. \tag{11.33}$$

The fact that for all $\varepsilon \in (0, \sqrt{v-u}]$, $\tau \in \mathbb{N}$ it holds that $\varepsilon^2 \leq 16(v-u)(\tau+1)^\tau R^\tau$ therefore shows that for all $\varepsilon \in (0, \sqrt{v-u}]$, $\tau \in \mathbb{N}$ it holds that

$$-\min\left\{1,\frac{\varepsilon^{2\mathfrak{d}}}{(16(v-u)(\tau+1)^\tau R^\tau)^{\mathfrak{d}}}\right\} = \frac{-\varepsilon^{2\mathfrak{d}}}{(16(v-u)(\tau+1)^\tau R^\tau)^{\mathfrak{d}}} \leq \frac{-\varepsilon^{2\mathfrak{d}}}{(c\tau)^{\tau\mathfrak{d}}}. \tag{11.34}$$

Furthermore, note that (11.31) implies that for all $\tau \in \mathbb{N}$ it holds that

$$128(\tau+1)^\tau R^\tau(v-u) \leq 128(2\tau)^\tau R^\tau(v-u) \leq (256R\tau(v-u+1))^\tau \leq (c\tau)^\tau. \tag{11.35}$$

The fact that for all $\varepsilon \in (0, \sqrt{v-u}]$, $\tau \in \mathbb{N}$ it holds that $\varepsilon^2 \leq 128(\tau+1)^\tau R^\tau(v-u)$ hence proves that for all $\varepsilon \in (0, \sqrt{v-u}]$, $\tau \in \mathbb{N}$ it holds that

$$\ln\left(\max\left\{1,\frac{128(\tau+1)^\tau R^\tau(v-u)}{\varepsilon^2}\right\}\right) = \ln\left(\frac{128(\tau+1)^\tau R^\tau(v-u)}{\varepsilon^2}\right) \leq \ln\left(\frac{(c\tau)^\tau}{\varepsilon^2}\right) \tag{11.36}$$

In addition, observe that (11.31) ensures that

$$\frac{-1}{32(v-u)^4} \leq \frac{-1}{c}. \tag{11.37}$$

Combining this, (11.34), and (11.36) with (11.32) proves that for all $\mathfrak{d}, K, M, \tau \in \mathbb{N}$, $\varepsilon \in (0, \sqrt{v-u}]$ with $\tau \geq 2d(2dL(b-a)\varepsilon^{-1}+2)^d$ and $\mathfrak{d} \geq \tau(d+1)+(\tau-3)\tau(\tau+1)+\tau+1$ it holds that

$$\mathbb{P}\left(\left[\int_{[u,v]^d}|\mathcal{N}_{u,v}^{\Xi_{\mathfrak{d}},K,M,\tau,\mathfrak{l}_\tau}(x)-\varphi(x)|^2\,\mathbb{P}_{X_1}(\mathrm{d}x)\right]^{1/2}>\varepsilon\right)$$
$$\leq \exp\left(\frac{-K\varepsilon^{2\mathfrak{d}}}{(c\tau)^{\tau\mathfrak{d}}}\right)+2\exp\left(\mathfrak{d}\ln\left(\frac{(c\tau)^\tau}{\varepsilon^2}\right)-\frac{\varepsilon^4 M}{c}\right). \tag{11.38}$$

The proof of Corollary 11.2.4 is thus complete. $\qquad\square$

**Corollary 11.2.5.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $d \in \mathbb{N}$, $L, a, u \in \mathbb{R}$, $b \in (a, \infty)$, $v \in (u, \infty)$, $R \in [\max\{1, L, |a|, |b|, 2|u|, 2|v|\}, \infty)$, let $X_m \colon \Omega \to [a, b]^d$, $m \in \mathbb{N}$, be i.i.d. random variables, let $\varphi \colon [a, b]^d \to [u, v]$ satisfy for all $x, y \in [a, b]^d$ that $|\varphi(x) - \varphi(y)| \leq$*

$L\|x - y\|_2$, let $\mathfrak{l}_\tau \in \mathbb{N}^\tau$, $\tau \in \mathbb{N}$, satisfy for all $\tau \in \mathbb{N} \cap [3, \infty)$ that $\mathfrak{l}_\tau = (d, \tau, \tau, \ldots, \tau, 1)$, let $\mathfrak{E}_{\mathfrak{d},M,\tau} \colon [-R, R]^{\mathfrak{d}} \times \Omega \to [0, \infty)$, $\mathfrak{d}, M, \tau \in \mathbb{N}$, satisfy for all $\mathfrak{d}, M \in \mathbb{N}$, $\tau \in \mathbb{N} \cap [3, \infty)$, $\theta \in [-R, R]^{\mathfrak{d}}$, $\omega \in \Omega$ with $\mathfrak{d} \geq \tau(d+1) + (\tau-3)\tau(\tau+1) + \tau + 1$ that

$$\mathfrak{E}_{\mathfrak{d},M,\tau}(\theta, \omega) = \frac{1}{M}\left[\sum_{m=1}^{M}|\mathscr{N}_{u,v}^{\theta,\mathfrak{l}_\tau}(X_m(\omega)) - \varphi(X_m(\omega))|^2\right], \tag{11.39}$$

*for every $\mathfrak{d} \in \mathbb{N}$ let $\Theta_{\mathfrak{d},k} \colon \Omega \to [-R, R]^{\mathfrak{d}}$, $k \in \mathbb{N}$, be i.i.d. random variables, assume for all $\mathfrak{d} \in \mathbb{N}$ that $\Theta_{\mathfrak{d},1}$ is continuous uniformly distributed on $[-R, R]^{\mathfrak{d}}$, and let $\Xi_{\mathfrak{d},K,M,\tau} \colon \Omega \to [-R, R]^{\mathfrak{d}}$, $\mathfrak{d}, K, M, \tau \in \mathbb{N}$, satisfy for all $\mathfrak{d}, K, M, \tau \in \mathbb{N}$ that $\Xi_{\mathfrak{d},K,M,\tau} = \Theta_{\mathfrak{d},\min\{k \in \{1,2,\ldots,K\} \colon \mathfrak{E}_{\mathfrak{d},M,\tau}(\Theta_{\mathfrak{d},k})=\min_{l \in \{}}$ (cf. Definition 2.1.27). Then there exists $c \in (0, \infty)$ such that for all $\mathfrak{d}, K, M, \tau \in \mathbb{N}$, $\varepsilon \in (0, \sqrt{v-u}]$ with $\tau \geq 2d(2dL(b-a)\varepsilon^{-1} + 2)^d$ and $\mathfrak{d} \geq \tau(d+1) + (\tau-3)\tau(\tau+1) + \tau + 1$ it holds that*

$$\mathbb{P}\left(\int_{[a,b]^d}|\mathscr{N}_{u,v}^{\Xi_{\mathfrak{d},K,M,\tau},\mathfrak{l}_\tau}(x) - \varphi(x)|\,\mathbb{P}_{X_1}(dx) > \varepsilon\right)$$
$$\leq \exp\left(-K(c\tau)^{-\tau\mathfrak{d}}\varepsilon^{2\mathfrak{d}}\right) + 2\exp\left(\mathfrak{d}\ln\left((c\tau)^\tau\varepsilon^{-2}\right) - c^{-1}\varepsilon^4 M\right). \tag{11.40}$$

*Proof of Corollary 11.2.5.* Note that Jensen's inequality shows that for all $f \in C([a,b]^d, \mathbb{R})$ it holds that

$$\int_{[a,b]^d}|f(x)|\,\mathbb{P}_{X_1}(dx) \leq \left[\int_{[a,b]^d}|f(x)|^2\,\mathbb{P}_{X_1}(dx)\right]^{\frac{1}{2}}. \tag{11.41}$$

Combining this with Corollary 11.2.4 proves (11.40). The proof of Corollary 11.2.5 is thus complete. □

## 11.2.2 Convergence rates for strong convergence

**Lemma 11.2.6.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $c \in [0, \infty)$, and let $X \colon \Omega \to [-c, c]$ be a random variable. Then it holds for all $\varepsilon, p \in (0, \infty)$ that*

$$\mathbb{E}[|X|^p] \leq \varepsilon^p\,\mathbb{P}(|X| \leq \varepsilon) + c^p\,\mathbb{P}(|X| > \varepsilon) \leq \varepsilon^p + c^p\,\mathbb{P}(|X| > \varepsilon). \tag{11.42}$$

*Proof of Lemma 11.2.6.* Observe that the assumption that for all $\omega \in \Omega$ it holds that $|X(\omega)| \leq c$ ensures that for all $\varepsilon, p \in (0, \infty)$ it holds that

$$\mathbb{E}[|X|^p] = \mathbb{E}\left[|X|^p\mathbb{1}_{\{|X|\leq\varepsilon\}}\right] + \mathbb{E}\left[|X|^p\mathbb{1}_{\{|X|>\varepsilon\}}\right] \leq \varepsilon^p\,\mathbb{P}(|X| \leq \varepsilon) + c^p\,\mathbb{P}(|X| > \varepsilon) \leq \varepsilon^p + c^p\,\mathbb{P}(|X| > \varepsilon). \tag{11.43}$$

The proof of Lemma 11.2.6 is thus complete. □

**Corollary 11.2.7.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $d \in \mathbb{N}$, $L, a, u \in \mathbb{R}$, $b \in (a, \infty)$, $v \in (u, \infty)$, $R \in [\max\{1, L, |a|, |b|, 2|u|, 2|v|\}, \infty)$, let $X_m \colon \Omega \to [a, b]^d$, $m \in \mathbb{N}$, be i.i.d. random variables, let $\varphi \colon [a, b]^d \to [u, v]$ satisfy for all $x, y \in [a, b]^d$ that $|\varphi(x) - \varphi(y)| \leq L\|x - y\|_2$, let $\mathfrak{l}_\tau \in \mathbb{N}^\tau$, $\tau \in \mathbb{N}$, satisfy for all $\tau \in \mathbb{N} \cap [3, \infty)$ that $\mathfrak{l}_\tau = (d, \tau, \tau, \ldots, \tau, 1)$, let $\mathfrak{E}_{\mathfrak{d},M,\tau} \colon [-R, R]^{\mathfrak{d}} \times \Omega \to [0, \infty)$, $\mathfrak{d}, M, \tau \in \mathbb{N}$, satisfy for all $\mathfrak{d}, M \in \mathbb{N}$, $\tau \in \mathbb{N} \cap [3, \infty)$, $\theta \in [-R, R]^{\mathfrak{d}}$, $\omega \in \Omega$ with $\mathfrak{d} \geq \tau(d+1) + (\tau-3)\tau(\tau+1) + \tau + 1$ that*

$$\mathfrak{E}_{\mathfrak{d},M,\tau}(\theta, \omega) = \frac{1}{M}\left[\sum_{m=1}^{M}|\mathscr{N}_{u,v}^{\theta,\mathfrak{l}_\tau}(X_m(\omega)) - \varphi(X_m(\omega))|^2\right], \tag{11.44}$$

*for every $\eth \in \mathbb{N}$ let $\Theta_{\eth,k}\colon \Omega \to [-R,R]^{\eth}$, $k \in \mathbb{N}$, be i.i.d. random variables, assume for all $\eth \in \mathbb{N}$ that $\Theta_{\eth,1}$ is continuous uniformly distributed on $[-R,R]^{\eth}$, and let $\Xi_{\eth,K,M,\tau}\colon \Omega \to [-R,R]^{\eth}$, $\eth, K, M, \tau \in \mathbb{N}$, satisfy for all $\eth, K, M, \tau \in \mathbb{N}$ that $\Xi_{\eth,K,M,\tau} = \Theta_{\eth,\min\{k\in\{1,2,...,K\}\colon \mathfrak{E}_{\eth,M,\tau}(\Theta_{\eth,k})=\min_{l\in\{}}$ (cf. Definition 2.1.27). Then there exists $c \in (0,\infty)$ such that for all $\eth, K, M, \tau \in \mathbb{N}$, $p \in [1,\infty)$, $\varepsilon \in (0, \sqrt{v-u}]$ with $\tau \geq 2d(2dL(b-a)\varepsilon^{-1} + 2)^d$ and $\eth \geq \tau(d+1) + (\tau - 3)\tau(\tau+1) + \tau + 1$ it holds that*

$$\left(\mathbb{E}\left[\left(\int_{[a,b]^d}|\mathcal{N}_{u,v}^{\Xi_{\eth,K,M,\tau},\mathfrak{l}_\tau}(x) - \varphi(x)|^2\,\mathbb{P}_{X_1}(\mathrm{d}x)\right)^{p/2}\right]\right)^{1/p}$$
$$\leq (v-u)\left[\exp\left(-K(c\tau)^{-\tau\eth}\varepsilon^{2\eth}\right) + 2\exp\left(\eth\ln\left((c\tau)^\tau\varepsilon^{-2}\right) - c^{-1}\varepsilon^4 M\right)\right]^{1/p} + \varepsilon. \tag{11.45}$$

*Proof of Corollary 11.2.7.* First, observe that Corollary 11.2.4 ensures that there exists $c \in (0,\infty)$ which satisfies for all $\eth, K, M, \tau \in \mathbb{N}$, $\varepsilon \in (0, \sqrt{v-u}]$ with $\tau \geq 2d(2dL(b-a)\varepsilon^{-1} + 2)^d$ and $\eth \geq \tau(d+1) + (\tau - 3)\tau(\tau+1) + \tau + 1$ that

$$\mathbb{P}\left(\left[\int_{[a,b]^d}|\mathcal{N}_{u,v}^{\Xi_{\eth,K,M,\tau},\mathfrak{l}_\tau}(x) - \varphi(x)|^2\,\mathbb{P}_{X_1}(\mathrm{d}x)\right]^{1/2} > \varepsilon\right)$$
$$\leq \exp\left(-K(c\tau)^{-\tau\eth}\varepsilon^{2\eth}\right) + 2\exp\left(\eth\ln\left((c\tau)^\tau\varepsilon^{-2}\right) - c^{-1}\varepsilon^4 M\right). \tag{11.46}$$

Lemma 11.2.6 (applied with $(\Omega, \mathcal{F}, \mathbb{P}) \curvearrowright (\Omega, \mathcal{F}, \mathbb{P})$, $c \curvearrowright v - u$, $X \curvearrowright (\Omega \ni \omega \mapsto [\int_{[a,b]^d}|\mathcal{N}_{u,v}^{\Xi_{\eth,K,M,\tau}(\omega),\mathfrak{l}_\tau}(x) - \varphi(x)|^2\,\mathbb{P}_{X_1}(\mathrm{d}x)]^{1/2} \in [u-v, v-u])$ in the notation of Lemma 11.2.6) hence ensures that for all $\eth, K, M, \tau \in \mathbb{N}$, $\varepsilon \in (0, \sqrt{v-u}]$, $p \in (0,\infty)$ with $\tau \geq 2d(2dL(b-a)\varepsilon^{-1} + 2)^d$ and $\eth \geq \tau(d+1) + (\tau - 3)\tau(\tau+1) + \tau + 1$ it holds that

$$\mathbb{E}\left[\left(\int_{[a,b]^d}|\mathcal{N}_{u,v}^{\Xi_{\eth,K,M,\tau},\mathfrak{l}_\tau}(x) - \varphi(x)|^2\,\mathbb{P}_{X_1}(\mathrm{d}x)\right)^{p/2}\right]$$
$$\leq \varepsilon^p + (v-u)^p\left[\exp\left(-K(c\tau)^{-\tau\eth}\varepsilon^{2\eth}\right) + 2\exp\left(\eth\ln\left((c\tau)^\tau\varepsilon^{-2}\right) - c^{-1}\varepsilon^4 M\right)\right]. \tag{11.47}$$

The fact that for all $p \in [1,\infty)$, $x, y \in [0,\infty)$ it holds that $(x+y)^{1/p} \leq x^{1/p} + y^{1/p}$ therefore establishes (11.45). The proof of Corollary 11.2.7 is thus complete. $\qquad\square$

# Chapter 12

# Machine learning for partial differential equations (PDEs)

## 12.1 Linear heat PDEs

This section is a modified extract from the article Beck et al. [1].

### 12.1.1 Stochastic optimization problems for expectations of random variables

**Lemma 12.1.1.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and let $X \colon \Omega \to \mathbb{R}$ be a random variable which satisfies $\mathbb{E}[|X|^2] < \infty$. Then*

*(i) it holds for every $y \in \mathbb{R}$ that*

$$\mathbb{E}\big[|X - y|^2\big] = \mathbb{E}\big[|X - \mathbb{E}[X]|^2\big] + |\mathbb{E}[X] - y|^2, \tag{12.1}$$

*(ii) it holds that there exists a unique real number $z \in \mathbb{R}$ such that*

$$\mathbb{E}\big[|X - z|^2\big] = \inf_{y \in \mathbb{R}} \mathbb{E}\big[|X - y|^2\big], \tag{12.2}$$

*and*

*(iii) it holds that*

$$\mathbb{E}\big[|X - \mathbb{E}[X]|^2\big] = \inf_{y \in \mathbb{R}} \mathbb{E}\big[|X - y|^2\big]. \tag{12.3}$$

*Proof of Lemma 12.1.1.* Observe that the fact that $\mathbb{E}[|X|] < \infty$ ensures that for every $y \in \mathbb{R}$ it holds that

$$
\begin{aligned}
\mathbb{E}\big[|X - y|^2\big] &= \mathbb{E}\big[|X - \mathbb{E}[X] + \mathbb{E}[X] - y|^2\big] \\
&= \mathbb{E}\big[|X - \mathbb{E}[X]|^2 + 2(X - \mathbb{E}[X])(\mathbb{E}[X] - y) + |\mathbb{E}[X] - y|^2\big] \\
&= \mathbb{E}\big[|X - \mathbb{E}[X]|^2\big] + 2(\mathbb{E}[X] - y)\mathbb{E}\big[X - \mathbb{E}[X]\big] + |\mathbb{E}[X] - y|^2 \\
&= \mathbb{E}\big[|X - \mathbb{E}[X]|^2\big] + |\mathbb{E}[X] - y|^2.
\end{aligned}
\tag{12.4}
$$

This establishes item (i). Items (ii) and (iii) are immediate consequences of item (i). The proof of Lemma 12.1.1 is thus complete. $\qquad\square$

## 12.1.2 Stochastic optimization problems for expectations of random fields

**Proposition 12.1.2.** *Let $d \in \mathbb{N}$, $a \in \mathbb{R}$, $b \in (a, \infty)$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $X = (X_x)_{x \in [a,b]^d} \colon [a,b]^d \times \Omega \to \mathbb{R}$ be a $(\mathcal{B}([a,b]^d) \otimes \mathcal{F})/\mathcal{B}(\mathbb{R})$-measurable function, assume for every $x \in [a,b]^d$ that $\mathbb{E}[|X_x|^2] < \infty$, and assume that the function $[a,b]^d \ni x \mapsto \mathbb{E}[X_x] \in \mathbb{R}$ is continuous. Then*

*(i) it holds that there exists a unique continuous function $u \colon [a,b]^d \to \mathbb{R}$ such that*

$$\int_{[a,b]^d} \mathbb{E}\big[|X_x - u(x)|^2\big] \, \mathrm{d}x = \inf_{v \in C([a,b]^d, \mathbb{R})} \left( \int_{[a,b]^d} \mathbb{E}\big[|X_x - v(x)|^2\big] \, \mathrm{d}x \right) \qquad (12.5)$$

*and*

*(ii) it holds for every $x \in [a,b]^d$ that $u(x) = \mathbb{E}[X_x]$.*

*Proof of Proposition 12.1.2.* Note that item (i) in Lemma 12.1.1 and the assumption that $\forall x \in [a,b]^d \colon \mathbb{E}[|X_x|^2] < \infty$ ensure that for every function $u \colon [a,b]^d \to \mathbb{R}$ and every $x \in [a,b]^d$ it holds that

$$\mathbb{E}\big[|X_x - u(x)|^2\big] = \mathbb{E}\big[|X_x - \mathbb{E}[X_x]|^2\big] + |\mathbb{E}[X_x] - u(x)|^2. \qquad (12.6)$$

Fubini's theorem (see, e.g., Klenke [19, Theorem 14.16]) hence proves that for every continuous function $u \colon [a,b]^d \to \mathbb{R}$ it holds that

$$\int_{[a,b]^d} \mathbb{E}\big[|X_x - u(x)|^2\big] \, \mathrm{d}x = \int_{[a,b]^d} \mathbb{E}\big[|X_x - \mathbb{E}[X_x]|^2\big] \, \mathrm{d}x + \int_{[a,b]^d} |\mathbb{E}[X_x] - u(x)|^2 \, \mathrm{d}x. \quad (12.7)$$

The assumption that the function $[a,b]^d \ni x \mapsto \mathbb{E}[X_x] \in \mathbb{R}$ is continuous therefore demonstrates that

$$
\begin{aligned}
&\int_{[a,b]^d} \mathbb{E}\big[|X_x - \mathbb{E}[X_x]|^2\big] \, \mathrm{d}x \\
&\geq \inf_{v \in C([a,b]^d, \mathbb{R})} \left( \int_{[a,b]^d} \mathbb{E}\big[|X_x - v(x)|^2\big] \, \mathrm{d}x \right) \\
&= \inf_{v \in C([a,b]^d, \mathbb{R})} \left( \int_{[a,b]^d} \mathbb{E}\big[|X_x - \mathbb{E}[X_x]|^2\big] \, \mathrm{d}x + \int_{[a,b]^d} |\mathbb{E}[X_x] - v(x)|^2 \, \mathrm{d}x \right) \\
&\geq \inf_{v \in C([a,b]^d, \mathbb{R})} \left( \int_{[a,b]^d} \mathbb{E}\big[|X_x - \mathbb{E}[X_x]|^2\big] \, \mathrm{d}x \right) \\
&= \int_{[a,b]^d} \mathbb{E}\big[|X_x - \mathbb{E}[X_x]|^2\big] \, \mathrm{d}x.
\end{aligned}
\qquad (12.8)
$$

Hence, we obtain that

$$\int_{[a,b]^d} \mathbb{E}\big[|X_x - \mathbb{E}[X_x]|^2\big] \, \mathrm{d}x = \inf_{v \in C([a,b]^d, \mathbb{R})} \left( \int_{[a,b]^d} \mathbb{E}\big[|X_x - v(x)|^2\big] \, \mathrm{d}x \right). \qquad (12.9)$$

Again the fact that the function $[a,b]^d \ni x \mapsto \mathbb{E}[X_x] \in \mathbb{R}$ is continuous therefore proves that there exists a continuous function $u \colon [a,b]^d \to \mathbb{R}$ such that

$$\int_{[a,b]^d} \mathbb{E}\big[|X_x - u(x)|^2\big] \, \mathrm{d}x = \inf_{v \in C([a,b]^d, \mathbb{R})} \left( \int_{[a,b]^d} \mathbb{E}\big[|X_x - v(x)|^2\big] \, \mathrm{d}x \right). \qquad (12.10)$$

Next observe that (12.7) and (12.9) yield that for every continuous function $u\colon [a,b]^d \to \mathbb{R}$ with

$$\int_{[a,b]^d} \mathbb{E}\big[|X_x - u(x)|^2\big]\,\mathrm{d}x = \inf_{v \in C([a,b]^d, \mathbb{R})} \left(\int_{[a,b]^d} \mathbb{E}\big[|X_x - v(x)|^2\big]\,\mathrm{d}x\right) \tag{12.11}$$

it holds that

$$\begin{aligned}
&\int_{[a,b]^d} \mathbb{E}\big[|X_x - \mathbb{E}[X_x]|^2\big]\,\mathrm{d}x \\
&= \inf_{v \in C([a,b]^d, \mathbb{R})} \left(\int_{[a,b]^d} \mathbb{E}\big[|X_x - v(x)|^2\big]\,\mathrm{d}x\right) = \int_{[a,b]^d} \mathbb{E}\big[|X_x - u(x)|^2\big]\,\mathrm{d}x \\
&= \int_{[a,b]^d} \mathbb{E}\big[|X_x - \mathbb{E}[X_x]|^2\big]\,\mathrm{d}x + \int_{[a,b]^d} |\mathbb{E}[X_x] - u(x)|^2\,\mathrm{d}x.
\end{aligned} \tag{12.12}$$

Hence, we obtain that for every continuous function $u\colon [a,b]^d \to \mathbb{R}$ with

$$\int_{[a,b]^d} \mathbb{E}\big[|X_x - u(x)|^2\big]\,\mathrm{d}x = \inf_{v \in C([a,b]^d, \mathbb{R})} \left(\int_{[a,b]^d} \mathbb{E}\big[|X_x - v(x)|^2\big]\,\mathrm{d}x\right) \tag{12.13}$$

it holds that

$$\int_{[a,b]^d} |\mathbb{E}[X_x] - u(x)|^2\,\mathrm{d}x = 0. \tag{12.14}$$

This and the assumption that the function $[a,b]^d \ni x \mapsto \mathbb{E}[X_x] \in \mathbb{R}$ is continuous yield that for every continuous function $u\colon [a,b]^d \to \mathbb{R}$ with

$$\int_{[a,b]^d} \mathbb{E}\big[|X_x - u(x)|^2\big]\,\mathrm{d}x = \inf_{v \in C([a,b]^d, \mathbb{R})} \left(\int_{[a,b]^d} \mathbb{E}\big[|X_x - v(x)|^2\big]\,\mathrm{d}x\right) \tag{12.15}$$

and every $x \in [a,b]^d$ it holds that $u(x) = \mathbb{E}[X_x]$. Combining this with (12.10) completes the proof of Proposition 12.1.2. $\qquad\square$

### 12.1.3   Feynman–Kac formulas

#### 12.1.3.1   Feynman–Kac formulas providing existence of solutions

**Proposition 12.1.3.** *Let $T \in (0,\infty)$, $d, m \in \mathbb{N}$, $B \in \mathbb{R}^{d \times m}$, $\varphi \in C^2(\mathbb{R}^d, \mathbb{R})$ satisfy $\sup_{x \in \mathbb{R}^d}\big[\sum_{i,j=1}^{d}(|\varphi(x)| + |(\frac{\partial}{\partial x_i}\varphi)(x)| + |(\frac{\partial^2}{\partial x_i \partial x_j}\varphi)(x)|)\big] < \infty$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $Z\colon \Omega \to \mathbb{R}^m$ be a standard normal random variable, and let $u\colon [0,T] \times \mathbb{R}^d \to \mathbb{R}$ satisfy for all $t \in [0,T]$, $x \in \mathbb{R}^d$ that*

$$u(t,x) = \mathbb{E}\Big[\varphi(x + \sqrt{t}BZ)\Big]. \tag{12.16}$$

*Then*

   (i) *it holds that $u \in C^{1,2}([0,T] \times \mathbb{R}^d, \mathbb{R})$ and*

   (ii) *it holds for all $t \in [0,T]$, $x \in \mathbb{R}^d$ that*

$$(\tfrac{\partial u}{\partial t})(t,x) = \tfrac{1}{2}\operatorname{Trace}\big(BB^*(\operatorname{Hess}_x u)(t,x)\big) \tag{12.17}$$

*(cf. Definition 2.2.29).*

*Proof of Proposition 12.1.3.* Throughout this proof let $e_1 = (1, 0, \ldots, 0), e_2 = (0, 1, \ldots, 0), \ldots, e_m = (0, \ldots, 0, 1) \in \mathbb{R}^m$, let $\langle \cdot, \cdot \rangle \colon (\cup_{k \in \mathbb{N}} (\mathbb{R}^k \times \mathbb{R}^k)) \to \mathbb{R}$ satisfy for all $k \in \mathbb{N}$, $x = (x_1, x_2, \ldots, x_k)$, $y = (y_1, y_2, \ldots, y_k) \in \mathbb{R}^k$ that $\langle x, y \rangle = \sum_{i=1}^{k} x_i y_i$, and let $\psi_{t,x} = (\psi_{t,x}(y))_{y \in \mathbb{R}^m} \colon \mathbb{R}^m \to \mathbb{R}$, $t \in [0, T]$, $x \in \mathbb{R}^d$, satisfy for all $t \in [0, T]$, $x \in \mathbb{R}^d$, $y \in \mathbb{R}^m$ that $\psi_{t,x}(y) = \varphi(x + \sqrt{t} By)$. Note that the assumption that $\varphi \in C^2(\mathbb{R}^d, \mathbb{R})$, the assumption that $\sup_{x \in \mathbb{R}^d} [\sum_{i,j=1}^{d} (|\varphi(x)| + |(\frac{\partial}{\partial x_i} \varphi)(x)| + |(\frac{\partial^2}{\partial x_i \partial x_j} \varphi)(x)|)] < \infty$, the chain rule, and Lebesgue's dominated convergence theorem ensure that

(I) for all $x \in \mathbb{R}^d$ it holds that $(0, T] \ni t \mapsto u(t, x) \in \mathbb{R}$ is differentiable,

(II) for all $t \in [0, T]$ it holds that $\mathbb{R}^d \ni x \mapsto u(t, x) \in \mathbb{R}$ is twice differentiable,

(III) for all $t \in (0, T]$, $x \in \mathbb{R}^d$ it holds that

$$(\tfrac{\partial u}{\partial t})(t, x) = \mathbb{E}\Big[\big\langle (\nabla \varphi)(x + \sqrt{t} BZ), \tfrac{1}{2\sqrt{t}} BZ \big\rangle\Big], \qquad (12.18)$$

and

(IV) for all $t \in [0, T]$, $x \in \mathbb{R}^d$ it holds that

$$(\mathrm{Hess}_x u)(t, x) = \mathbb{E}\Big[(\mathrm{Hess}\,\varphi)(x + \sqrt{t} BZ)\Big]. \qquad (12.19)$$

Observe that items (III) and (IV), the assumption that $\varphi \in C^2(\mathbb{R}^d, \mathbb{R})$, the assumption that $\sup_{x \in \mathbb{R}^d} [\sum_{i,j=1}^{d} (|\varphi(x)| + |(\frac{\partial}{\partial x_i} \varphi)(x)| + |(\frac{\partial^2}{\partial x_i \partial x_j} \varphi)(x)|)] < \infty$, the fact that $\mathbb{E}[\|Z\|_2] < \infty$, and Lebesgue's dominated convergence theorem prove that $(0, T] \times \mathbb{R}^d \ni (t, x) \mapsto (\frac{\partial u}{\partial t})(t, x) \in \mathbb{R}$ and $[0, T] \times \mathbb{R}^d \ni (t, x) \mapsto (\mathrm{Hess}_x u)(t, x) \in \mathbb{R}^{d \times d}$ are continuous (cf. Definition 3.1.16). Next note that item (IV) and the fact that for all $X \in \mathbb{R}^{m \times d}$, $Y \in \mathbb{R}^{d \times m}$ it holds that $\mathrm{Trace}(XY) = \mathrm{Trace}(YX)$ imply that for all $t \in (0, T]$, $x \in \mathbb{R}^d$ it holds that

$$\tfrac{1}{2} \mathrm{Trace}\big(BB^*(\mathrm{Hess}_x u)(t, x)\big) = \mathbb{E}\Big[\tfrac{1}{2} \mathrm{Trace}\big(BB^*(\mathrm{Hess}\,\varphi)(x + \sqrt{t} BZ)\big)\Big]$$

$$= \tfrac{1}{2} \mathbb{E}\Big[\mathrm{Trace}\big(B^*(\mathrm{Hess}\,\varphi)(x + \sqrt{t} BZ)B\big)\Big] = \tfrac{1}{2} \mathbb{E}\Big[\sum_{k=1}^{m} \langle e_k, B^*(\mathrm{Hess}\,\varphi)(x + \sqrt{t} BZ)Be_k \rangle\Big]$$

$$= \tfrac{1}{2} \mathbb{E}\Big[\sum_{k=1}^{m} \langle Be_k, (\mathrm{Hess}\,\varphi)(x + \sqrt{t} BZ)Be_k \rangle\Big] = \tfrac{1}{2} \mathbb{E}\Big[\sum_{k=1}^{m} \varphi''(x + \sqrt{t} BZ)(Be_k, Be_k)\Big]$$

$$= \tfrac{1}{2t} \mathbb{E}\Big[\sum_{k=1}^{m} (\psi_{t,x})''(Z)(e_k, e_k)\Big] = \tfrac{1}{2t} \mathbb{E}\Big[\sum_{k=1}^{m} (\tfrac{\partial^2}{\partial y_k^2} \psi_{t,x})(Z)\Big] = \tfrac{1}{2t} \mathbb{E}[(\Delta \psi_{t,x})(Z)]$$

$$(12.20)$$

(cf. Definition 2.2.29). The assumption that $Z \colon \Omega \to \mathbb{R}^m$ is a standard normal random variable and integration by parts hence ensure that for all $t \in (0, T]$, $x \in \mathbb{R}^d$ it holds that

$$\tfrac{1}{2} \mathrm{Trace}\big(BB^*(\mathrm{Hess}_x u)(t, x)\big)$$

$$= \frac{1}{2t} \int_{\mathbb{R}^m} (\Delta \psi_{t,x})(y) \left[\frac{\exp(-\frac{\langle y, y \rangle}{2})}{(2\pi)^{m/2}}\right] dy = \frac{1}{2t} \int_{\mathbb{R}^m} \langle (\nabla \psi_{t,x})(y), y \rangle \left[\frac{\exp(-\frac{\langle y, y \rangle}{2})}{(2\pi)^{m/2}}\right] dy$$

$$= \frac{1}{2\sqrt{t}} \int_{\mathbb{R}^m} \Big\langle B^*(\nabla \varphi)(x + \sqrt{t} By), y \Big\rangle \left[\frac{\exp(-\frac{\langle y, y \rangle}{2})}{(2\pi)^{m/2}}\right] dy \qquad (12.21)$$

$$= \frac{1}{2\sqrt{t}} \mathbb{E}\Big[\langle B^*(\nabla \varphi)(x + \sqrt{t} BZ), Z \rangle\Big] = \mathbb{E}\Big[\big\langle (\nabla \varphi)(x + \sqrt{t} BZ), \tfrac{1}{2\sqrt{t}} BZ \big\rangle\Big].$$

Item (III) therefore proves that for all $t \in (0, T]$, $x \in \mathbb{R}^d$ it holds that

$$(\tfrac{\partial u}{\partial t})(t, x) = \tfrac{1}{2} \operatorname{Trace}\big(BB^*(\operatorname{Hess}_x u)(t, x)\big). \tag{12.22}$$

The fundamental theorem of calculus hence implies that for all $t, s \in (0, T]$, $x \in \mathbb{R}^d$ it holds that

$$u(t, x) - u(s, x) = \int_s^t (\tfrac{\partial u}{\partial t})(r, x)\, dr = \int_s^t \tfrac{1}{2} \operatorname{Trace}\big(BB^*(\operatorname{Hess}_x u)(r, x)\big)\, dr. \tag{12.23}$$

The fact that $[0, T] \times \mathbb{R}^d \ni (t, x) \mapsto (\operatorname{Hess}_x u)(t, x) \in \mathbb{R}^{d \times d}$ is continuous therefore ensures for all $t \in (0, T]$, $x \in \mathbb{R}^d$ that

$$\frac{u(t, x) - u(0, x)}{t} = \lim_{s \searrow 0}\left[\frac{u(t, x) - u(s, x)}{t}\right] = \frac{1}{t}\int_0^t \tfrac{1}{2} \operatorname{Trace}\big(BB^*(\operatorname{Hess}_x u)(r, x)\big)\, dr. \tag{12.24}$$

This and the fact that $[0, T] \times \mathbb{R}^d \ni (t, x) \mapsto (\operatorname{Hess}_x u)(t, x) \in \mathbb{R}^{d \times d}$ is continuous imply that for all $x \in \mathbb{R}^d$ it holds that

$$\begin{aligned}
&\limsup_{t \searrow 0}\left|\frac{u(t, x) - u(0, x)}{t} - \tfrac{1}{2}\operatorname{Trace}\big(BB^*(\operatorname{Hess}_x u)(0, x)\big)\right| \\
&\leq \limsup_{t \searrow 0}\left[\frac{1}{t}\int_0^t \left|\tfrac{1}{2}\operatorname{Trace}\big(BB^*(\operatorname{Hess}_x u)(s, x)\big) - \tfrac{1}{2}\operatorname{Trace}\big(BB^*(\operatorname{Hess}_x u)(0, x)\big)\right| ds\right] \\
&\leq \limsup_{t \searrow 0}\left[\sup_{s \in [0, t]}\left|\tfrac{1}{2}\operatorname{Trace}\big(BB^*\big((\operatorname{Hess}_x u)(s, x) - (\operatorname{Hess}_x u)(0, x)\big)\big)\right|\right] = 0.
\end{aligned} \tag{12.25}$$

Item (I) hence establishes that for all $x \in \mathbb{R}^d$ it holds that $[0, T] \ni t \mapsto u(t, x) \in \mathbb{R}$ is differentiable. Combining this with (12.25) and (12.22) ensures that for all $t \in [0, T]$, $x \in \mathbb{R}^d$ it holds that

$$(\tfrac{\partial u}{\partial t})(t, x) = \tfrac{1}{2}\operatorname{Trace}\big(BB^*(\operatorname{Hess}_x u)(t, x)\big). \tag{12.26}$$

This and the fact that $[0, T] \times \mathbb{R}^d \ni (t, x) \mapsto (\operatorname{Hess}_x u)(t, x) \in \mathbb{R}^{d \times d}$ is continuous establish item (i). In addition, note that (12.26) establishes item (ii). The proof of Proposition 12.1.3 is thus complete. $\square$

**Definition 12.1.4.** *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space. We say that $W$ is an $m$-dimensional $\mathbb{P}$-standard Brownian motion (we say that $W$ is a $\mathbb{P}$-standard Brownian motion, we say that $W$ is a standard Brownian motion) if and only if there exists $T \in (0, \infty)$ such that*

*(i) it holds that $m \in \mathbb{N}$,*

*(ii) it holds that $W : [0, T] \times \Omega \times \mathbb{R}^m$ is a function,*

*(iii) it holds for all $\omega \in \Omega$ that $[0, T] \ni s \mapsto W_s(\omega) \in \mathbb{R}^m$ is continuous,*

*(iv) it holds for all $\omega \in \Omega$ that $W_0(\omega) = 0 \in \mathbb{R}^m$,*

*(v) it holds for all $t_1 \in [0, T]$, $t_2 \in [0, T]$ with $t_1 < t_2$ that $\Omega \ni \omega \mapsto (t_2 - t_1)^{-1/2}(W_{t_2}(\omega) - W_{t_1}(\omega)) \in \mathbb{R}^m$ is a standard normal random variable, and*

Figure 12.1: Four trajectories of a 1-dimensional standard Brownian motion

*(vi) it holds for all $n \in \{3, 4, 5, \dots\}$, $t_1, t_2, \dots, t_n \in [0, T]$ with $t_1 \leq t_2 \leq \cdots \leq t_n$ that $W_{t_2} - W_{t_1}, W_{t_3} - W_{t_2}, \dots, W_{t_n} - W_{t_{n-1}}$ are independent.*

**Corollary 12.1.5.** *Let $T \in (0, \infty)$, $d, m \in \mathbb{N}$, $B \in \mathbb{R}^{d \times m}$, $\varphi \in C^2(\mathbb{R}^d, \mathbb{R})$ satisfy $\sup_{x \in \mathbb{R}^d} \left[ \sum_{i,j=1}^d \left( |\varphi(x)| + |(\frac{\partial}{\partial x_i} \varphi)(x)| + |(\frac{\partial^2}{\partial x_i \partial x_j} \varphi)(x)| \right) \right] < \infty$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $W: [0, T] \times \Omega \to \mathbb{R}^m$ be a standard Brownian motion, and let $u: [0, T] \times \mathbb{R}^d \to \mathbb{R}$ satisfy for all $t \in [0, T]$, $x \in \mathbb{R}^d$ that*

$$u(t, x) = \mathbb{E}\Big[ \varphi(x + BW_t) \Big] \tag{12.27}$$

*(cf. Definition 12.1.4). Then*

*(i) it holds that $u \in C^{1,2}([0, T] \times \mathbb{R}^d, \mathbb{R})$ and*

*(ii) it holds for all $t \in [0, T]$, $x \in \mathbb{R}^d$ that*

$$\big(\tfrac{\partial u}{\partial t}\big)(t, x) = \tfrac{1}{2} \operatorname{Trace}\big( BB^* (\operatorname{Hess}_x u)(t, x) \big) \tag{12.28}$$

*(cf. Definition 2.2.29).*

*Proof of Corollary 12.1.5.* First, observe that the assumption that $W: [0, T] \times \Omega \to \mathbb{R}^m$ is a standard Brownian motion ensures that for all $t \in [0, T]$, $x \in \mathbb{R}^d$ it holds that

$$u(t, x) = \mathbb{E}[\varphi(x + BW_t)] = \mathbb{E}\left[ \varphi\left( x + \sqrt{t} B \frac{W_T}{\sqrt{T}} \right) \right]. \tag{12.29}$$

The fact that $\frac{W_T}{\sqrt{T}} \colon \Omega \to \mathbb{R}^m$ is a standard normal random variable and Proposition 12.1.3 hence establish items (i) and (ii). The proof of Corollary 12.1.5 is thus complete. $\quad\square$

### 12.1.3.2 Feynman–Kac formulas providing uniqueness of solutions

**Lemma 12.1.6** (A special case of Vitali's convergence theorem). *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $X_n \colon \Omega \to \mathbb{R}$, $n \in \mathbb{N}_0$, be random variables with $\mathbb{P}(\limsup_{n\to\infty} |X_n - X_0| = 0) = 1$, and let $p \in (1, \infty)$ satisfy $\sup_{n\in\mathbb{N}} \mathbb{E}[|X_n|^p] < \infty$. Then*

(i) *it holds that $\limsup_{n\to\infty} \mathbb{E}[|X_n - X_0|] = 0$,*

(ii) *it holds that $\mathbb{E}[|X_0|] < \infty$, and*

(iii) *it holds that $\limsup_{n\to\infty} |\mathbb{E}[X_n] - \mathbb{E}[X_0]| = 0$.*

**Proposition 12.1.7.** *Let $d \in \mathbb{N}$, $T, \rho \in (0, \infty)$, let $f \in C([0, T] \times \mathbb{R}^d, \mathbb{R})$, let $u \in C^{1,2}([0, T] \times \mathbb{R}^d, \mathbb{R})$ have at most polynomially growing partial derivatives, assume for all $t \in [0, T]$, $x \in \mathbb{R}^d$ that*

$$(\tfrac{\partial u}{\partial t})(t, x) = \rho\, (\Delta_x u)(t, x) + f(t, x), \tag{12.30}$$

*let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, and let $W \colon [0, T] \times \Omega \to \mathbb{R}^d$ be a standard Brownian motion (cf. Definition 12.1.4). Then it holds for all $t \in [0, T]$, $x \in \mathbb{R}^d$ that*

$$u(t, x) = \mathbb{E}\!\left[ u(0, x + \sqrt{2\rho}\,W_t) + \int_0^t f(t - s, x + \sqrt{2\rho}\,W_s)\, ds \right]. \tag{12.31}$$

*Proof of Proposition 12.1.7.* Throughout this proof let $\langle \cdot, \cdot \rangle \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ satisfy for all $x = (x_1, x_2, \ldots, x_d)$, $y = (y_1, y_2, \ldots, y_d) \in \mathbb{R}^d$ that $\langle x, y \rangle = \sum_{i=1}^d x_i y_i$, let $D_1 \colon [0, T] \times \mathbb{R}^d \to \mathbb{R}$ satisfy for all $t \in [0, T]$, $x \in \mathbb{R}^d$ that $D_1(t, x) = (\tfrac{\partial u}{\partial t})(t, x)$, let $D_2 = (D_{2,1}, D_{2,2}, \ldots, D_{2,d}) \colon [0, T] \times \mathbb{R}^d \to \mathbb{R}^d$ satisfy for all $t \in [0, T]$, $x \in \mathbb{R}^d$ that $D_2(t, x) = (\nabla_x u)(t, x)$, let $H = (H_{i,j})_{i,j\in\{1,2,\ldots,d\}} \colon [0, T] \times \mathbb{R}^d \to \mathbb{R}^{d \times d}$ satisfy for all $t \in [0, T]$, $x \in \mathbb{R}^d$ that $H(t, x) = (\mathrm{Hess}_x u)(t, x)$, let $\gamma \colon \mathbb{R}^d \to \mathbb{R}$ satisfy for all $z \in \mathbb{R}^d$ that

$$\gamma(z) = (2\pi)^{-d/2} \exp\!\left( -\tfrac{\|z\|_2^2}{2} \right), \tag{12.32}$$

and let $v_{t,x} \colon [0, t] \to \mathbb{R}$, $t \in [0, T]$, $x \in \mathbb{R}^d$, satisfy for all $t \in [0, T]$, $x \in \mathbb{R}^d$, $s \in [0, t]$ that

$$v_{t,x}(s) = \mathbb{E}\!\left[ u(s, x + \sqrt{2\rho}\,W_{t-s}) \right]. \tag{12.33}$$

Note that the assumption that $W$ is a standard Brownian motion implies that for all $t \in (0, T]$, $s \in [0, t)$ it holds that $(t - s)^{-1/2} W_{t-s} \colon \Omega \to \mathbb{R}^d$ is a standard normal random variable. This ensures that for all $t \in (0, T]$, $x \in \mathbb{R}^d$, $s \in [0, t)$ it holds that

$$v_{t,x}(s) = \mathbb{E}\!\left[ u(s, x + \sqrt{2\rho(t - s)}\,(t - s)^{-1/2} W_{t-s}) \right] = \int_{\mathbb{R}^d} u(s, x + \sqrt{2\rho(t - s)}\,z)\gamma(z)\, dz. \tag{12.34}$$

The assumption that $u$ has at most polynomially growing partial derivatives, the fact that $(0, \infty) \ni s \mapsto \sqrt{s} \in (0, \infty)$ is differentiable, the chain rule, and Vitali's convergence theorem hence show that for all $t \in (0, T]$, $x \in \mathbb{R}^d$, $s \in [0, t)$ it holds that $v_{t,x}|_{[0,t)} \in C^1([0, t), \mathbb{R})$ and

$$(v_{t,x})'(s) = \int_{\mathbb{R}^d} \left[ D_1(s, x + \sqrt{2\rho(t - s)}\,z) + \left\langle D_2(s, x + \sqrt{2\rho(t - s)}\,z), \frac{-\rho z}{\sqrt{2\rho(t - s)}} \right\rangle \right] \gamma(z)\, dz. \tag{12.35}$$

Next note that the fact that for all $z \in \mathbb{R}^d$ it holds that $(\nabla \gamma)(z) = -\gamma(z)z$ implies that for all $t \in (0,T]$, $x \in \mathbb{R}^d$, $s \in [0,t)$ it holds that

$$\int_{\mathbb{R}^d} \left\langle D_2(s, x + \sqrt{2\rho(t-s)}z), \frac{-\rho z}{\sqrt{2\rho(t-s)}} \right\rangle \gamma(z)\, dz = \int_{\mathbb{R}^d} \left\langle D_2(s, x + \sqrt{2\rho(t-s)}z), \frac{\rho(\nabla \gamma)(z)}{\sqrt{2\rho(t-s)}} \right\rangle dz$$

$$= \frac{\rho}{\sqrt{2\rho(t-s)}} \sum_{i=1}^{d} \left[ \int_{\mathbb{R}^d} D_{2,i}(s, x + \sqrt{2\rho(t-s)}z)(\tfrac{\partial \gamma}{\partial z_i})(z_1, z_2, \ldots, z_d)\, dz \right].$$

$$(12.36)$$

Next observe that integration by parts proves that for all $t \in (0,T]$, $x \in \mathbb{R}^d$, $s \in [0,t)$, $i \in \{1, 2, \ldots, d\}$, $a \in \mathbb{R}$, $b \in (a, \infty)$ it holds that

$$\int_a^b D_{2,i}(s, x + \sqrt{2\rho(t-s)}(z_1, z_2, \ldots, z_d))(\tfrac{\partial \gamma}{\partial z_i})(z_1, z_2, \ldots, z_d)\, dz_i$$

$$= \left[ D_{2,i}(s, x + \sqrt{2\rho(t-s)}(z_1, z_2, \ldots, z_d))\gamma(z_1, z_2, \ldots, z_d) \right]_{z_i=a}^{z_i=b} \qquad (12.37)$$

$$- \int_a^b \sqrt{2\rho(t-s)} H_{i,i}(s, x + \sqrt{2\rho(t-s)}(z_1, z_2, \ldots, z_d))\gamma(z_1, z_2, \ldots, z_d)\, dz_i.$$

The assumption that $u$ has at most polynomially growing derivatives hence implies that for all $t \in (0,T]$, $x \in \mathbb{R}^d$, $s \in [0,t)$, $i \in \{1, 2, \ldots, d\}$ it holds that

$$\int_{\mathbb{R}} D_{2,i}(s, x + \sqrt{2\rho(t-s)}(z_1, z_2, \ldots, z_d))(\tfrac{\partial \gamma}{\partial z_i})(z_1, z_2, \ldots, z_d)\, dz_i$$

$$= -\sqrt{2\rho(t-s)} \int_{\mathbb{R}} H_{i,i}(s, x + \sqrt{2\rho(t-s)}(z_1, z_2, \ldots, z_d))\gamma(z_1, z_2, \ldots, z_d)\, dz_i.$$

$$(12.38)$$

Combining this with (12.36) and Fubini's theorem ensures that for all $t \in (0,T]$, $x \in \mathbb{R}^d$, $s \in [0,t)$ it holds that

$$\int_{\mathbb{R}^d} \left\langle D_2(s, x + \sqrt{2\rho(t-s)}z), \frac{-\rho z}{\sqrt{2\rho(t-s)}} \right\rangle \gamma(z)\, dz = -\rho \sum_{i=1}^{d} \int_{\mathbb{R}^d} H_{i,i}(s, x + \sqrt{2\rho(t-s)}(z))\gamma(z)\, dz$$

$$= -\int_{\mathbb{R}^d} \rho \operatorname{Trace}\big(H(s, x + \sqrt{2\rho(t-s)}(z))\big)\gamma(z)\, dz.$$

$$(12.39)$$

This, (12.35), (12.30), and the fact that for all $t \in (0,T]$, $s \in [0,t)$ it holds that $(t-s)^{-1/2}W_{t-s} \colon \Omega \to \mathbb{R}^d$ is a standard normal random variable imply that for all $t \in (0,T]$, $x \in \mathbb{R}^d$, $s \in [0,t)$ it holds that

$$(v_{t,x})'(s) = \int_{\mathbb{R}^d} \left[ D_1(s, x + \sqrt{2\rho(t-s)}z) - \rho \operatorname{Trace}\big(H(s, x + \sqrt{2\rho(t-s)}z)\big) \right] \gamma(z)\, dz$$

$$= \int_{\mathbb{R}^d} f(s, x + \sqrt{2\rho(t-s)}z)\gamma(z)\, dz = \mathbb{E}\Big[ f(s, x + \sqrt{2\rho}W_{t-s}) \Big].$$

$$(12.40)$$

The fact that $W_0 = 0$, the fact that for all $t \in [0,T]$, $x \in \mathbb{R}^d$ it holds that $v_{t,x} \colon [0,t] \to \mathbb{R}$ is a continuous function, and the fundamental theorem of calculus therefore demonstrate

that for all $t \in [0, T]$, $x \in \mathbb{R}^d$ it holds that

$$
\begin{aligned}
u(t, x) &= \mathbb{E}\Big[u(t, x + \sqrt{2\rho} W_{t-t})\Big] = v_{t,x}(t) = v_{t,x}(0) + \int_0^t (v_{t,x})'(s) \, ds \\
&= \mathbb{E}\Big[u(0, x + \sqrt{2\rho} W_t)\Big] + \int_0^t \mathbb{E}\Big[f(s, x + \sqrt{2\rho} W_{t-s})\Big] \, ds.
\end{aligned}
\tag{12.41}
$$

Fubini's theorem and the fact that $u$ and $f$ are at most polynomially growing hence establish (12.31). This completes the proof of Proposition 12.1.7. $\qquad\square$

**Corollary 12.1.8.** *Let $d \in \mathbb{N}$, $T, \rho \in (0, \infty)$, $\varrho = \sqrt{2\rho T}$, $a \in \mathbb{R}$, $b \in (a, \infty)$, let $\varphi \colon \mathbb{R}^d \to \mathbb{R}$ be a function, let $u \in C^{1,2}([0,T] \times \mathbb{R}^d, \mathbb{R})$ have at most polynomially growing partial derivatives, assume for all $t \in [0, T]$, $x \in \mathbb{R}^d$ that $u(0, x) = \varphi(x)$ and*

$$
(\tfrac{\partial u}{\partial t})(t, x) = \rho \, (\Delta_x u)(t, x),
\tag{12.42}
$$

*let $(\Omega, \mathbb{F}, \mathbb{P})$ be a probability space, and let $\mathbb{W} \colon \Omega \to \mathbb{R}^d$ be a standard normal random variable. Then*

  *(i) it holds that the function $\varphi \colon \mathbb{R}^d \to \mathbb{R}$ is twice continuously differentiable with at most polynomially growing derivatives and*

  *(ii) it holds for every $x \in \mathbb{R}^d$ that $u(T, x) = \mathbb{E}\big[\varphi(\varrho \mathbb{W} + x)\big]$.*

*Proof of Corollary 12.1.8.* Note that the assumption that $u \in C^{1,2}([0,T] \times \mathbb{R}^d, \mathbb{R})$ has at most polynomially growing partial derivatives and the fact that for all $x \in \mathbb{R}^d$ it holds that $\varphi(x) = u(0, x)$ establish item (i). Next observe that Proposition 12.1.7 proves item (ii). The proof of Corollary 12.1.8 is thus complete. $\qquad\square$

**Definition 12.1.9.** *Let $d \in \mathbb{N}$ and let $f \colon \mathbb{R}^d \to \mathbb{R}$ and $g \colon \mathbb{R}^d \to \mathbb{R}$ be $\mathcal{B}(\mathbb{R}^d)/\mathcal{B}(\mathbb{R})$-measurable functions. Then we denote by $f * g \colon \big\{x \in \mathbb{R}^d \colon \min\{\int_{\mathbb{R}^d} \max\{0, f(x-y)g(y)\} \, dy, - \int_{\mathbb{R}^d} \min\{0, y)g(y)\} \, dy\} < \infty\big\} \to [-\infty, \infty]$ the function which satisfies for all $x \in \mathbb{R}^d$ with $\min\{\int_{\mathbb{R}^d} \max\{0, f(x-y)g(y)\} \, dy, - \int_{\mathbb{R}^d} \min\{0, f(x-y)g(y)\} \, dy\} < \infty$ that*

$$
(f * g)(x) = \int_{\mathbb{R}^d} f(x - y)g(y) \, dy.
\tag{12.43}
$$

**Exercise 12.1.1.** *Let $d \in \mathbb{N}$, $T \in (0, \infty)$, let $\gamma_\sigma \colon \mathbb{R}^d \to \mathbb{R}$, $\sigma \in (0, \infty)$, satisfy for all $\sigma \in (0, \infty)$, $x \in \mathbb{R}^d$ that*

$$
\gamma_\sigma(x) = (2\pi\sigma^2)^{-\frac{d}{2}} \exp\left(\frac{-\|x\|_2^2)}{2\sigma^2}\right),
\tag{12.44}
$$

*and for every $\rho \in (0, \infty)$ and $\varphi \in C^2(\mathbb{R}^d, \mathbb{R})$ with $\sup_{x \in \mathbb{R}^d} \big[\sum_{i,j=1}^d \big(|\varphi(x)| + |(\tfrac{\partial}{\partial x_i}\varphi)(x)| + |(\tfrac{\partial^2}{\partial x_i \partial x_j}\varphi)(x)|\big)\big] < \infty$ let $u_{\rho,\varphi} \colon [0, T] \times \mathbb{R}^d \to \mathbb{R}$ satisfy for all $t \in (0, T]$, $x \in \mathbb{R}^d$ that*

$$
u_{\rho,\varphi}(0, x) = \varphi(x) \qquad and \qquad u_{\rho,\varphi}(t, x) = (\varphi * \gamma_{\sqrt{2t\rho}})(x).
\tag{12.45}
$$

*Prove or disprove the following statement: For every $\rho \in (0, \infty)$ and for every $\varphi \in C^2(\mathbb{R}^d, \mathbb{R})$ with $\sup_{x \in \mathbb{R}^d} \big[\sum_{i,j=1}^d \big(|\varphi(x)| + |(\tfrac{\partial}{\partial x_i}\varphi)(x)| + |(\tfrac{\partial^2}{\partial x_i \partial x_j}\varphi)(x)|\big)\big] < \infty$ it holds for all $t \in (0, T]$, $x \in \mathbb{R}^d$ that $u_{\rho,\varphi} \in C^{1,2}([0, T] \times \mathbb{R}^d, \mathbb{R})$ and*

$$
(\tfrac{\partial u_{\rho,\varphi}}{\partial t})(t, x) = \rho(\Delta_x u_{\rho,\varphi})(t, x).
\tag{12.46}
$$

**Exercise 12.1.2.** *Prove or disprove the following statement: For every $x \in \mathbb{R}$ it holds that*

$$e^{-x^2/2} = \frac{1}{\sqrt{2\pi}}\left[\int_{\mathbb{R}} e^{-t^2/2} e^{-\mathrm{i}xt}\, dt\right]. \tag{12.47}$$

**Exercise 12.1.3.** *Let $d \in \mathbb{N}$, $T \in (0, \infty)$, let $\gamma_\sigma \colon \mathbb{R}^d \to \mathbb{R}$, $\sigma \in (0, \infty)$, satisfy for all $\sigma \in (0, \infty)$, $x \in \mathbb{R}^d$ that*

$$\gamma_\sigma(x) = (2\pi\sigma^2)^{-\frac{d}{2}} \exp\left(\frac{-\|x\|_2^2}{2\sigma^2}\right), \tag{12.48}$$

*for every $\varphi \in C^2(\mathbb{R}^d, \mathbb{R})$ with $\sup_{x \in \mathbb{R}^d}\left[\sum_{i,j=1}^{d}\left(|\varphi(x)| + |(\frac{\partial}{\partial x_i}\varphi)(x)| + |(\frac{\partial^2}{\partial x_i \partial x_j}\varphi)(x)|\right)\right] < \infty$ let $u_\varphi \colon [0,T] \times \mathbb{R}^d \to \mathbb{R}$ satisfy for all $t \in (0,T]$, $x \in \mathbb{R}^d$ that*

$$u_\varphi(0,x) = \varphi(x) \qquad and \qquad u_\varphi(t,x) = (\varphi * \gamma_{\sqrt{2t}})(x), \tag{12.49}$$

*and let $\psi_i \colon \mathbb{R}^d \to \mathbb{R}$, $i \in \mathbb{N}^d$ satisfy for all $i = (i_1, i_2, \ldots, i_d) \in \mathbb{N}^d$, $x = (x_1, x_2, \ldots, x_d) \in \mathbb{R}^d$ that*

$$\psi_i(x) = 2^{\frac{d}{2}}\left[\prod_{k=1}^{d} \sin(i_k \pi x_k)\right]. \tag{12.50}$$

*Prove or disprove the following statement: For all $i = (i_1, i_2, \ldots, i_d) \in \mathbb{N}^d$, $t \in [0,T]$, $x \in \mathbb{R}^d$ it holds that*

$$u_{\psi_i}(t,x) = \exp\left(-\pi^2\left[\sum_{k=1}^{d} |i_k|^2\right]t\right)\psi_i(x). \tag{12.51}$$

**Exercise 12.1.4.** *Let $d \in \mathbb{N}$, $T \in (0, \infty)$, let $\gamma_\sigma \colon \mathbb{R}^d \to \mathbb{R}$, $\sigma \in (0, \infty)$, satisfy for all $\sigma \in (0, \infty)$, $x \in \mathbb{R}^d$ that*

$$\gamma_\sigma(x) = (2\pi\sigma^2)^{-\frac{d}{2}} \exp\left(\frac{-\|x\|_2^2}{2\sigma^2}\right), \tag{12.52}$$

*and let $\psi_i \colon \mathbb{R}^d \to \mathbb{R}$, $i \in \mathbb{N}^d$, satisfy for all $i = (i_1, i_2, \ldots, i_d) \in \mathbb{N}^d$, $x = (x_1, x_2, \ldots, x_d) \in \mathbb{R}^d$ that*

$$\psi_i(x) = 2^{\frac{d}{2}}\left[\prod_{k=1}^{d} \sin(i_k \pi x_k)\right]. \tag{12.53}$$

*Prove or disprove the following statement: For every $i \in \mathbb{N}^d$, $s \in [0,T]$, $y \in \mathbb{R}^d$, and every function $u \in C^{1,2}([0,T] \times \mathbb{R}^d, \mathbb{R})$ with at most polynomially growing derivatives such that it holds for all $t \in (0,T)$, $x \in \mathbb{R}^d$ that $u(0,x) = \psi_i(x)$ and*

$$\left(\tfrac{\partial u}{\partial t}\right)(t,x) = (\Delta_x u)(t,x) \tag{12.54}$$

*it holds that*

$$u(s,y) = \exp\left(-\pi^2\left[\sum_{k=1}^{d} |i_k|^2\right]s\right)\psi_i(y). \tag{12.55}$$

## 12.1.4 Stochastic optimization problems for PDEs

The proof of Proposition 12.1.10 is based on an application of Proposition 12.1.2 and Proposition 12.1.7. A detailed proof of Proposition 12.1.10 can be found in Beck et al. [1, Proposition 2.7].

**Proposition 12.1.10.** *Let $d \in \mathbb{N}$, $T, \rho \in (0, \infty)$, $\varrho = \sqrt{2\rho T}$, $a \in \mathbb{R}$, $b \in (a, \infty)$, let $\varphi \colon \mathbb{R}^d \to \mathbb{R}$ be a function, let $u \in C^{1,2}([0,T] \times \mathbb{R}^d, \mathbb{R})$ have at most polynomially growing partial derivatives, assume for all $t \in [0,T]$, $x \in \mathbb{R}^d$ that $u(0,x) = \varphi(x)$ and*

$$(\tfrac{\partial u}{\partial t})(t, x) = \rho \, (\Delta_x u)(t, x), \tag{12.56}$$

*let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $\mathbb{W} \colon \Omega \to \mathbb{R}^d$ be a standard normal random variable, let $\xi \colon \Omega \to [a,b]^d$ be a continuous uniformly distributed random variable, and assume that $\mathbb{W}$ and $\xi$ are independent. Then*

*(i) it holds that the function $\varphi \colon \mathbb{R}^d \to \mathbb{R}$ is twice continuously differentiable with at most polynomially growing derivatives,*

*(ii) it holds that there exists a unique continuous function $U \colon [a,b]^d \to \mathbb{R}$ such that*

$$\mathbb{E}\big[|\varphi(\varrho\mathbb{W} + \xi) - U(\xi)|^2\big] = \inf_{v \in C([a,b]^d, \mathbb{R})} \mathbb{E}\big[|\varphi(\varrho\mathbb{W} + \xi) - v(\xi)|^2\big], \tag{12.57}$$

*and*

*(iii) it holds for every $x \in [a,b]^d$ that $U(x) = u(T, x)$.*

*Proof of Proposition 12.1.10.* First, note that (12.56), the assumption that $\mathbb{W}$ is a standard normal random variable, and Corollary 12.1.8 prove that for all $x \in \mathbb{R}^d$ it holds that the function $\varphi \colon \mathbb{R}^d \to \mathbb{R}$ is twice continuously differentiable with at most polynomially growing derivatives and

$$u(T, x) = \mathbb{E}\big[u(0, \varrho\mathbb{W} + x)\big] = \mathbb{E}\big[\varphi(\varrho\mathbb{W} + x)\big]. \tag{12.58}$$

Moreover, observe that the assumption that $\mathbb{W}$ is a standard normal random variable, the fact that $\varphi$ is continuous, and the fact that $\varphi$ is at most polynomially growing and continuous ensure that

(I) it holds that $[a,b]^d \times \Omega \ni (x, \omega) \mapsto \varphi(\varrho\mathbb{W}(\omega) + x) \in \mathbb{R}$ is $(\mathcal{B}([a,b]^d) \otimes \mathcal{F})/\mathcal{B}(\mathbb{R})$-measurable and

(II) it holds for all $x \in [a,b]^d$ that $\mathbb{E}[|\varphi(\varrho\mathbb{W} + x)|^2] < \infty$.

Proposition 12.1.2 and (12.58) hence ensure that

(A) there exists a unique continuous function $U \colon [a,b]^d \to \mathbb{R}$ which satisfies that

$$\int_{[a,b]^d} \mathbb{E}\big[|\varphi(\varrho\mathbb{W} + x) - U(x)|^2\big] \, \mathrm{d}x = \inf_{v \in C([a,b]^d, \mathbb{R})} \left( \int_{[a,b]^d} \mathbb{E}\big[|\varphi(\varrho\mathbb{W} + x) - v(x)|^2\big] \, \mathrm{d}x \right) \tag{12.59}$$

*and*

(B) it holds for all $x \in [a,b]^d$ that $U(x) = u(T, x)$.

Next note that the assumption that $\mathbb{W}$ and $\xi$ are independent, item (I), and the assumption that $\xi$ is continuously uniformly distributed on $[a,b]^d$ imply that for all $v \in C([a,b]^d, \mathbb{R})$ it holds that

$$\mathbb{E}\big[|\varphi(\varrho\mathbb{W} + \xi) - v(\xi)|^2\big] = \frac{1}{(b-a)^d} \int_{[a,b]^d} \mathbb{E}\big[|\varphi(\varrho\mathbb{W} + x) - v(x)|^2\big] \, \mathrm{d}x. \tag{12.60}$$

Combining this with item (A) establishes item (ii). In addition, observe that items (A) and (B) and (12.60) establish item (iii). This completes the proof of Proposition 12.1.10. $\square$

### 12.1.5 Towards a deep learning scheme for PDEs

Let $d \in \mathbb{N}$, $T, \rho \in (0, \infty)$, $\varrho = \sqrt{2\rho T}$, $a \in \mathbb{R}$, $b \in (a, \infty)$, let $\varphi \colon \mathbb{R}^d \to \mathbb{R}$ be a function, let $u \in C^{1,2}([0,T] \times \mathbb{R}^d, \mathbb{R})$ have at most polynomially growing partial derivatives, assume for all $t \in [0,T]$, $x \in \mathbb{R}^d$ that $u(0,x) = \varphi(x)$ and

$$(\tfrac{\partial u}{\partial t})(t,x) = \rho\,(\Delta_x u)(t,x) \tag{12.61}$$

let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $\mathbb{W} \colon \Omega \to \mathbb{R}^d$ be a standard normally distributed random variable, let $\xi \colon \Omega \to [a,b]^d$ be a continuous uniformly distributed random variable, assume that $\mathbb{W}$ and $\xi$ are independent. Proposition 12.1.10 then ensures that the solution $u$ of the heat equation in (12.61) at time $T$ on $[a,b]^d$ is the unique global minimizer of the function

$$C([a,b]^d, \mathbb{R}) \ni v \mapsto \mathbb{E}\big[|\varphi(\varrho\mathbb{W} + \xi) - v(\xi)|^2\big] \in [0, \infty). \tag{12.62}$$

Now an idea of a simply machine learning based approximation method for PDEs (see [1]) is to approximate the set $C([a,b]^d, \mathbb{R})$ of all continuous functions from $[a,b]^d$ to $\mathbb{R}$ through the set of all deep artificial neural networks with a fixed sufficiently large architecture. More formally, let $L \in \mathbb{N}$, $l_1, l_2, \dots, l_L \in \mathbb{N}$, $\mathfrak{d} = (dl_1 + l_1) + (\sum_{k=2}^{L} l_k(l_{k-1} + 1)) + (l_L + 1)$ and consider the function

$$\left\{ w \in C([a,b]^d, \mathbb{R}) \colon \begin{bmatrix} \exists\,\theta \in \mathbb{R}^{\mathfrak{d}} \colon \forall\, x \in [a,b]^d \colon \\ w(x) = (\mathcal{N}^{\theta,d}_{\mathfrak{R}_{l_1},\dots,\mathfrak{R}_{l_L},\mathrm{id}_{\mathbb{R}}})(x) \end{bmatrix} \right\} \ni v \mapsto \mathbb{E}\big[|\varphi(\varrho\mathbb{W} + \xi) - v(\xi)|^2\big] \in [0, \infty) \tag{12.63}$$

(cf. Definition 2.1.2). The approach of the machine learning scheme in Beck et al. [1] is then to approximatively compute a suitable minimizer of the function in (12.63) and to view the resulting approximation of a suitable minimizer of the function in (12.63) as an approximation of the solution $u$ of the heat equation in (12.61) at time $T$ on $[a,b]^d$. To approximatively compute a suitable minimizer of the function in (12.63), we reformulate (12.63) by employing the parametrization function induced by artificial neural networks to obtain the function

$$\mathbb{R}^{\mathfrak{d}} \ni \theta \mapsto \mathbb{E}\Big[\big|\varphi(\varrho\mathbb{W} + \xi) - (\mathcal{N}^{\theta,d}_{\mathfrak{R}_{l_1},\dots,\mathfrak{R}_{l_L},\mathrm{id}_{\mathbb{R}}})(\xi)\big|^2\Big] \in [0, \infty). \tag{12.64}$$

A suitable minimizer of the function in (12.64) can then be approximatively computed by means of stochastic gradient descent optimization algorithms. We refer to Beck et al. [1] for numerical simulations.

## 12.2 Nonlinear PDEs

### 12.2.1 Splitting approximations

**Theorem 12.2.1.** *Let $T \in (0, \infty)$, $p \in [1, \infty)$, $f \in C^2(\mathbb{R}, \mathbb{R})$, let $u_d \in C^{1,2}([0,T] \times \mathbb{R}^d, \mathbb{R})$, $d \in \mathbb{N}$, satisfy for all $d \in \mathbb{N}$, $t \in [0,T]$, $x \in \mathbb{R}^d$ that*

$$(\tfrac{\partial}{\partial t} u_d)(t,x) = (\Delta_x u_d)(t,x) + f(u_d(t,x)), \tag{12.65}$$

*and assume for all $d \in \mathbb{N}$, $i,j \in \{1,2,\dots,d\}$ that $\sup_{t \in [0,T]} \sup_{x = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d}\big[(1 + \sum_{k=1}^{d} |x_k|)^{-p}(|(\tfrac{\partial^2}{\partial x_i \partial x_j} u_d)(t,x)| + |(\tfrac{\partial}{\partial t} u_d)(t,x)| + |f''(x_1)| + |f'(x_1)|]\big] < \infty$. Then*

(i) *there exist unique at most polynomially growing* $\mathcal{U}_n^{d,N} \in C^{1,2}([\frac{(n-1)T}{N}, \frac{nT}{N}] \times \mathbb{R}^d, \mathbb{R})$, $d, N \in \mathbb{N}$, $n \in \{0, 1, \dots, N\}$, *which satisfy for all* $d, N \in \mathbb{N}$, $n \in \{0, 1, \dots, N-1\}$, $t \in [\frac{nT}{N}, \frac{(n+1)T}{N}]$, $s \in [\frac{-T}{N}, 0]$, $x \in \mathbb{R}^d$ *that* $\mathcal{U}_{n+1}^{d,N}(\frac{nT}{N}, x) = \mathcal{U}_n^{d,N}(\frac{nT}{N}, x) + \frac{T}{N} f(\mathcal{U}_n^{d,N}(\frac{nT}{N}, x))$, $\mathcal{U}_0^{d,N}(s, x) = u_d(0, x)$, *and*

$$(\tfrac{\partial}{\partial t}\mathcal{U}_{n+1}^{d,N})(t, x) = (\Delta_x \mathcal{U}_{n+1}^{d,N})(t, x) \tag{12.66}$$

*and*

(ii) *there exists* $c \in \mathbb{R}$ *such that for all* $d, N \in \mathbb{N}$, $x = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$ *it holds that*

$$|\mathcal{U}_N^{d,N}(T, x) - u_d(T, x)| \le cd^{p+1}N^{-1/2}\big(1 + \textstyle\sum_{i=1}^d |x_i|\big)^p. \tag{12.67}$$

## 12.2.2 DNN approximation result

The next result, Theorem 12.2.1 below, establishes a DNN approximation result for nonlinear PDEs (cf. Grohs et al. [11] and Hutzenthaler et al. [15]).

**Theorem 12.2.2.** *Let* $T, p, \kappa \in (0, \infty)$, $(\mathfrak{g}_{d,\varepsilon})_{d\in\mathbb{N},\varepsilon\in(0,1]} \subseteq \mathbf{N}$, $(c_d)_{d\in\mathbb{N}} \subseteq (0, \infty)$ *(cf. Definition 2.2.1), let* $f\colon \mathbb{R} \to \mathbb{R}$ *be globally Lipschitz continuous, for every* $d \in \mathbb{N}$ *let* $u_d \in C^{1,2}([0, T] \times \mathbb{R}^d, \mathbb{R})$, *and assume for all* $d \in \mathbb{N}$, $x \in \mathbb{R}^d$, $\varepsilon \in (0, 1]$, $t \in (0, T)$ *that* $\mathcal{R}_{\mathfrak{r}}(\mathfrak{g}_{d,\varepsilon}) \in C(\mathbb{R}^d, \mathbb{R})$, $|(\mathcal{R}_{\mathfrak{r}}(\mathfrak{g}_{d,\varepsilon}))(x)| \le \kappa d^\kappa (1 + \|x\|_2^\kappa)$, $|u_d(0, x) - (\mathcal{R}_{\mathfrak{r}}(\mathfrak{g}_{d,\varepsilon}))(x)| \le \varepsilon\kappa d^\kappa (1 + \|x\|_2^\kappa)$, $\mathcal{P}(\mathfrak{g}_{d,\varepsilon}) \le \kappa d^\kappa \varepsilon^{-\kappa}$, $|u_d(t, x)| \le c_d(1 + \|x\|_2^{c_d})$, *and*

$$(\tfrac{\partial u_d}{\partial t})(t, x) = (\Delta_x u_d)(t, x) + f(u_d(t, x)) \tag{12.68}$$

*(cf. Definitions 2.1.6 and 2.2.3). Then there exist* $(\mathfrak{u}_{d,\varepsilon})_{d\in\mathbb{N},\varepsilon\in(0,1]} \subseteq \mathbf{N}$ *and* $\eta \in (0, \infty)$ *such that for all* $d \in \mathbb{N}$, $\varepsilon \in (0, 1]$ *it holds that* $\mathcal{P}(\mathfrak{u}_{d,\varepsilon}) \le \eta d^\eta \varepsilon^{-\eta}$ *and*

$$\left[\int_{[0,T]\times[0,1]^d} |u_d(y) - (\mathcal{R}_{\mathfrak{r}}(\mathfrak{u}_{d,\varepsilon}))(y)|^p \mathrm{d}y\right]^{1/p} \le \varepsilon. \tag{12.69}$$

Numerical simulations for deep learning based approximation schemes for nonlinear PDEs can, e.g., be found in [9, 12] and the references mentioned in [1].

# Chapter 13

# Optimization through flows of ordinary differential equations

## 13.1 Introductory comments for the training of artificial neural networks

In this section we briefly sketch how gradient descent type optimization methods factor into machine learning problems. To do this, we now recall the deep supervised learning framework sketched in Section 1.1 above. Let $d, M \in \mathbb{N}$, $\mathcal{E} \in C(\mathbb{R}^d, \mathbb{R})$, $x_1, x_2, \ldots, x_{M+1} \in \mathbb{R}^d$, $y_1, y_2, \ldots, y_M \in \mathbb{R}$ satisfy for all $m \in \{1, 2, \ldots, M\}$ that

$$y_m = \mathcal{E}(x_m), \tag{13.1}$$

and let $\Phi \colon C(\mathbb{R}^d, \mathbb{R}) \to [0, \infty)$ satisfy for all $\phi \in C(\mathbb{R}^d, \mathbb{R})$ that

$$\Phi(\phi) = \sum_{m=1}^{M} |\phi(x_m) - y_m|^2. \tag{13.2}$$

As in Section 1.1 we think of $M \in \mathbb{N}$ as the number of available input-output data pairs, we think of $d \in \mathbb{N}$ as the dimension of the input data, we think of $\mathcal{E} \colon \mathbb{R}^d \to \mathbb{R}$ as an unknown function which relates input and output data through (13.1), we think of $x_1, x_2, \ldots, x_{M+1} \in \mathbb{R}^d$ as the available known input data, we think of $y_1, y_2, \ldots, y_M \in \mathbb{R}$ as the available known output data, and the function $\Phi \colon C(\mathbb{R}^d, \mathbb{R}) \to [0, \infty)$ is the objective function in the optimization problem associated to the supervised learning problem in (13.2) above (cf. (1.2) in Section 1.1 above). In particular, observe that $\Phi(\mathcal{E}) = 0$ and we are trying to approximate the function $\mathcal{E}$ by approximatively computing a global minimizer of the function $\Phi \colon \mathbb{R}^d \to \mathbb{R}$. In order to make this problem amenable to discrete numerical computations, we consider a spatially discretized version of the problem, where we compute minimizers of the function $\Phi$ restricted to a set of realization functions of neural networks. To do this, let $h \in \mathbb{N}$, $l_1, l_2, \ldots, l_h, \mathbf{d} \in \mathbb{N}$ satisfy $\mathbf{d} = l_1(d + 1) + \left[\sum_{k=2}^{h} l_k(l_{k-1} + 1)\right] + l_h + 1$, and let

$$\mathfrak{N} = \left\{ \left(\mathbb{R}^d \ni x \mapsto \mathcal{N}_{\mathfrak{S}_{l_1}, \mathfrak{S}_{l_2}, \ldots, \mathfrak{S}_{l_h}, \mathrm{id}_{\mathbb{R}}}^{\theta, d}(x) \in \mathbb{R}\right) \colon \theta \in \mathbb{R}^{\mathbf{d}}\right\} \subseteq C(\mathbb{R}^d, \mathbb{R}) \tag{13.3}$$

(cf. Definitions 2.1.2 and 2.1.15). We think of $h$ as the number of hidden layers of the neural networks we use as approximators, for every $i \in \{1, 2, \ldots, h\}$ we think of $l_i \in \mathbb{N}$

as the number of neurons in the $i$-th hidden layer of the neural networks we use as approximators, we think of $\mathbf{d}$ as the number of real parameters necessary to describe the neural networks we use as approximators, and we think of $\mathfrak{N}$ as the set of realization functions of the neural networks we use as approximators.

We can now reformulate the optimization problem as the problem of approximately computing minima of the function $f\colon \mathbb{R}^{\mathbf{d}} \to [0,\infty)$ which satisfies for all $\theta \in \mathbb{R}^{\mathbf{d}}$ that

$$f(\theta) = \left[\sum_{m=1}^{M} \left|\left(\mathcal{N}^{\theta,d}_{\mathfrak{S}_{l_1},\mathfrak{S}_{l_2},\ldots,\mathfrak{S}_{l_h},\mathrm{id}_{\mathbb{R}}}\right)(x_m) - y_m\right|^2\right] \tag{13.4}$$

and this optimization is now accessible to discrete numerical computations.

Let $\xi \in \mathbb{R}^d$ and let $\Theta = (\Theta_t)_{t\in[0,\infty)}\colon [0,\infty) \to \mathbb{R}^{\mathbf{d}}$ be a continuously differentiable function which satisfies for all $t \in [0,\infty)$ that

$$\Theta_0 = \xi \qquad \text{and} \qquad \dot{\Theta}_t = -(\nabla f)(\Theta_t). \tag{13.5}$$

Let $(\gamma_n)_{n\in\mathbb{N}} \subseteq [0,\infty)$ and let $\theta = (\theta_n)_{n\in\mathbb{N}_0}\colon \mathbb{N}_0 \to \mathbb{R}^{\mathbf{d}}$ satisfy for all $n \in \mathbb{N}$ that

$$\theta_0 = \xi \qquad \text{and} \qquad \theta_n = \theta_{n-1} - \gamma_n(\nabla f)(\theta_{n-1}). \tag{13.6}$$

## 13.2 Auxiliary results

### 13.2.1 A Gronwall differential inequality

The following lemma, Lemma 13.2.1 below, is referred to as a Gronwall inequality in the literature (cf., e.g., Henry [13, Chapter 7]). Gronwall inequalities are powerful tools to study dynamical systems and, especially, solutions of differential equations.

**Lemma 13.2.1** (Gronwall inequality)**.** *Let $\alpha \in \mathbb{R}$, $T \in (0,\infty)$, $\epsilon \in C^1([0,T],\mathbb{R})$ satisfy for all $t \in (0,T)$ that*

$$\epsilon'(t) \leq \alpha\epsilon(t). \tag{13.7}$$

*Then it holds for all $t \in [0,T]$ that*

$$\epsilon(t) \leq \epsilon(0)e^{\alpha t}. \tag{13.8}$$

*Proof of Lemma 13.2.1.* Throughout this proof let $u\colon [0,T] \to \mathbb{R}$ satisfy for all $t \in [0,T]$ that

$$u(t) = \frac{\epsilon(t)}{e^{\alpha t}} = \epsilon(t)e^{-\alpha t}. \tag{13.9}$$

Observe that the assumption that $\epsilon \in C^1([0,T],\mathbb{R})$ implies that $u \in C^1([0,T],\mathbb{R})$. Moreover, note that (13.7) assures that for all $t \in (0,T)$ it holds that

$$u'(t) = \epsilon'(t)e^{-\alpha t} - \epsilon(t)\alpha e^{-\alpha t} \leq \alpha\epsilon(t)e^{-\alpha t} - \epsilon(t)\alpha e^{-\alpha t} = 0. \tag{13.10}$$

The fundamental theorem of calculus hence demonstrates that for all $t \in [0,T]$ it holds that

$$\frac{\epsilon(t)}{e^{\alpha t}} = u(t) = u(0) + \int_0^t u'(s)\,\mathrm{d}s \leq u(0) + \int_0^t 0\,\mathrm{d}s = u(0) = \epsilon(0). \tag{13.11}$$

Therefore, we obtain for all $t \in [0,T]$ that

$$\epsilon(t) \leq \epsilon(0)e^{\alpha t}. \tag{13.12}$$

The proof of Lemma 13.2.1 is thus complete. $\square$

## 13.2.2   Lyapunov-type functions for ordinary differential equations

**Definition 13.2.2.** *We denote by $\langle \cdot, \cdot \rangle \colon \left[ \bigcup_{d \in \mathbb{N}} (\mathbb{R}^d \times \mathbb{R}^d) \right] \to \mathbb{R}$ the function which satisfies for all $d \in \mathbb{N}$, $x = (x_1, x_2, \ldots, x_d)$, $y = (y_1, y_2, \ldots, y_d) \in \mathbb{R}^d$ that*

$$\langle x, y \rangle = \sum_{i=1}^{d} x_i y_i. \tag{13.13}$$

**Lemma 13.2.3** (Lyapunov-type functions for ordinary differential equations). *Let $d \in \mathbb{N}$, $\alpha \in \mathbb{R}$, $T \in (0, \infty)$, let $O \subseteq \mathbb{R}^d$ be an open set, let $g \in C(O, \mathbb{R}^d)$, $V \in C^1(O, \mathbb{R})$ satisfy for all $\theta \in O$ that*

$$V'(\theta)g(\theta) = \langle (\nabla V)(\theta), g(\theta) \rangle \leq \alpha V(\theta), \tag{13.14}$$

*and let $\Theta \in C([0, T], O)$ satisfy for all $t \in [0, T]$ that $\Theta_t = \Theta_0 + \int_0^t g(\Theta_s) \, \mathrm{d}s$ (cf. Definition 13.2.2). Then it holds for all $t \in [0, T]$ that*

$$V(\Theta_t) \leq e^{\alpha t} V(\Theta_0). \tag{13.15}$$

*Proof of Lemma 13.2.3.* Throughout this proof let $\epsilon \colon [0, T] \to \mathbb{R}$ satisfy for all $t \in [0, T]$ that $\epsilon(t) = V(\Theta_t)$. Observe that the fundamental theorem of calculus, the chain rule, and (13.14) ensure that for all $t \in [0, T]$ it holds that $\epsilon \in C^1([0, T], \mathbb{R})$ and

$$\begin{aligned}
\epsilon'(t) = \tfrac{\mathrm{d}}{\mathrm{d}t}(V(\Theta_t)) &= V'(\Theta_t)\big(\tfrac{\mathrm{d}}{\mathrm{d}t}(\Theta_t)\big) \\
&= V'(\Theta_t)g(\Theta_t) \leq \alpha V(\Theta_t) = \alpha\epsilon(t).
\end{aligned} \tag{13.16}$$

The Gronwall inequality, e.g., in Lemma 13.2.1 (applied with $\alpha \curvearrowright \alpha$, $T \curvearrowright T$, $\epsilon \curvearrowright \epsilon$ in the notation of Lemma 13.2.1) hence demonstrates that for all $t \in [0, T]$ it holds that

$$V(\Theta_t) = \epsilon(t) \leq \epsilon(0)e^{\alpha t} = e^{\alpha t} V(\Theta_0). \tag{13.17}$$

The proof of Lemma 13.2.3 is thus complete. $\qquad\square$

## 13.2.3   On quadratic Lyapunov-type functions and coercivity-type conditions

**Lemma 13.2.4** (Derivative of the standard norm). *Let $d \in \mathbb{N}$, $\vartheta \in \mathbb{R}^d$ and let $f \colon \mathbb{R}^d \to \mathbb{R}$ satisfy for all $\theta \in \mathbb{R}^d$ that*

$$f(\theta) = \|\theta - \vartheta\|_2^2 \tag{13.18}$$

*(cf. Definition 3.1.16). Then it holds for all $\theta \in \mathbb{R}^d$ that $f \in C^\infty(\mathbb{R}^d, \mathbb{R})$ and*

$$(\nabla f)(\theta) = 2(\theta - \vartheta). \tag{13.19}$$

*Proof of Lemma 13.2.4.* Throughout this proof let $\vartheta_1, \vartheta_2, \ldots, \vartheta_d \in \mathbb{R}$ satisfy $\vartheta = (\vartheta_1, \vartheta_2, \ldots, \vartheta_d)$. Note that the fact that for all $\theta = (\theta_1, \theta_2, \ldots, \theta_d) \in \mathbb{R}^d$ it holds that

$$f(\theta) = \sum_{i=1}^{d} |\theta_i - \vartheta_i|^2 \tag{13.20}$$

implies that for all $\theta = (\theta_1, \theta_2, \ldots, \theta_d) \in \mathbb{R}^d$ it holds that $f \in C^\infty(\mathbb{R}^d, \mathbb{R})$ and

$$(\nabla f)(\theta) = \begin{pmatrix} \left(\frac{\partial f}{\partial \theta_1}\right)(\theta) \\ \vdots \\ \left(\frac{\partial f}{\partial \theta_d}\right)(\theta) \end{pmatrix} = \begin{pmatrix} 2(\theta_1 - \vartheta_1) \\ \vdots \\ 2(\theta_d - \vartheta_d) \end{pmatrix} = 2(\theta - \vartheta). \tag{13.21}$$

The proof of Lemma 13.2.4 is thus complete. □

**Corollary 13.2.5** (On quadratic Lyapunov-type functions and coercivity-type conditions). *Let $d \in \mathbb{N}$, $c \in \mathbb{R}$, $T \in (0, \infty)$, $\vartheta \in \mathbb{R}^d$, let $O \subseteq \mathbb{R}^d$ be an open set, let $f \in C^1(O, \mathbb{R})$ satisfy for all $\theta \in O$ that*

$$\langle \theta - \vartheta, (\nabla f)(\theta) \rangle \geq c \|\theta - \vartheta\|_2^2, \tag{13.22}$$

*and let $\Theta \in C([0, T], O)$ satisfy for all $t \in [0, T]$ that $\Theta_t = \Theta_0 - \int_0^t (\nabla f)(\Theta_s) \, ds$ (cf. Definitions 3.1.16 and 13.2.2). Then it holds for all $t \in [0, T]$ that*

$$\|\Theta_t - \vartheta\|_2 \leq e^{-ct} \|\Theta_0 - \vartheta\|_2. \tag{13.23}$$

*Proof of Corollary 13.2.5.* Throughout this proof let $g \colon O \to \mathbb{R}^d$ satisfy for all $\theta \in O$ that

$$g(\theta) = -(\nabla f)(\theta) \tag{13.24}$$

and let $V \colon O \to \mathbb{R}$ satisfy for all $\theta \in O$ that

$$V(\theta) = \|\theta - \vartheta\|_2^2. \tag{13.25}$$

Observe that Lemma 13.2.4 and (13.22) ensure that for all $\theta \in O$ it holds that $V \in C^1(O, \mathbb{R})$ and

$$\begin{aligned} V'(\theta)g(\theta) &= \langle (\nabla V)(\theta), g(\theta) \rangle = \langle 2(\theta - \vartheta), g(\theta) \rangle \\ &= -2\langle (\theta - \vartheta), (\nabla f)(\theta) \rangle \leq -2c\|\theta - \vartheta\|_2^2 = -2cV(\theta). \end{aligned} \tag{13.26}$$

Lemma 13.2.3 hence proves that for all $t \in [0, T]$ it holds that

$$\|\Theta_t - \vartheta\|_2^2 = V(\Theta_t) \leq e^{-2ct} V(\Theta_0) = e^{-2ct} \|\Theta_0 - \vartheta\|_2^2. \tag{13.27}$$

The proof of Corollary 13.2.5 is thus complete. □

### 13.2.4 Sufficient and necessary conditions for local minima

**Lemma 13.2.6.** *Let $d \in \mathbb{N}$, let $O \subseteq \mathbb{R}^d$ be an open set, let $\vartheta \in O$, let $f \colon O \to \mathbb{R}$ be a function, assume that $f$ is differentiable at $\vartheta$, and assume that $(\nabla f)(\vartheta) \neq 0$. Then there exists $\theta \in O$ such that $f(\theta) < f(\vartheta)$.*

*Proof of Lemma 13.2.6.* Throughout this proof let $v \in \mathbb{R}^d \backslash \{0\}$ satisfy $v = -(\nabla f)(\vartheta)$, let $\delta \in (0, \infty)$ satisfy for all $t \in (-\delta, \delta)$ that

$$\vartheta + tv = \vartheta - t(\nabla f)(\vartheta) \in O, \tag{13.28}$$

and let $g \colon (-\delta, \delta) \to \mathbb{R}$ satisfy for all $t \in (-\delta, \delta)$ that

$$g(t) = f(\vartheta + tv). \tag{13.29}$$

Note that for all $t \in (0, \delta)$ it holds that

$$
\begin{aligned}
\left| \left[ \frac{g(t) - g(0)}{t} \right] + \|v\|_2^2 \right| &= \left| \left[ \frac{f(\vartheta + tv) - f(\vartheta)}{t} \right] + \|(\nabla f)(\vartheta)\|_2^2 \right| \\
&= \left| \left[ \frac{f(\vartheta + tv) - f(\vartheta)}{t} \right] + \langle (\nabla f)(\vartheta), (\nabla f)(\vartheta) \rangle \right| \\
&= \left| \left[ \frac{f(\vartheta + tv) - f(\vartheta)}{t} \right] - \langle (\nabla f)(\vartheta), v \rangle \right|.
\end{aligned}
\tag{13.30}
$$

Therefore, we obtain that for all $t \in (0, \delta)$ it holds that

$$
\begin{aligned}
\left| \left[ \frac{g(t) - g(0)}{t} \right] + \|v\|_2^2 \right| &= \left| \left[ \frac{f(\vartheta + tv) - f(\vartheta)}{t} \right] - f'(\vartheta)v \right| \\
&= \left| \frac{f(\vartheta + tv) - f(\vartheta) - f'(\vartheta)tv}{t} \right| = \frac{|f(\vartheta + tv) - f(\vartheta) - f'(\vartheta)tv|}{t}.
\end{aligned}
\tag{13.31}
$$

The assumption that $f$ is differentiable at $\vartheta$ hence demonstrates that

$$
\limsup_{t \searrow 0} \left| \left[ \frac{g(t) - g(0)}{t} \right] + \|v\|_2^2 \right| = 0.
\tag{13.32}
$$

The fact that $\|v\|_2^2 > 0$ therefore demonstrates that there exists $t \in (0, \delta)$ such that

$$
\left| \left[ \frac{g(t) - g(0)}{t} \right] + \|v\|_2^2 \right| < \frac{\|v\|_2^2}{2}.
\tag{13.33}
$$

The triangle inequality and the fact that $\|v\|_2^2 > 0$ hence prove that

$$
\begin{aligned}
\frac{g(t) - g(0)}{t} &= \left[ \frac{g(t) - g(0)}{t} + \|v\|_2^2 \right] - \|v\|_2^2 \le \left| \left[ \frac{g(t) - g(0)}{t} \right] + \|v\|_2^2 \right| - \|v\|_2^2 \\
&< \frac{\|v\|_2^2}{2} - \|v\|_2^2 = -\frac{\|v\|_2^2}{2} < 0.
\end{aligned}
\tag{13.34}
$$

This ensures that

$$
f(\vartheta + tv) = g(t) < g(0) = f(\vartheta).
\tag{13.35}
$$

The proof of Lemma 13.2.6 is thus complete. $\qquad \square$

**Lemma 13.2.7** (A necessary condition for a local minimum). *Let $d \in \mathbb{N}$, let $O \subseteq \mathbb{R}^d$ be an open set, let $\vartheta \in O$, let $f \colon O \to \mathbb{R}$ be a function, assume that $f$ is differentiable at $\vartheta$, and assume*

$$
f(\vartheta) = \inf_{\theta \in O} f(\theta).
\tag{13.36}
$$

*Then $(\nabla f)(\vartheta) = 0$.*

*Proof of Lemma 13.2.7.* We prove Lemma 13.2.7 by contradiction. We thus assume that $(\nabla f)(\vartheta) \neq 0$. Lemma 13.2.6 then implies that there exists $\theta \in O$ such that $f(\theta) < f(\vartheta)$. Combining this with (13.36) shows that

$$
f(\theta) < f(\vartheta) = \inf_{w \in O} f(w) \le f(\theta).
\tag{13.37}
$$

The proof of Lemma 13.2.7 is thus complete. $\qquad \square$

**Lemma 13.2.8** (A sufficient condition for a local minimum)**.** *Let $d \in \mathbb{N}$ and let $c \in (0, \infty)$, $r \in (0, \infty]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{w \in \mathbb{R}^d \colon \|w - \vartheta\|_2 \leq r\}$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ satisfy for all $\theta \in \mathbb{B}$ that*

$$\langle \theta - \vartheta, (\nabla f)(\theta) \rangle \geq c \|\theta - \vartheta\|_2^2 \tag{13.38}$$

*(cf. Definitions 3.1.16 and 13.2.2). Then*

*(i) it holds for all $\theta \in \mathbb{B}$ that $f(\theta) - f(\vartheta) \geq \frac{c}{2} \|\theta - \vartheta\|_2^2$,*

*(ii) it holds that $\{\theta \in \mathbb{B} \colon f(\theta) = \inf_{w \in \mathbb{B}} f(w)\} = \{\vartheta\}$, and*

*(iii) it holds that $(\nabla f)(\vartheta) = 0$.*

*Proof of Lemma 13.2.8.* Throughout this proof let $B$ be the set given by

$$B = \{w \in \mathbb{R}^d \colon \|w - \vartheta\|_2 < r\}. \tag{13.39}$$

Note that (13.38) implies that for all $v \in \mathbb{R}^d$ with $\|v\| \leq r$ it holds that

$$\langle (\nabla f)(\vartheta + v), v \rangle \geq c \|v\|_2^2. \tag{13.40}$$

The fundamental theorem of calculus hence demonstrates that for all $\theta \in \mathbb{B}$ it holds that

$$\begin{aligned}
f(\theta) - f(\vartheta) &= \left[ f(\vartheta + t(\theta - \vartheta)) \right]_{t=0}^{t=1} \\
&= \int_0^1 f'(\vartheta + t(\theta - \vartheta))(\theta - \vartheta) \, dt \\
&= \int_0^1 \langle (\nabla f)(\vartheta + t(\theta - \vartheta)), t(\theta - \vartheta) \rangle \frac{1}{t} \, dt \\
&\geq \int_0^1 c \|t(\theta - \vartheta)\|_2^2 \frac{1}{t} \, dt = c \|\theta - \vartheta\|_2^2 \left[ \int_0^1 t \, dt \right] = \frac{c}{2} \|\theta - \vartheta\|_2^2.
\end{aligned} \tag{13.41}$$

This proves item (i). Next observe that (13.41) ensures that for all $\theta \in \mathbb{B} \backslash \{\vartheta\}$ it holds that

$$f(\theta) \geq f(\vartheta) + \tfrac{c}{2} \|\theta - \vartheta\|^2 > f(\vartheta). \tag{13.42}$$

Hence, we obtain for all $\theta \in \mathbb{B} \backslash \{\vartheta\}$ that

$$\inf_{w \in \mathbb{B}} f(w) = f(\vartheta) < f(\theta). \tag{13.43}$$

This establishes item (ii). It thus remains thus remains to prove item (iii). For this observe that item (ii) ensures that

$$\{\theta \in B \colon f(\theta) = \inf_{w \in B} f(w)\} = \{\vartheta\}. \tag{13.44}$$

Combining this, the fact that $B$ is an open set, and Lemma 13.2.7 (applied with $d \curvearrowright d$, $O \curvearrowright B$, $\vartheta \curvearrowright \vartheta$, $f \curvearrowright f|_B$ in the notation of Lemma 13.2.7) assures that $(\nabla f)(\vartheta) = 0$. This establishes item (iii). The proof of Lemma 13.2.8 is thus complete. $\square$

## 13.2.5   On a linear growth condition

**Lemma 13.2.9** (On a linear growth condition)**.** *Let* $d \in \mathbb{N}$, $L \in (0, \infty)$, $r \in (0, \infty]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{w \in \mathbb{R}^d \colon \|w - \vartheta\|_2 \leq r\}$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ *satisfy for all* $\theta \in \mathbb{B}$ *that*

$$\|(\nabla f)(\theta)\|_2 \leq L\|\theta - \vartheta\|_2 \tag{13.45}$$

*(cf. Definition 3.1.16).   Then it holds for all* $\theta \in \mathbb{B}$ *that*

$$f(\theta) - f(\vartheta) \leq \tfrac{L}{2}\|\theta - \vartheta\|_2^2. \tag{13.46}$$

*Proof of Lemma 13.2.9.* Observe that (13.45), the Cauchy-Schwarz inequality, and the fundamental theorem of calculus ensure that for all $\theta \in \mathbb{B}$ it holds that

$$\begin{aligned}
f(\theta) - f(\vartheta) &= \big[f(\vartheta + t(\theta - \vartheta))\big]_{t=0}^{t=1} \\
&= \int_0^1 f'(\vartheta + t(\theta - \vartheta))(\theta - \vartheta)\,\mathrm{d}t \\
&= \int_0^1 \langle (\nabla f)(\vartheta + t(\theta - \vartheta)), \theta - \vartheta\rangle\,\mathrm{d}t \\
&\leq \int_0^1 \|(\nabla f)(\vartheta + t(\theta - \vartheta))\|_2 \|\theta - \vartheta\|_2\,\mathrm{d}t \\
&\leq \int_0^1 L\|\vartheta + t(\theta - \vartheta) - \vartheta\|_2\|\theta - \vartheta\|_2\,\mathrm{d}t \\
&= L\|\theta - \vartheta\|_2^2\bigg[\int_0^1 t\,\mathrm{d}t\bigg] = \tfrac{L}{2}\|\theta - \vartheta\|_2^2
\end{aligned} \tag{13.47}$$

(cf. Definition 13.2.2).   The proof of Lemma 13.2.9 is thus complete. $\square$

# 13.3   Optimization through flows of ordinary differential equations (ODEs)

## 13.3.1   Approximation of local minima through gradient flows

**Proposition 13.3.1** (Approximation of local minima through gradient flows)**.** *Let* $d \in \mathbb{N}$, $c, T \in (0, \infty)$, $r \in (0, \infty]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{w \in \mathbb{R}^d \colon \|w - \vartheta\|_2 \leq r\}$, $\xi \in \mathbb{B}$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ *satisfy for all* $\theta \in \mathbb{B}$ *that*

$$\langle \theta - \vartheta, (\nabla f)(\theta)\rangle \geq c\|\theta - \vartheta\|_2^2, \tag{13.48}$$

*and let* $\Theta \in C([0, T], \mathbb{R}^d)$ *satisfy for all* $t \in [0, T]$ *that* $\Theta_t = \xi - \int_0^t (\nabla f)(\Theta_s)\,\mathrm{d}s$ *(cf. Definitions 3.1.16 and 13.2.2). Then*

*(i)  it holds that* $\{\theta \in \mathbb{B} \colon f(\theta) = \inf_{w \in \mathbb{B}} f(w)\} = \{\vartheta\}$,

*(ii)  it holds for all* $t \in [0, T]$ *that* $\|\Theta_t - \vartheta\|_2 \leq e^{-ct}\|\xi - \vartheta\|_2$, *and*

*(iii)  it holds for all* $t \in [0, T]$ *that*

$$0 \leq \tfrac{c}{2}\|\Theta_t - \vartheta\|_2^2 \leq f(\Theta_t) - f(\vartheta). \tag{13.49}$$

*Proof of Proposition 13.3.1.* Throughout this proof let $V \colon \mathbb{R}^d \to [0, \infty)$ satisfy for all $\theta \in \mathbb{R}^d$ that $V(\theta) = \|\theta - \vartheta\|_2^2$, let $\epsilon \colon [0, T] \to [0, \infty)$ satisfy for all $t \in [0, T]$ that $\epsilon(t) = \|\Theta_t - \vartheta\|_2^2 = V(\Theta_t)$, and let $\tau \in [0, T]$ be the real number given by

$$\tau = \inf(\{t \in [0, T] \colon \Theta_t \notin \mathbb{B}\} \cup \{T\}) = \inf(\{t \in [0, T] \colon \epsilon(t) > r^2\} \cup \{T\}). \tag{13.50}$$

Note that (13.48) and item (ii) in Lemma 13.2.8 establish item (i). Next observe that Lemma 13.2.4 implies that for all $\theta \in \mathbb{R}^d$ it holds that $V \in C^1(\mathbb{R}^d, [0, \infty))$ and

$$(\nabla V)(\theta) = 2(\theta - \vartheta). \tag{13.51}$$

Moreover, observe that the fundamental theorem of calculus (see, e.g., Coleman [5, Theorem 3.9]) and the fact that $\mathbb{R}^d \ni v \mapsto (\nabla f)(v) \in \mathbb{R}^d$ and $\Theta \colon [0, T] \to \mathbb{R}^d$ are continuous functions ensure that for all $t \in [0, T]$ it holds that $\Theta \in C^1([0, T], \mathbb{R}^d)$ and

$$\tfrac{\mathrm{d}}{\mathrm{d}t}(\Theta_t) = -(\nabla f)(\Theta_t). \tag{13.52}$$

Combining (13.48) and (13.51) hence demonstrates that for all $t \in [0, \tau]$ it holds that $\epsilon \in C^1([0, T], [0, \infty))$ and

$$
\begin{aligned}
\epsilon'(t) &= \tfrac{\mathrm{d}}{\mathrm{d}t}\big(V(\Theta_t)\big) = V'(\Theta_t)\big(\tfrac{\mathrm{d}}{\mathrm{d}t}(\Theta_t)\big) \\
&= \langle (\nabla V)(\Theta_t), \tfrac{\mathrm{d}}{\mathrm{d}t}(\Theta_t)\rangle \\
&= \langle 2(\Theta_t - \vartheta), -(\nabla f)(\Theta_t)\rangle \\
&= -2\langle (\Theta_t - \vartheta), (\nabla f)(\Theta_t)\rangle \\
&\leq -2c\|\Theta_t - \vartheta\|_2^2 = -2c\epsilon(t).
\end{aligned}
\tag{13.53}
$$

The Gronwall inequality, e.g., in Lemma 13.2.1 therefore implies that for all $t \in [0, \tau]$ it holds that

$$\epsilon(t) \leq \epsilon(0)e^{-2ct}. \tag{13.54}$$

Hence, we obtain for all $t \in [0, \tau]$ that

$$\|\Theta_t - \vartheta\|_2 = \sqrt{\epsilon(t)} \leq \sqrt{\epsilon(0)}e^{-ct} = \|\Theta_0 - \vartheta\|_2 e^{-ct} = \|\xi - \vartheta\|_2 e^{-ct}. \tag{13.55}$$

Next note that the assumption that $r \in (0, \infty]$ and the fact that $\epsilon \colon [0, T] \to [0, \infty)$ is a continuous function show that for all $t \in (\epsilon^{-1}(\{0\})) \cap [0, T) = \{s \in [0, T) \colon \epsilon(s) = 0\}$ it holds that

$$\inf(\{s \in [t, T] \colon \epsilon(s) > r^2\} \cup \{T\}) > t. \tag{13.56}$$

Hence, we obtain that for all $t \in (\epsilon^{-1}(\{0\})) \cap \{0\}$ it holds that

$$\tau = \inf(\{s \in [0, T] \colon \epsilon(s) > r^2\} \cup \{T\}) > 0. \tag{13.57}$$

In addition, observe that (13.53) and the assumption that $c \in (0, \infty)$ assure that for all $t \in [0, \tau]$ with $\epsilon(t) > 0$ it holds that

$$\epsilon'(t) \leq -2c\epsilon(t) < 0. \tag{13.58}$$

The fact that $\epsilon' \colon [0, T] \to [0, \infty)$ is a continuous function therefore demonstrates that for all $t \in [0, \tau] \cap [0, T)$ with $\varepsilon(t) > 0$ it holds that

$$\inf(\{u \in [t, T] \colon \epsilon'(u) > 0\} \cup \{T\}) > t. \tag{13.59}$$

This shows that for all $t \in (\epsilon^{-1}((0,\infty))) \cap \{0\}$ it holds that

$$\inf(\{u \in [0,T] \colon \epsilon'(u) > 0\} \cup \{T\}) > 0. \tag{13.60}$$

Next note that the fundamental theorem of calculus and the assumption that $\xi \in \mathbb{B}$ imply that for all $s \in [0,T]$ with $s < \inf(\{u \in [0,T] \colon \epsilon'(u) > 0\} \cup \{T\})$ it holds that

$$\epsilon(s) = \epsilon(0) + \int_0^s \epsilon'(u)\,\mathrm{d}u \leq \epsilon(0) = \|\xi - \vartheta\|_2^2 \leq r^2. \tag{13.61}$$

Combining this with (13.60) establishes that for all $t \in (\epsilon^{-1}((0,\infty))) \cap \{0\}$ it holds that $\tau > 0$. This and (13.57) ensure that

$$\tau > 0. \tag{13.62}$$

Combining this and (13.55) demonstrates that

$$\|\Theta_\tau - \vartheta\|_2 \leq \|\xi - \vartheta\|_2 e^{-c\tau} < r. \tag{13.63}$$

The fact that $\epsilon \colon [0,T] \to [0,\infty)$ is a continuous function and (13.62) hence assure that $\tau = T$. Combining this with (13.55) proves that for all $t \in [0,T]$ it holds that

$$\|\Theta_t - \vartheta\|_2 \leq \|\xi - \vartheta\|_2 e^{-ct}. \tag{13.64}$$

This establishes item (ii). It thus remains to prove item (iii). For this observe that (13.48) and item (i) in Lemma 13.2.8 demonstrate that for all $\theta \in \mathbb{B}$ it holds that

$$0 \leq \tfrac{c}{2}\|\theta - \vartheta\|_2^2 \leq f(\theta) - f(\vartheta). \tag{13.65}$$

Combining this, (13.64), and item (ii) implies that for all $t \in [0,T]$ it holds that

$$0 \leq \tfrac{c}{2}\|\Theta_t - \vartheta\|_2^2 \leq f(\Theta_t) - f(\vartheta) \tag{13.66}$$

This establishes item (iii). The proof of Proposition 13.3.1 is thus complete. $\qquad\square$

## 13.3.2 Existence and uniqueness of solutions of ODEs

**Lemma 13.3.2** (Local existence of maximal solution of ordinary differential equations)**.** *Let $d \in \mathbb{N}$, $\xi \in \mathbb{R}^d$, $T \in (0,\infty)$, let $\|\cdot\| \colon \mathbb{R}^d \to [0,\infty)$ be a norm, and let $g \colon \mathbb{R}^d \to \mathbb{R}^d$ be a locally Lipschitz continuous function. Then there exist a unique real number $\tau \in (0,T]$ and a unique continuous function $\Theta \colon [0,\tau) \to \mathbb{R}^d$ such that for all $t \in [0,\tau)$ it holds that*

$$\liminf_{s \nearrow \tau} \left[ \|\Theta_s\| + \tfrac{1}{(T-s)} \right] = \infty \qquad and \qquad \Theta_t = \xi + \int_0^t g(\Theta_s)\,\mathrm{d}s. \tag{13.67}$$

**Lemma 13.3.3** (Local existence of maximal solution of ordinary differential equations on an infinite time interval)**.** *Let $d \in \mathbb{N}$, $\xi \in \mathbb{R}^d$, let $\|\cdot\| \colon \mathbb{R}^d \to [0,\infty)$ be a norm, and let $g \colon \mathbb{R}^d \to \mathbb{R}^d$ be a locally Lipschitz continuous function. Then there exist a unique extended real number $\tau \in (0,\infty]$ and a unique continuous function $\Theta \colon [0,\tau) \to \mathbb{R}^d$ such that for all $t \in [0,\tau)$ it holds that*

$$\liminf_{s \nearrow \tau} \left[ \|\Theta_s\| + s \right] = \infty \qquad and \qquad \Theta_t = \xi + \int_0^t g(\Theta_s)\,\mathrm{d}s. \tag{13.68}$$

*Proof of Lemma 13.3.3.* First, observe that Lemma 13.3.2 implies that there exist unique real numbers $\tau_n \in (0, n]$, $n \in \mathbb{N}$, and unique continuous functions $\Theta^{(n)} \colon [0, \tau_n) \to \mathbb{R}^d$, $n \in \mathbb{N}$, such that for all $n \in \mathbb{N}$, $t \in [0, \tau_n)$ it holds that

$$\liminf_{s \nearrow \tau_n} \left[ \|\|\Theta_s^{(n)}\|\| + \tfrac{1}{(n-s)} \right] = \infty \qquad \text{and} \qquad \Theta_t^{(n)} = \xi + \int_0^t g(\Theta_s^{(n)}) \, \mathrm{d}s. \tag{13.69}$$

This shows that for all $n \in \mathbb{N}$, $t \in [0, \min\{\tau_{n+1}, n\})$ it holds that

$$\liminf_{s \nearrow \tau_{n+1}} \left[ \|\|\Theta_s^{(n+1)}\|\| + \tfrac{1}{(n+1-s)} \right] = \infty \qquad \text{and} \qquad \Theta_t^{(n+1)} = \xi + \int_0^t g(\Theta_s^{(n+1)}) \, \mathrm{d}s. \tag{13.70}$$

Hence, we obtain that for all $n \in \mathbb{N}$, $t \in [0, \min\{\tau_{n+1}, n\})$ it holds that

$$\liminf_{s \nearrow \min\{\tau_{n+1}, n\}} \left[ \|\|\Theta_s^{(n+1)}\|\| + \tfrac{1}{(n-s)} \right] = \infty \tag{13.71}$$

$$\text{and} \qquad \Theta_t^{(n+1)} = \xi + \int_0^t g(\Theta_s^{(n+1)}) \, \mathrm{d}s. \tag{13.72}$$

Combining this with (13.69) demonstrates that for all $n \in \mathbb{N}$ it holds that

$$\tau_n = \min\{\tau_{n+1}, n\} \qquad \text{and} \qquad \Theta^{(n)} = \Theta^{(n+1)}\big|_{[0, \min\{\tau_{n+1}, n\})}. \tag{13.73}$$

Therefore, we obtain that for all $n \in \mathbb{N}$ it holds that

$$\tau_n \leq \tau_{n+1} \qquad \text{and} \qquad \Theta^{(n)} = \Theta^{(n+1)}\big|_{[0, \tau_n)}. \tag{13.74}$$

Next let $\mathfrak{t} \in (0, \infty]$ be the extended real number given by

$$\mathfrak{t} = \lim_{n \to \infty} \tau_n \tag{13.75}$$

and let $\boldsymbol{\Theta} \colon [0, \mathfrak{t}) \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$, $t \in [0, \tau_n)$ that

$$\boldsymbol{\Theta}_t = \Theta_t^{(n)}. \tag{13.76}$$

Observe that for all $t \in [0, \mathfrak{t})$ there exists $n \in \mathbb{N}$ such that $t \in [0, \tau_n)$. This, (13.69), and (13.74) assure that for all $t \in [0, \mathfrak{t})$ it holds that $\boldsymbol{\Theta} \in C([0, \mathfrak{t}), \mathbb{R}^d)$ and

$$\boldsymbol{\Theta}_t = \xi + \int_0^t g(\boldsymbol{\Theta}_s) \, \mathrm{d}s. \tag{13.77}$$

In addition, note that (13.73) ensures that for all $n \in \mathbb{N}$, $k \in \{n, n+1, n+2, \ldots\}$ it holds that

$$\min\{\tau_{k+1}, n\} = \min\{\tau_{k+1}, k, n\} = \min\{\min\{\tau_{k+1}, k\}, n\} = \min\{\tau_k, n\}. \tag{13.78}$$

This shows that for all $n \in \mathbb{N}$, $k \in \{n+1, n+2, n+3, \ldots\}$ it holds that $\min\{\tau_k, n\} = \min\{\tau_{k-1}, n\}$. Hence, we obtain that for all $n \in \mathbb{N}$, $k \in \{n+1, n+2, n+3, \ldots\}$ it holds that

$$\min\{\tau_k, n\} = \min\{\tau_{k-1}, n\} = \ldots = \min\{\tau_{n+1}, n\} = \min\{\tau_n, n\} = \tau_n. \tag{13.79}$$

Combining this with the fact that $(\tau_n)_{n\in\mathbb{N}} \subseteq [0,\infty)$ is a non-decreasing sequence implies that for all $n \in \mathbb{N}$ it holds that

$$\min\{\mathfrak{t}, n\} = \min\left\{\lim_{k\to\infty}\tau_k, n\right\} = \lim_{k\to\infty}\left(\min\{\tau_k, n\}\right) = \lim_{k\to\infty}\tau_n = \tau_n. \tag{13.80}$$

Therefore, we obtain that for all $n \in \mathbb{N}$ with $\mathfrak{t} < n$ it holds that

$$\tau_n = \min\{\mathfrak{t}, n\} = \mathfrak{t}. \tag{13.81}$$

This, (13.69), and (13.76) demonstrate that for all $n \in \mathbb{N}$ with $\mathfrak{t} < n$ it holds that

$$\liminf_{s\nearrow\mathfrak{t}}\|\!|\!|\boldsymbol{\Theta}_s|\!|\!| = \liminf_{s\nearrow\tau_n}\|\!|\!|\boldsymbol{\Theta}_s|\!|\!| = \liminf_{s\nearrow\tau_n}\|\!|\!|\Theta_s^{(n)}|\!|\!|$$
$$= -\tfrac{1}{(n-\mathfrak{t})} + \liminf_{s\nearrow\tau_n}\left[\|\!|\!|\Theta_s^{(n)}|\!|\!| + \tfrac{1}{(n-\mathfrak{t})}\right] \tag{13.82}$$
$$= -\tfrac{1}{(n-\mathfrak{t})} + \liminf_{s\nearrow\tau_n}\left[\|\!|\!|\Theta_s^{(n)}|\!|\!| + \tfrac{1}{(n-s)}\right] = \infty.$$

Therefore, we obtain that

$$\liminf_{s\nearrow\mathfrak{t}}\left[\|\!|\!|\boldsymbol{\Theta}_s|\!|\!| + s\right] = \infty. \tag{13.83}$$

Next note that for all $\hat{\mathfrak{t}} \in (0,\infty]$, $\hat{\boldsymbol{\Theta}} \in C([0,\hat{\mathfrak{t}}),\mathbb{R}^d)$, $n \in \mathbb{N}$, $t \in [0,\min\{\hat{\mathfrak{t}},n\})$ with $\liminf_{s\nearrow\hat{\mathfrak{t}}}[\|\!|\!|\hat{\boldsymbol{\Theta}}_s|\!|\!| + s] = \infty$ and $\forall s \in [0,\hat{\mathfrak{t}})\colon \hat{\boldsymbol{\Theta}}_s = \xi + \int_0^s g(\hat{\boldsymbol{\Theta}}_u)\,\mathrm{d}u$ it holds that

$$\liminf_{s\nearrow\min\{\hat{\mathfrak{t}},n\}}\left[\|\!|\!|\hat{\boldsymbol{\Theta}}_s|\!|\!| + \tfrac{1}{(n-s)}\right] = \infty \qquad \text{and} \qquad \hat{\boldsymbol{\Theta}}_t = \xi + \int_0^t g(\hat{\boldsymbol{\Theta}}_s)\,\mathrm{d}s. \tag{13.84}$$

This and (13.69) demonstrate that for all $\hat{\mathfrak{t}} \in (0,\infty]$, $\hat{\boldsymbol{\Theta}} \in C([0,\hat{\mathfrak{t}}),\mathbb{R}^d)$, $n \in \mathbb{N}$ with $\liminf_{t\nearrow\hat{\mathfrak{t}}}[\|\!|\!|\hat{\boldsymbol{\Theta}}_t|\!|\!| + t] = \infty$ and $\forall t \in [0,\hat{\mathfrak{t}})\colon \hat{\boldsymbol{\Theta}}_t = \xi + \int_0^t g(\hat{\boldsymbol{\Theta}}_s)\,\mathrm{d}s$ it holds that

$$\min\{\hat{\mathfrak{t}}, n\} = \tau_n \qquad \text{and} \qquad \hat{\boldsymbol{\Theta}}|_{[0,\tau_n)} = \Theta^{(n)}. \tag{13.85}$$

Combining (13.77) and (13.83) hence assures that for all $\hat{\mathfrak{t}} \in (0,\infty]$, $\hat{\boldsymbol{\Theta}} \in C([0,\hat{\mathfrak{t}}),\mathbb{R}^d)$, $n \in \mathbb{N}$ with $\liminf_{t\nearrow\hat{\mathfrak{t}}}[\|\!|\!|\hat{\boldsymbol{\Theta}}_t|\!|\!| + t] = \infty$ and $\forall t \in [0,\hat{\mathfrak{t}})\colon \hat{\boldsymbol{\Theta}}_t = \xi + \int_0^t g(\hat{\boldsymbol{\Theta}}_s)\,\mathrm{d}s$ it holds that

$$\min\{\hat{\mathfrak{t}}, n\} = \tau_n = \min\{\mathfrak{t}, n\} \qquad \text{and} \qquad \hat{\boldsymbol{\Theta}}|_{[0,\tau_n)} = \Theta^{(n)} = \boldsymbol{\Theta}|_{[0,\tau_n)}. \tag{13.86}$$

This and (13.75) show that for all $\hat{\mathfrak{t}} \in (0,\infty]$, $\hat{\boldsymbol{\Theta}} \in C([0,\hat{\mathfrak{t}}),\mathbb{R}^d)$ with $\liminf_{t\nearrow\hat{\mathfrak{t}}}[\|\!|\!|\hat{\boldsymbol{\Theta}}_t|\!|\!| + t] = \infty$ and $\forall t \in [0,\hat{\tau})\colon \hat{\boldsymbol{\Theta}}_t = \xi + \int_0^t g(\hat{\boldsymbol{\Theta}}_s)\,\mathrm{d}s$ it holds that

$$\hat{\mathfrak{t}} = \mathfrak{t} \qquad \text{and} \qquad \hat{\boldsymbol{\Theta}} = \boldsymbol{\Theta}. \tag{13.87}$$

Combining this, (13.77), and (13.83) completes the proof of Lemma 13.3.3. $\qquad\square$

### 13.3.3 Approximation of local minima through gradient flows revisited

**Theorem 13.3.4** (Approximation of local minima through gradient flows revisited)**.** *Let* $d \in \mathbb{N}$ *and let* $c \in (0,\infty)$, $r \in (0,\infty]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{w \in \mathbb{R}^d\colon \|w - \vartheta\|_2 \leq r\}$, $\xi \in \mathbb{B}$, $f \in C^2(\mathbb{R}^d, \mathbb{R})$ *satisfy for all* $\theta \in \mathbb{B}$ *that*

$$\langle \theta - \vartheta, (\nabla f)(\theta)\rangle \geq c\|\theta - \vartheta\|_2^2 \tag{13.88}$$

*(cf. Definitions 3.1.16 and 13.2.2). Then*

(i) *there exists a unique continuous function* $\Theta\colon [0,\infty) \to \mathbb{R}^d$ *such that for all* $t \in [0,\infty)$ *it holds that*

$$\Theta_t = \xi - \int_0^t (\nabla f)(\Theta_s)\,\mathrm{d}s, \tag{13.89}$$

(ii) *it holds that* $\{\theta \in \mathbb{B}\colon f(\theta) = \inf_{w\in\mathbb{B}} f(w)\} = \{\vartheta\}$,

(iii) *it holds for all* $t \in [0,\infty)$ *that* $\|\Theta_t - \vartheta\|_2 \le e^{-ct}\|\xi - \vartheta\|_2$, *and*

(iv) *it holds for all* $t \in [0,\infty)$ *that*

$$0 \le \tfrac{c}{2}\|\Theta_t - \vartheta\|_2^2 \le f(\Theta_t) - f(\vartheta). \tag{13.90}$$

*Proof of Theorem 13.3.4.* First, observe that the assumption that $f \in C^2(\mathbb{R}^d, \mathbb{R})$ ensures that $\mathbb{R}^d \ni \theta \mapsto -(\nabla f)(\theta) \in \mathbb{R}^d$ is a continuously differentiable function. The fundamental theorem of calculus hence implies that $\mathbb{R}^d \ni \theta \mapsto -(\nabla f)(\theta) \in \mathbb{R}^d$ is a locally Lipschitz continuous function. Combining this with Lemma 13.3.3 (applied with $g(\theta) \curvearrowright -(\nabla f)(\theta)$ for $\theta \in \mathbb{R}^d$ in the notation of Lemma 13.3.3) proves that there exists a unique extended real number $\tau \in (0,\infty]$ and a unique continuous function $\Theta\colon [0,\tau) \to \mathbb{R}^d$ such that for all $t \in [0,\tau)$ it holds that

$$\liminf_{s\nearrow\tau}\big[\|\Theta_s\|_2 + s\big] = \infty \qquad \text{and} \qquad \Theta_t = \xi - \int_0^t (\nabla f)(\Theta_s)\,\mathrm{d}s. \tag{13.91}$$

Next observe that Proposition 13.3.1 proves that for all $t \in [0,\tau)$ it holds that

$$\|\Theta_t - \vartheta\|_2 \le e^{-ct}\|\xi - \vartheta\|_2. \tag{13.92}$$

This implies that

$$\begin{aligned}
\liminf_{s\nearrow\tau}\|\Theta_s\|_2 &\le \Big[\liminf_{s\nearrow\tau}\|\Theta_s - \vartheta\|_2\Big] + \|\vartheta\|_2 \\
&\le \Big[\liminf_{s\nearrow\tau} e^{-ct}\|\xi - \vartheta\|_2\Big] + \|\vartheta\|_2 \le \|\xi - \vartheta\|_2 + \|\vartheta\|_2 < \infty.
\end{aligned} \tag{13.93}$$

This and (13.91) demonstrate that $\tau = \infty$. This proves item (i). Moreover, note that Proposition 13.3.1 and item (i) establish items (ii), (iii), and (iv). The proof of Theorem 13.3.4 is thus complete. $\qquad\square$

### 13.3.4 Approximation error with respect to the objective function

**Corollary 13.3.5** (Approximation error with respect to the objective function)**.** *Let* $d \in \mathbb{N}$ *and let* $c, L \in (0,\infty)$, $r \in (0,\infty]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{w \in \mathbb{R}^d\colon \|w - \vartheta\|_2 \le r\}$, $\xi \in \mathbb{B}$, $f \in C^2(\mathbb{R}^d, \mathbb{R})$ *satisfy for all* $\theta \in \mathbb{B}$ *that*

$$\langle \theta - \vartheta, (\nabla f)(\theta)\rangle \ge c\|\theta - \vartheta\|_2^2 \qquad \text{and} \qquad \|(\nabla f)(\theta)\|_2 \le L\|\theta - \vartheta\|_2 \tag{13.94}$$

*(cf. Definitions 3.1.16 and 13.2.2). Then*

(i) *there exists a unique continuous function* $\Theta\colon [0,\infty) \to \mathbb{R}^d$ *such that for all* $t \in [0,\infty)$ *it holds that*

$$\Theta_t = \xi - \int_0^t (\nabla f)(\Theta_s)\, ds, \tag{13.95}$$

(ii) *it holds that* $\{\theta \in \mathbb{B}\colon f(\theta) = \inf_{w \in \mathbb{B}} f(w)\} = \{\vartheta\}$,

(iii) *it holds for all* $t \in [0,\infty)$ *that* $\|\Theta_t - \vartheta\|_2 \le e^{-ct}\|\xi - \vartheta\|_2$, *and*

(iv) *it holds for all* $t \in [0,\infty)$ *that*

$$0 \le \tfrac{c}{2}\|\Theta_t - \vartheta\|_2^2 \le f(\Theta_t) - f(\vartheta) \le \tfrac{L}{2}\|\Theta_t - \vartheta\|_2^2 \le \tfrac{L}{2}e^{-2ct}\|\xi - \vartheta\|_2^2. \tag{13.96}$$

*Proof of Corollary 13.3.5.* Theorem 13.3.4 and Lemma 13.2.9 establish items (i)–(iv). The proof of Corollary 13.3.5 is thus complete. □

# Chapter 14

# Deterministic gradient descent type optimization methods

## 14.1 The gradient descent optimization method

In this section we review and study the classical plain vanilla GD optimization method (cf., for example, Nesterov [23, Section 1.2.3], Boyd & Vandenberghe [2, Section 9.3], and Bubeck [3, Chapter 3]). A simple intuition behind the GD optimization method is the idea to solve a minimization problem by performing successive steps in direction of the steepest descents of the objective function, that is, by performing successive steps in the opposite direction of the gradients of the objective function. A slightly different and maybe a bit more accurate perspective for the GD optimization method is to view the GD optimization method as a plain vanilla Euler discretization of the gradient flow ODE in Theorem 13.3.4 in Chapter 13.

**Definition 14.1.1** (Gradient descent optimization method). *Let $d \in \mathbb{N}$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \infty)$, $\xi \in \mathbb{R}^d$ and let $f \colon \mathbb{R}^d \to \mathbb{R}$ and $g \colon \mathbb{R}^d \to \mathbb{R}^d$ satisfy for all $\theta \in \{v \in \mathbb{R}^d \colon (f$ is differentiable at $v)\}$ that*

$$g(\theta) = (\nabla f)(\theta). \tag{14.1}$$

*Then we say that $\Theta$ is the gradient descent process for the objective function $f$ with generalized gradient $g$, learning rates $(\gamma_n)_{n \in \mathbb{N}}$, and initial value $\xi$ (we say that $\Theta$ is the gradient descent process for the objective function $f$ with learning rates $(\gamma_n)_{n \in \mathbb{N}}$ and initial value $\xi$) if and only if it holds that $\Theta \colon \mathbb{N}_0 \to \mathbb{R}^d$ is the function from $\mathbb{N}_0$ to $\mathbb{R}^d$ which satisfies for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \gamma_n g(\Theta_{n-1}). \tag{14.2}$$

### 14.1.1 Lyapunov-type stability for GD type optimization methods

Lemma 13.2.3 in Subsection 13.2.2 and Corollary 13.2.5 in Subsection 13.2.3 in Chapter 13 above, in particular, illustrate how Lyapunov-type functions can be employed to establish convergence properties for gradient flows. The next two results, Proposition 14.1.2 and Corollary 14.1.3 below, are, roughly speaking, the time-discrete anologon of Lemma 13.2.3 and Corollary 13.2.5, respectively.

**Proposition 14.1.2** (Lyapunov-type stability for discrete-time dynamical systems). *Let $d \in \mathbb{N}$, $\xi \in \mathbb{R}^d$, $c \in (0, \infty)$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, c]$, let $V : \mathbb{R}^d \to \mathbb{R}$, $\Phi : \mathbb{R}^d \times [0, \infty) \to \mathbb{R}^d$, and $\varepsilon : [0, c] \to [0, \infty)$ satisfy for all $\theta \in \mathbb{R}^d$, $t \in [0, c]$ that*

$$V(\Phi(\theta, t)) \leq \varepsilon(t) V(\theta), \tag{14.3}$$

*and let $\Theta : \mathbb{N}_0 \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Phi(\Theta_{n-1}, \gamma_n). \tag{14.4}$$

*Then it holds for all $n \in \mathbb{N}_0$ that*

$$V(\Theta_n) \leq \left[ \prod_{k=1}^{n} \varepsilon(\gamma_k) \right] V(\xi). \tag{14.5}$$

*Proof of Proposition 14.1.2.* We prove (14.5) by induction on $n \in \mathbb{N}_0$. For the base case $n = 0$ note that the assumption that $\Theta_0 = \xi$ ensures that $V(\Theta_0) = V(\xi)$. This establishes (14.5) in the base case $n = 0$. For the induction step observe (14.4) and (14.3) ensure that for all $n \in \mathbb{N}_0$ with $V(\Theta_n) \leq (\prod_{k=1}^{n} \varepsilon(\gamma_k)) V(\xi)$ it holds that

$$\begin{aligned} V(\Theta_{n+1}) = V(\Phi(\Theta_n, \gamma_{n+1})) &\leq \varepsilon(\gamma_{n+1}) V(\Theta_n) \\ &\leq \varepsilon(\gamma_{n+1}) \left( \left[ \prod_{k=1}^{n} \varepsilon(\gamma_k) \right] V(\xi) \right) = \left[ \prod_{k=1}^{n+1} \varepsilon(\gamma_k) \right] V(\xi). \end{aligned} \tag{14.6}$$

Induction thus establishes (14.5). This completes the proof of Proposition 14.1.2. $\qquad \square$

**Corollary 14.1.3** (On quadratic Lyapunov-type functions for the GD optimization method). *Let $d \in \mathbb{N}$, $c \in (0, \infty)$, $\vartheta, \xi \in \mathbb{R}^d$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, c]$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$, let $\|\cdot\| : \mathbb{R}^d \to [0, \infty)$ be a norm, let $\varepsilon : [0, c] \to [0, \infty)$ satisfy for all $\theta \in \mathbb{R}^d$, $t \in [0, c]$ that*

$$\|\theta - t(\nabla f)(\theta) - \vartheta\|^2 \leq \varepsilon(t) \|\theta - \vartheta\|^2, \tag{14.7}$$

*and let $\Theta : \mathbb{N}_0 \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \gamma_n(\nabla f)(\Theta_{n-1}). \tag{14.8}$$

*Then it holds for all $n \in \mathbb{N}_0$ that*

$$\|\Theta_n - \vartheta\| \leq \left[ \prod_{k=1}^{n} [\varepsilon(\gamma_k)]^{1/2} \right] \|\xi - \vartheta\|. \tag{14.9}$$

*Proof of Corollary 14.1.3.* Throughout this proof let $V : \mathbb{R}^d \to \mathbb{R}$ satisfy for all $\theta \in \mathbb{R}^d$ that

$$V(\theta) = \|\theta - \vartheta\|^2. \tag{14.10}$$

Observe that Proposition 14.1.2 (applied with $V \curvearrowleft V$ in the notation of Proposition 14.1.2) implies that for all $n \in \mathbb{N}_0$ it holds that

$$\|\Theta_n - \vartheta\|^2 = V(\Theta_n) \leq \left[ \prod_{k=1}^{n} \varepsilon(\gamma_k) \right] V(\xi) = \left[ \prod_{k=1}^{n} \varepsilon(\gamma_k) \right] \|\xi - \vartheta\|^2. \tag{14.11}$$

This establishes (14.9). The proof of Corollary 14.1.3 is thus complete. $\qquad \square$

Corollary 14.1.3, in particular, illustrates that the one-step Lyapunov stability assumption in (14.7) may provide us suitable estimates for the approximation errors associated to the GD optimization method; see (14.9) above. The next result, Lemma 14.1.4 below, now provides us sufficient conditions which ensure that the one-step Lyapunov stability condition in (14.7) is satisfied so that we are in the position to apply Corollary 14.1.3 above to obtain estimates for the approximation errors associated to the GD optimization method. Lemma 14.1.4 employs the growth condition and the coercivity-type condition in (13.94) in Corollary 13.3.5 above. Results similar to Lemma 14.1.4 can, e.g., be found in Dereich & Müller-Gronbach [7, Remark 2.1] and Jentzen et al. [16, Lemma 2.1]. We will employ the statement of Lemma 14.1.4 in our error analysis for the GD optimization method in Subsection 14.1.2 below.

**Lemma 14.1.4** (Sufficient conditions for a one-step Lyapunov-type stability condition). *Let $d \in \mathbb{N}$, let $\langle\!\langle \cdot, \cdot \rangle\!\rangle \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ be a scalar product, let $\|\!\|\cdot\|\!\| \colon \mathbb{R}^d \to [0, \infty)$ satisfy for all $v \in \mathbb{R}^d$ that $\|\!\|v\|\!\| = \sqrt{\langle\!\langle v, v \rangle\!\rangle}$, and let $c, L \in (0, \infty)$, $r \in (0, \infty]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{w \in \mathbb{R}^d \colon \|\!\|w - \vartheta\|\!\| \leq r\}$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ satisfy for all $\theta \in \mathbb{B}$ that*

$$\langle\!\langle \theta - \vartheta, (\nabla f)(\theta) \rangle\!\rangle \geq c\|\!\|\theta - \vartheta\|\!\|^2 \qquad and \qquad \|\!\|(\nabla f)(\theta)\|\!\| \leq L\|\!\|\theta - \vartheta\|\!\|. \qquad (14.12)$$

*Then*

(i) *it holds that $c \leq L$,*

(ii) *it holds for all $\theta \in \mathbb{B}$, $\gamma \in [0, \infty)$ that*

$$\|\!\|\theta - \gamma(\nabla f)(\theta) - \vartheta\|\!\|^2 \leq (1 - 2\gamma c + \gamma^2 L^2)\|\!\|\theta - \vartheta\|\!\|^2, \qquad (14.13)$$

(iii) *it holds for all $\gamma \in (0, \frac{2c}{L^2})$ that $0 \leq 1 - 2\gamma c + \gamma^2 L^2 < 1$, and*

(iv) *it holds for all $\theta \in \mathbb{B}$, $\gamma \in [0, \frac{c}{L^2}]$ that*

$$\|\!\|\theta - \gamma(\nabla f)(\theta) - \vartheta\|\!\|^2 \leq (1 - c\gamma)\|\!\|\theta - \vartheta\|\!\|^2. \qquad (14.14)$$

*Proof of Lemma 14.1.4.* First of all, note that (14.12) ensures that for all $\theta \in \mathbb{B}$, $\gamma \in [0, \infty)$ it holds that

$$\begin{aligned}
\|\!\|\theta - \gamma(\nabla f)(\theta) - \vartheta\|\!\|^2 &= \|\!\|(\theta - \vartheta) - \gamma(\nabla f)(\theta)\|\!\|^2 \\
&= \|\!\|\theta - \vartheta\|\!\|^2 - 2\gamma \langle\!\langle \theta - \vartheta, (\nabla f)(\theta) \rangle\!\rangle + \gamma^2 \|\!\|(\nabla f)(\theta)\|\!\|^2 \\
&\leq \|\!\|\theta - \vartheta\|\!\|^2 - 2\gamma c\|\!\|\theta - \vartheta\|\!\|^2 + \gamma^2 L^2 \|\!\|\theta - \vartheta\|\!\|^2 \\
&= (1 - 2\gamma c + \gamma^2 L^2)\|\!\|\theta - \vartheta\|\!\|^2.
\end{aligned} \qquad (14.15)$$

This establishes item (ii). Moreover, note that the fact that $\mathbb{B} \backslash \{\vartheta\} \neq \emptyset$ and (14.15) assure that for all $\gamma \in [0, \infty)$ it holds that

$$1 - 2\gamma c + \gamma^2 L^2 \geq 0. \qquad (14.16)$$

Hence, we obtain that

$$\begin{aligned}
1 - \frac{c^2}{L^2} = 1 - \frac{2c^2}{L^2} + \frac{c^2}{L^2} &= 1 - 2\left[\frac{c}{L^2}\right]c + \left[\frac{c^2}{L^4}\right]L^2 \\
&= 1 - 2\left[\frac{c}{L^2}\right]c + \left[\frac{c}{L^2}\right]^2 L^2 \geq 0.
\end{aligned} \qquad (14.17)$$

This implies that $\frac{c^2}{L^2} \leq 1$. Therefore, we obtain that $c^2 \leq L^2$. This establishes item (i). Furthermore, observe that (14.16) ensures that for all $\gamma \in (0, \frac{2c}{L^2})$ it holds that

$$0 \leq 1 - 2\gamma c + \gamma^2 L^2 = 1 - \underbrace{\gamma}_{>0} \underbrace{(2c - \gamma L^2)}_{>0} < 1. \tag{14.18}$$

This proves item (iii). In addition, note that for all $\gamma \in [0, \frac{c}{L^2}]$ it holds that

$$1 - 2\gamma c + \gamma^2 L^2 \leq 1 - 2\gamma c + \gamma\left[\tfrac{c}{L^2}\right]L^2 = 1 - c\gamma. \tag{14.19}$$

Combining this with (14.15) establishes item (iv). The proof of Lemma 14.1.4 is thus complete. $\qquad\square$

**Exercise 14.1.1.** *Prove or disprove the following statement: There exist $d \in \mathbb{N}$, $\gamma \in (0, \infty)$, $\varepsilon \in (0, 1)$, $r \in (0, \infty]$, $\vartheta, \theta \in \mathbb{R}^d$ and there exists a function $g\colon \mathbb{R}^d \to \mathbb{R}^d$ such that $\|\theta - \vartheta\|_2 \leq r$, $\forall\, \xi \in \{w \in \mathbb{R}^d\colon \|w - \vartheta\|_2 \leq r\}\colon \|\xi - \gamma g(\xi) - \vartheta\|_2 \leq \varepsilon\|\xi - \vartheta\|_2$, and*

$$\langle \theta - \vartheta, g(\theta)\rangle < \min\left\{\tfrac{1-\varepsilon^2}{2\gamma}, \tfrac{\gamma}{2}\right\} \max\left\{\|\theta - \vartheta\|_2^2, \|g(\theta)\|_2^2\right\}. \tag{14.20}$$

**Exercise 14.1.2.** *Prove or disprove the following statement: For all $d \in \mathbb{N}$, $r \in (0, \infty]$, $\vartheta \in \mathbb{R}^d$ and for every function $g\colon \mathbb{R}^d \to \mathbb{R}^d$ which satisfies $\forall\, \theta \in \{w \in \mathbb{R}^d\colon \|w - \vartheta\|_2 \leq r\}\colon \langle \theta - \vartheta, g(\theta)\rangle \geq \frac{1}{2}\max\{\|\theta - \vartheta\|_2^2, \|g(\theta)\|_2^2\}$ it holds that*

$$\forall\, \theta \in \{w \in \mathbb{R}^d\colon \|w - \vartheta\|_2 \leq r\}\colon \big(\langle \theta - \vartheta, g(\theta)\rangle \geq \tfrac{1}{2}\|\theta - \vartheta\|_2^2 \wedge \|g(\theta)\|_2 \leq 2\|\theta - \vartheta\|_2\big). \tag{14.21}$$

**Exercise 14.1.3.** *Prove or disprove the following statement: For all $d \in \mathbb{N}$, $c \in (0, \infty)$, $r \in (0, \infty]$, $\vartheta, v \in \mathbb{R}^d$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$, $s, t \in [0, 1]$ such that $\|v\|_2 \leq r$, $s \leq t$, and $\forall\, \theta \in \{w \in \mathbb{R}^d\colon \|w - \vartheta\|_2 \leq r\}\colon \langle \theta - \vartheta, (\nabla f)(\theta)\rangle \geq c\|\theta - \vartheta\|_2^2$ it holds that*

$$f(\vartheta + tv) - f(\vartheta + sv) \geq \tfrac{c}{2}(t^2 - s^2)\|v\|_2^2. \tag{14.22}$$

**Exercise 14.1.4.** *Prove or disprove the following statement: For every $d \in \mathbb{N}$, $c \in (0, \infty)$, $r \in (0, \infty]$, $\vartheta \in \mathbb{R}^d$ and for every $f \in C^1(\mathbb{R}^d, \mathbb{R})$ which satisfies for all $v \in \mathbb{R}^d$, $s, t \in [0, 1]$ with $\|v\|_2 \leq r$ and $s \leq t$ that $f(\vartheta + tv) - f(\vartheta + sv) \geq c(t^2 - s^2)\|v\|_2^2$ it holds that*

$$\forall\, \theta \in \{w \in \mathbb{R}^d\colon \|w - \vartheta\|_2 \leq r\}\colon \langle \theta - \vartheta, (\nabla f)(\theta)\rangle \geq 2c\|\theta - \vartheta\|_2^2. \tag{14.23}$$

**Exercise 14.1.5.** *Let $d \in \mathbb{N}$ and for every $v \in \mathbb{R}^d$, $R \in [0, \infty]$ let $\mathbb{B}_R(v) = \{w \in \mathbb{R}^d\colon \|w - v\|_2 \leq R\}$. Prove or disprove the following statement: For all $r \in (0, \infty]$, $\vartheta \in \mathbb{R}^d$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ the following two statements are equivalent:*

*(i) There exists $c \in (0, \infty)$ such that for all $\theta \in \mathbb{B}_r(\vartheta)$ it holds that*

$$\langle \theta - \vartheta, (\nabla f)(\theta)\rangle \geq c\|\theta - \vartheta\|_2^2. \tag{14.24}$$

*(ii) There exists $c \in (0, \infty)$ such that for all $v, w \in \mathbb{B}_r(\vartheta)$, $s, t \in [0, 1]$ with $s \leq t$ it holds that*

$$f(\vartheta + t(v - \vartheta)) - f(\vartheta + s(v - \vartheta)) \geq c(t^2 - s^2)\|v - \vartheta\|_2^2. \tag{14.25}$$

**Exercise 14.1.6.** *Let $d \in \mathbb{N}$ and for every $v \in \mathbb{R}^d$, $R \in [0, \infty]$ let $\mathbb{B}_R(v) = \{w \in \mathbb{R}^d\colon \|v - w\|_2 \leq R\}$. Prove or disprove the following statement: For all $r \in (0, \infty]$, $\vartheta \in \mathbb{R}^d$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ the following three statements are equivalent:*

(i) *There exist $c, L \in (0, \infty)$ such that for all $\theta \in \mathbb{B}_r(\vartheta)$ it holds that*

$$\langle \theta - \vartheta, (\nabla f)(\theta) \rangle \geq c \|\theta - \vartheta\|_2^2 \qquad and \qquad \|(\nabla f)(\theta)\|_2 \leq L \|\theta - \vartheta\|_2. \qquad (14.26)$$

(ii) *There exist $\gamma \in (0, \infty)$, $\varepsilon \in (0, 1)$ such that for all $\theta \in \mathbb{B}_r(\vartheta)$ it holds that*

$$\|\theta - \gamma(\nabla f)(\theta) - \vartheta\|_2 \leq \varepsilon \|\theta - \vartheta\|_2. \qquad (14.27)$$

(iii) *There exists $c \in (0, \infty)$ such that for all $\theta \in \mathbb{B}_r(\vartheta)$ it holds that*

$$\langle \theta - \vartheta, (\nabla f)(\theta) \rangle \geq c \max\{ \|\theta - \vartheta\|_2^2, \|(\nabla f)(\theta)\|_2^2 \}. \qquad (14.28)$$

## 14.1.2 Error analysis for the GD optimization method

In this subsection we provide an error analysis for the GD optimization method. In particular, we show under suitable hypotheses (cf. Proposition 14.1.5 below) that the GD optimization method (cf. Definition 14.1.1 above) converges to a local minimum of the objective function of the considered optimization problem.

### 14.1.2.1 Error estimates for the GD optimization method

**Proposition 14.1.5** (Error estimates for the GD optimization method). *Let $d \in \mathbb{N}$, $c, L \in (0, \infty)$, $r \in (0, \infty]$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \frac{2c}{L^2}]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{w \in \mathbb{R}^d \colon \|w - \vartheta\| \leq r\}$, $\xi \in \mathbb{B}$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ satisfy for all $\theta \in \mathbb{B}$ that*

$$\langle \theta - \vartheta, (\nabla f)(\theta) \rangle \geq c \|\theta - \vartheta\|_2^2 \qquad and \qquad \|(\nabla f)(\theta)\|_2 \leq L \|\theta - \vartheta\|_2, \qquad (14.29)$$

*and let $\Theta \colon \mathbb{N}_0 \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \gamma_n (\nabla f)(\Theta_{n-1}). \qquad (14.30)$$

*Then*

(i) *it holds that $\{\theta \in \mathbb{B} \colon f(\theta) = \inf_{w \in \mathbb{B}} f(w)\} = \{\vartheta\}$,*

(ii) *it holds for all $n \in \mathbb{N}$ that $0 \leq 1 - 2c\gamma_n + (\gamma_n)^2 L^2 \leq 1$,*

(iii) *it holds for all $n \in \mathbb{N}$ that $\|\Theta_n - \vartheta\|_2 \leq (1 - 2c\gamma_n + (\gamma_n)^2 L^2)^{1/2} \|\Theta_{n-1} - \vartheta\|_2 \leq r$,*

(iv) *it holds for all $n \in \mathbb{N}_0$ that*

$$\|\Theta_n - \vartheta\|_2 \leq \left[ \prod_{k=1}^n (1 - 2c\gamma_k + (\gamma_k)^2 L^2)^{1/2} \right] \|\xi - \vartheta\|_2, \qquad (14.31)$$

*and*

(v) *it holds for all $n \in \mathbb{N}_0$ that*

$$0 \leq f(\Theta_n) - f(\vartheta) \leq \tfrac{L}{2} \|\Theta_n - \vartheta\|_2^2 \leq \tfrac{L}{2} \left[ \prod_{k=1}^n (1 - 2c\gamma_k + (\gamma_k)^2 L^2) \right] \|\xi - \vartheta\|_2^2. \quad (14.32)$$

*Proof of Proposition 14.1.5.* First, note that (14.29) and item (ii) in Lemma 13.2.8 prove item (i). Moreover, observe that (14.29), item (iii) in Lemma 14.1.4, the assumption that for all $n \in \mathbb{N}$ it holds that $\gamma_n \in [0, \frac{2c}{L^2}]$, and the fact that

$$1 - 2c\left[\frac{2c}{L^2}\right] + \left[\frac{2c}{L^2}\right]^2 L^2 = 1 - \frac{4c^2}{L^2} + \left[\frac{4c^2}{L^4}\right]L^2 = 1 - \frac{4c^2}{L^2} + \frac{4c^2}{L^2} = 1 \qquad (14.33)$$

and establish item (ii). Next we claim that for all $n \in \mathbb{N}$ it holds that

$$\|\Theta_n - \vartheta\|_2 \leq (1 - 2c\gamma_n + (\gamma_n)^2 L^2)^{1/2}\|\Theta_{n-1} - \vartheta\|_2 \leq r. \qquad (14.34)$$

We now prove (14.34) by induction on $n \in \mathbb{N}$. For the base case $n = 1$ observe that the assumption that $\Theta_0 = \xi \in \mathbb{B}$, item (ii) in Lemma 14.1.4, and item (ii) ensure that

$$\begin{aligned}
\|\Theta_1 - \vartheta\|_2^2 &= \|\Theta_0 - \gamma_1(\nabla f)(\Theta_0) - \vartheta\|_2^2 \\
&\leq (1 - 2c\gamma_1 + (\gamma_1)^2 L^2)\|\Theta_0 - \vartheta\|_2^2 \\
&\leq \|\Theta_0 - \vartheta\|_2^2 \leq r^2.
\end{aligned} \qquad (14.35)$$

This establishes (14.34) in the base case $n = 1$. For the induction step observe that item (ii) in Lemma 14.1.4 and item (ii) imply that for all $n \in \mathbb{N}$ with $\Theta_n \in \mathbb{B}$ it holds that

$$\begin{aligned}
\|\Theta_{n+1} - \vartheta\|_2^2 &= \|\Theta_n - \gamma_{n+1}(\nabla f)(\Theta_n) - \vartheta\|_2^2 \\
&\leq \underbrace{(1 - 2c\gamma_{n+1} + (\gamma_{n+1})^2 L^2)}_{\in[0,1]}\|\Theta_n - \vartheta\|_2^2 \\
&\leq \|\Theta_n - \vartheta\|_2^2 \leq r^2.
\end{aligned} \qquad (14.36)$$

This demonstrates that for all $n \in \mathbb{N}$ with $\|\Theta_n - \vartheta\|_2 \leq r$ it holds that

$$\|\Theta_{n+1} - \vartheta\|_2 \leq (1 - 2c\gamma_{n+1} + (\gamma_{n+1})^2 L^2)^{1/2}\|\Theta_n - \vartheta\|_2 \leq r. \qquad (14.37)$$

Induction thus proves (14.34). Next observe that (14.34) establishes item (iii). Moreover, note that induction and item (iii) prove item (iv). Furthermore, note that item (iii) and the fact that $\Theta_0 = \xi \in \mathbb{B}$ ensure that for all $n \in \mathbb{N}_0$ it holds that $\Theta_n \in \mathbb{B}$. Combining this, (14.29), and Lemma 13.2.9 with items (i) and (iv) establishes item (v). The proof of Proposition 14.1.5 is thus complete. $\square$

### 14.1.2.2 Size of the learning rates

In the next result, Corollary 14.1.6 below, we, roughly speaking, specialize Proposition 14.1.5 to the case where the learning rates $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \frac{2c}{L^2}]$ are a constant sequence.

**Corollary 14.1.6** (Convergence of gradient descent for constant learning rates). *Let $d \in \mathbb{N}$, $c, L \in (0, \infty)$, $r \in (0, \infty]$, $\gamma \in (0, \frac{2c}{L^2})$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{w \in \mathbb{R}^d : \|w - \vartheta\| \leq r\}$, $\xi \in \mathbb{B}$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ satisfy for all $\theta \in \mathbb{B}$ that*

$$\langle \theta - \vartheta, (\nabla f)(\theta) \rangle \geq c\|\theta - \vartheta\|_2^2 \qquad \text{and} \qquad \|(\nabla f)(\theta)\|_2 \leq L\|\theta - \vartheta\|_2, \qquad (14.38)$$

*and let $\Theta \colon \mathbb{N}_0 \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad \text{and} \qquad \Theta_n = \Theta_{n-1} - \gamma(\nabla f)(\Theta_{n-1}). \qquad (14.39)$$

*Then*

(i) *it holds that* $\{\theta \in \mathbb{B} \colon f(\theta) = \inf_{w \in \mathbb{B}} f(w)\} = \{\vartheta\}$,

(ii) *it holds that* $0 \leq 1 - 2c\gamma + \gamma^2 L^2 < 1$,

(iii) *it holds for all* $n \in \mathbb{N}_0$ *that*

$$\|\Theta_n - \vartheta\|_2 \leq \left[1 - 2c\gamma + \gamma^2 L^2\right]^{n/2} \|\xi - \vartheta\|_2, \tag{14.40}$$

*and*

(iv) *it holds for all* $n \in \mathbb{N}_0$ *that*

$$0 \leq f(\Theta_n) - f(\vartheta) \leq \tfrac{L}{2}\|\Theta_n - \vartheta\|_2^2 \leq \tfrac{L}{2}\left[1 - 2c\gamma + \gamma^2 L^2\right]^n \|\xi - \vartheta\|_2^2. \tag{14.41}$$

*Proof of Corollary 14.1.6.* Observe that item (iii) in Lemma 14.1.4 proves item (ii). In addition, note that Proposition 14.1.5 establishes items (i), (iii), and (iv). The proof of Corollary 14.1.6 is thus complete. □

Corollary 14.1.6 above establishes under suitable hypotheses convergence of the GD optimization method in the case where the learning rates are constant and strictly smaller than $\frac{2c}{L^2}$. The next result, Lemma 14.1.7 below, demonstrates that the condition that the learning rates are strictly smaller than $\frac{2c}{L^2}$ in Corollary 14.1.6 can, in general, not be relaxed.

**Lemma 14.1.7** (Sharp bounds on the learning rate for the convergence of gradient descent). *Let* $d \in \mathbb{N}$, $\alpha \in (0, \infty)$, $\gamma \in \mathbb{R}$, $\vartheta \in \mathbb{R}^d$, $\xi \in \mathbb{R}^d \backslash \{\vartheta\}$, *let* $f \colon \mathbb{R}^d \to \mathbb{R}$ *satisfy for all* $\theta \in \mathbb{R}^d$ *that*

$$f(\theta) = \tfrac{\alpha}{2}\|\theta - \vartheta\|_2^2, \tag{14.42}$$

*and let* $\Theta \colon \mathbb{N}_0 \to \mathbb{R}^d$ *satisfy for all* $n \in \mathbb{N}$ *that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \gamma(\nabla f)(\Theta_{n-1}). \tag{14.43}$$

*Then*

(i) *it holds for all* $\theta \in \mathbb{R}^d$ *that* $\langle \theta - \vartheta, (\nabla f)(\theta) \rangle = \alpha \|\theta - \vartheta\|_2^2$,

(ii) *it holds for all* $\theta \in \mathbb{R}^d$ *that* $\|(\nabla f)(\theta)\|_2 = \alpha \|\theta - \vartheta\|_2$,

(iii) *it holds for all* $n \in \mathbb{N}_0$ *that* $\|\Theta_n - \vartheta\|_2 = |1 - \gamma\alpha|^n \|\xi - \vartheta\|_2$, *and*

(iv) *it holds that*

$$\liminf_{n \to \infty}\|\Theta_n - \vartheta\|_2 = \limsup_{n \to \infty}\|\Theta_n - \vartheta\|_2 = \begin{cases} 0 & : \gamma \in (0, 2/\alpha) \\ \|\xi - \vartheta\|_2 & : \gamma \in \{0, 2/\alpha\} \\ \infty & : \gamma \in \mathbb{R} \setminus [0, 2/\alpha]. \end{cases} \tag{14.44}$$

*Proof of Lemma 14.1.7.* First of all, note that Lemma 13.2.4 ensures that for all $\theta \in \mathbb{R}^d$ it holds that $f \in C^\infty(\mathbb{R}^d, \mathbb{R})$ and

$$(\nabla f)(\theta) = \tfrac{\alpha}{2}(2(\theta - \vartheta)) = \alpha(\theta - \vartheta). \tag{14.45}$$

This proves item (ii). Moreover, observe that (14.45) assures that for all $\theta \in \mathbb{R}^d$ it holds that

$$\langle \theta - \vartheta, (\nabla f)(\theta) \rangle = \langle \theta - \vartheta, \alpha(\theta - \vartheta) \rangle = \alpha \|\theta - \vartheta\|^2. \tag{14.46}$$

This establishes item (i). Next note that (14.43) and (14.45) demonstrate that for all $n \in \mathbb{N}$ it holds that

$$\begin{aligned}
\Theta_n - \vartheta &= \Theta_{n-1} - \gamma(\nabla f)(\Theta_{n-1}) - \vartheta \\
&= \Theta_{n-1} - \gamma\alpha(\Theta_{n-1} - \vartheta) - \vartheta \\
&= (1 - \gamma\alpha)(\Theta_{n-1} - \vartheta).
\end{aligned} \tag{14.47}$$

The assumption that $\Theta_0 = \xi$ and induction hence prove that for all $n \in \mathbb{N}_0$ it holds that

$$\Theta_n - \vartheta = (1 - \gamma\alpha)^n(\Theta_0 - \vartheta) = (1 - \gamma\alpha)^n(\xi - \vartheta). \tag{14.48}$$

Therefore, we obtain for all $n \in \mathbb{N}_0$ that

$$\|\Theta_n - \vartheta\|_2 = |1 - \gamma\alpha|^n \|\xi - \vartheta\|_2. \tag{14.49}$$

This establishes item (iii). Combining item (iii) with the fact that for all $t \in (0, 2/\alpha)$ it holds that $|1 - t\alpha| \in [0, 1)$, the fact that for all $t \in \{0, 2/\alpha\}$ it holds that $|1 - t\alpha| = 1$, the fact that for all $t \in \mathbb{R} \setminus [0, 2/\alpha]$ it holds that $|1 - t\alpha| \in (1, \infty)$, and the fact that $\|\xi - \vartheta\|_2 > 0$ establishes item (iv). The proof of Lemma 14.1.7 is thus complete. $\qquad\square$

### 14.1.2.3 Convergence rates

The next result, Corollary 14.1.8 below, establishes a convergence rate for the GD optimization method in the case of possibly non-constant learning rates. We prove Corollary 14.1.8 through an application of Proposition 14.1.5 above.

**Corollary 14.1.8** (Qualitative convergence of gradient descent). *Let $d \in \mathbb{N}$, $c, L \in (0, \infty)$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \infty)$, $\xi, \vartheta \in \mathbb{R}^d$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ satisfy for all $\theta \in \mathbb{R}^d$ that*

$$\langle \theta - \vartheta, (\nabla f)(\theta) \rangle \geq c\|\theta - \vartheta\|_2^2, \qquad \|(\nabla f)(\theta)\|_2 \leq L\|\theta - \vartheta\|_2, \tag{14.50}$$

$$\text{and} \qquad 0 < \liminf_{n \to \infty} \gamma_n \leq \limsup_{n \to \infty} \gamma_n < \tfrac{2c}{L^2}, \tag{14.51}$$

*and let $\Theta \colon \mathbb{N}_0 \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad \text{and} \qquad \Theta_n = \Theta_{n-1} - \gamma_n(\nabla f)(\Theta_{n-1}). \tag{14.52}$$

*Then*

(i) *it holds that $\{\theta \in \mathbb{R}^d \colon f(\theta) = \inf_{w \in \mathbb{R}^d} f(w)\} = \{\vartheta\}$,*

(ii) *there exist $\epsilon \in (0, 1)$, $C \in \mathbb{R}$ such that for all $n \in \mathbb{N}_0$ it holds that*

$$\|\Theta_n - \vartheta\|_2 \leq \epsilon^n C, \tag{14.53}$$

*and*

(iii) *there exist $\epsilon \in (0, 1)$, $C \in \mathbb{R}$ such that for all $n \in \mathbb{N}_0$ it holds that*

$$0 \leq f(\Theta_n) - f(\vartheta) \leq \epsilon^n C. \tag{14.54}$$

*Proof of Corollary 14.1.8.* Throughout this proof let $\alpha, \beta \in \mathbb{R}$ satisfy

$$0 < \alpha < \liminf_{n \to \infty} \gamma_n \le \limsup_{n \to \infty} \gamma_n < \beta < \tfrac{2c}{L^2} \tag{14.55}$$

(cf. (14.51)), let $m \in \mathbb{N}$ satisfy for all $n \in \mathbb{N}$ that $\gamma_{m+n} \in [\alpha, \beta]$, and let $h \colon \mathbb{R} \to \mathbb{R}$ satisfy for all $t \in \mathbb{R}$ that

$$h(t) = 1 - 2ct + t^2 L^2. \tag{14.56}$$

Observe that (14.50) and item (ii) in Lemma 13.2.8 prove item (i). In addition, observe that the fact that for all $t \in \mathbb{R}$ it holds that $h'(t) = -2c + 2tL^2$ implies that for all $t \in (-\infty, \tfrac{c}{L^2}]$ it holds that

$$h'(t) \le -2c + 2\big[\tfrac{c}{L^2}\big]L^2 = 0. \tag{14.57}$$

The fundamental theorem of calculus hence assures that for all $t \in [\alpha, \beta] \cap [0, \tfrac{c}{L^2}]$ it holds that

$$h(t) = h(\alpha) + \int_\alpha^t h'(s)\,\mathrm{d}s \le h(\alpha) + \int_\alpha^t 0\,\mathrm{d}s = h(\alpha) \le \max\{h(\alpha), h(\beta)\}. \tag{14.58}$$

Furthermore, observe that the fact that for all $t \in \mathbb{R}$ it holds that $h'(t) = -2c + 2tL^2$ implies that for all $t \in [\tfrac{c}{L^2}, \infty)$ it holds that

$$h'(t) \ge -2c + 2\big[\tfrac{c}{L^2}\big]L^2 = 0. \tag{14.59}$$

The fundamental theorem of calculus hence ensures that for all $t \in [\alpha, \beta] \cap [\tfrac{c}{L^2}, \infty)$ it holds that

$$\max\{h(\alpha), h(\beta)\} \ge h(\beta) = h(t) + \int_t^\beta h'(s)\,\mathrm{d}s \ge h(t) + \int_t^\beta 0\,\mathrm{d}s = h(t). \tag{14.60}$$

Combining this and (14.58) establishes that for all $t \in [\alpha, \beta]$ it holds that

$$h(t) \le \max\{h(\alpha), h(\beta)\}. \tag{14.61}$$

Moreover, observe that the fact that $\alpha, \beta \in (0, \tfrac{2c}{L^2})$ and item (iii) in Lemma 14.1.4 ensure that

$$\{h(\alpha), h(\beta)\} \subseteq [0, 1). \tag{14.62}$$

Hence, we obtain that

$$\max\{h(\alpha), h(\beta)\} \in [0, 1). \tag{14.63}$$

This implies that there exists $\varepsilon \in \mathbb{R}$ such that

$$0 \le \max\{h(\alpha), h(\beta)\} < \varepsilon < 1. \tag{14.64}$$

Next note that the fact that for all $n \in \mathbb{N}$ it holds that $\gamma_{m+n} \in [\alpha, \beta] \subseteq [0, \tfrac{2c}{L^2}]$, items (ii) and (iv) in Proposition 14.1.5 (applied with $d \curvearrowright d$, $c \curvearrowright c$, $L \curvearrowright L$, $r \curvearrowright \infty$, $(\gamma_n)_{n \in \mathbb{N}} \curvearrowright (\gamma_{m+n})_{n \in \mathbb{N}}$, $\vartheta \curvearrowright \vartheta$, $\xi \curvearrowright \Theta_m$, $f \curvearrowright f$ in the notation of Proposition 14.1.5), (14.50), (14.52), and (14.61) demonstrate that for all $n \in \mathbb{N}$ it holds that

$$
\begin{aligned}
\|\Theta_{m+n} - \vartheta\|_2 &\le \left[\prod_{k=1}^n (1 - 2c\gamma_{m+k} + (\gamma_{m+k})^2 L^2)^{1/2}\right] \|\Theta_m - \vartheta\|_2 \\
&= \left[\prod_{k=1}^n (h(\gamma_{m+k}))^{1/2}\right] \|\Theta_m - \vartheta\|_2 \\
&\le (\max\{h(\alpha), h(\beta)\})^{n/2} \|\Theta_m - \vartheta\|_2 \\
&\le \varepsilon^{n/2} \|\Theta_m - \vartheta\|_2.
\end{aligned}
\tag{14.65}
$$

This shows that for all $n \in \mathbb{N}$ with $n > m$ it holds that

$$\|\Theta_n - \vartheta\|_2 \leq \varepsilon^{(n-m)/2} \|\Theta_m - \vartheta\|_2. \tag{14.66}$$

The fact that for all $n \in \mathbb{N}_0$ with $n \leq m$ it holds that

$$\|\Theta_n - \vartheta\|_2 = \left[\frac{\|\Theta_n - \vartheta\|_2}{\varepsilon^{n/2}}\right]\varepsilon^{n/2} \leq \left[\max\left\{\frac{\|\Theta_k - \vartheta\|_2}{\varepsilon^{k/2}} : k \in \{0, 1, \ldots, m\}\right\}\right]\varepsilon^{n/2} \tag{14.67}$$

hence assures that for all $n \in \mathbb{N}_0$ it holds that

$$\begin{aligned}
&\|\Theta_n - \vartheta\|_2 \\
&\leq \max\left\{\left[\max\left\{\frac{\|\Theta_k - \vartheta\|_2}{\varepsilon^{k/2}} : k \in \{0, 1, \ldots, m\}\right\}\right]\varepsilon^{n/2}, \varepsilon^{(n-m)/2}\|\Theta_m - \vartheta\|_2\right\} \\
&= (\varepsilon^{1/2})^n\left[\max\left\{\max\left\{\frac{\|\Theta_k - \vartheta\|_2}{\varepsilon^{k/2}} : k \in \{0, 1, \ldots, m\}\right\}, \varepsilon^{-m/2}\|\Theta_m - \vartheta\|_2\right\}\right] \\
&= (\varepsilon^{1/2})^n\left[\max\left\{\frac{\|\Theta_k - \vartheta\|_2}{\varepsilon^{k/2}} : k \in \{0, 1, \ldots, m\}\right\}\right].
\end{aligned} \tag{14.68}$$

This proves item (ii). In addition, note that Lemma 13.2.9, item (i), and (14.68) assure that for all $n \in \mathbb{N}_0$ it holds that

$$\begin{aligned}
0 &\leq f(\Theta_n) - f(\vartheta) \leq \tfrac{L}{2}\|\Theta_n - \vartheta\|_2^2 \\
&\leq \frac{\varepsilon^n L}{2}\left[\max\left\{\frac{\|\Theta_k - \vartheta\|_2^2}{\varepsilon^k} : k \in \{0, 1, \ldots, m\}\right\}\right].
\end{aligned} \tag{14.69}$$

This establishes item (iii). The proof of Corollary 14.1.8 is thus complete. $\qquad\square$

### 14.1.2.4 Error estimates in the case of small learning rates

Inequality (14.31) in item (iv) in Proposition 14.1.5 above provides us an error estimate for the GD optimization method in the case where the learning rates $(\gamma_n)_{n\in\mathbb{N}}$ in Proposition 14.1.5 satisfy that for all $n \in \mathbb{N}$ it holds that $\gamma_n \leq \frac{2c}{L^2}$. The error estimate in (14.31) can be simplified in the special case where the learning rates $(\gamma_n)_{n\in\mathbb{N}}$ satisfy the more restrictive condition that for all $n \in \mathbb{N}$ it holds that $\gamma_n \leq \frac{c}{L^2}$. This is the subject of the next result, Corollary 14.1.9 below. We prove Corollary 14.1.9 through an application of Proposition 14.1.5 above.

**Corollary 14.1.9** (Error estimates in the case of small learning rates). *Let $d \in \mathbb{N}$, $c, L \in (0, \infty)$, $r \in (0, \infty]$, $(\gamma_n)_{n\in\mathbb{N}} \subseteq [0, \frac{c}{L^2}]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{w \in \mathbb{R}^d \colon \|w - \vartheta\| \leq r\}$, $\xi \in \mathbb{B}$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ satisfy for all $\theta \in \mathbb{B}$ that*

$$\langle \theta - \vartheta, (\nabla f)(\theta)\rangle \geq c\|\theta - \vartheta\|_2^2 \qquad and \qquad \|(\nabla f)(\theta)\|_2 \leq L\|\theta - \vartheta\|_2, \tag{14.70}$$

*and let $\Theta \colon \mathbb{N}_0 \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \gamma_n(\nabla f)(\Theta_{n-1}). \tag{14.71}$$

*Then*

*(i) it holds that $\{\theta \in \mathbb{B} \colon f(\theta) = \inf_{w \in \mathbb{B}} f(w)\} = \{\vartheta\}$,*

*(ii) it holds for all $n \in \mathbb{N}$ that $0 \leq 1 - c\gamma_n \leq 1$,*

*(iii) it holds for all $n \in \mathbb{N}_0$ that*

$$\|\Theta_n - \vartheta\|_2 \leq \left[\prod_{k=1}^{n} (1 - c\gamma_k)^{1/2}\right] \|\xi - \vartheta\|_2, \tag{14.72}$$

*and*

*(iv) it holds for all $n \in \mathbb{N}_0$ that*

$$0 \leq f(\Theta_n) - f(\vartheta) \leq \frac{L}{2} \left[\prod_{k=1}^{n} (1 - c\gamma_k)\right] \|\xi - \vartheta\|_2^2. \tag{14.73}$$

*Proof of Corollary 14.1.9.* Note that item (ii) in Proposition 14.1.5 and the assumption that for all $n \in \mathbb{N}$ it holds that $\gamma_n \in [0, \frac{c}{L^2}]$ ensure that for all $n \in \mathbb{N}$ it holds that

$$0 \leq 1 - 2c\gamma_n + (\gamma_n)^2 L^2 \leq 1 - 2c\gamma_n + \gamma_n \left[\frac{c}{L^2}\right] L^2 = 1 - 2c\gamma_n + \gamma_n c = 1 - c\gamma_n \leq 1. \tag{14.74}$$

This proves item (ii). Moreover, note that (14.74) and Proposition 14.1.5 establish items (i), (iii), and (iv). The proof of Corollary 14.1.9 is thus complete. $\square$

In the next result, Corollary 14.1.10 below, we, roughly speaking, specialize Corollary 14.1.9 above to the case where the learning rates $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \frac{c}{L^2}]$ are a constant sequence.

**Corollary 14.1.10** (Error estimates in the case of small and constant learning rates).
*Let $d \in \mathbb{N}$, $c, L \in (0, \infty)$, $r \in (0, \infty]$, $\gamma \in (0, \frac{c}{L^2}]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{w \in \mathbb{R}^d \colon \|w - \vartheta\| \leq r\}$, $\xi \in \mathbb{B}$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ satisfy for all $\theta \in \mathbb{B}$ that*

$$\langle \theta - \vartheta, (\nabla f)(\theta) \rangle \geq c \|\theta - \vartheta\|_2^2 \qquad and \qquad \|(\nabla f)(\theta)\|_2 \leq L \|\theta - \vartheta\|_2, \tag{14.75}$$

*and let $\Theta \colon \mathbb{N}_0 \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \gamma(\nabla f)(\Theta_{n-1}). \tag{14.76}$$

*Then*

*(i) it holds that $\{\theta \in \mathbb{B} \colon f(\theta) = \inf_{w \in \mathbb{B}} f(w)\} = \{\vartheta\}$,*

*(ii) it holds that $0 \leq 1 - c\gamma < 1$,*

*(iii) it holds for all $n \in \mathbb{N}_0$ that $\|\Theta_n - \vartheta\|_2 \leq (1 - c\gamma)^{n/2} \|\xi - \vartheta\|_2$, and*

*(iv) it holds for all $n \in \mathbb{N}_0$ that $0 \leq f(\Theta_n) - f(\vartheta) \leq \frac{L}{2} (1 - c\gamma)^n \|\xi - \vartheta\|_2^2$.*

*Proof of Corollary 14.1.10.* Note that Corollary 14.1.10 is an immediate consequence of Corollary 14.1.9. The proof of Corollary 14.1.10 is thus complete. $\square$

### 14.1.2.5   On the spectrum of the Hessian of the objective function at a local minimum

A crucial ingredient in our error analysis for the GD optimization method in Subsections 14.1.2.1–14.1.2.4 above is to employ the growth and the coercivity-type hypotheses, e.g., in (14.29) in Proposition 14.1.5 above. In this subsection we disclose in Lemma 14.1.12 below suitable conditions on the Hessians of the objective function of the considered optimization problem which are sufficient to ensure that (14.29) is satisfied so that we are in the position to apply the error analysis in Subsections 14.1.2.1–14.1.2.4 above (cf. Corollary 14.1.13 below). Our proof of Lemma 14.1.12 employs the following classical result (see Lemma 14.1.11 below) for symmetric matrices with real entries.

**Lemma 14.1.11** (Properties of the spectrum of real symmetric matrices)**.** *Let* $d \in \mathbb{N}$, *let* $A \in \mathbb{R}^{d \times d}$ *be a symmetric matrix, and let* $\mathscr{S} = \{\lambda \in \mathbb{C} \colon (\exists\, v \in \mathbb{C}^d \backslash \{0\} \colon Av = \lambda v)\}$. *Then*

*(i) it holds that* $\mathscr{S} = \{\lambda \in \mathbb{R} \colon (\exists\, v \in \mathbb{R}^d \backslash \{0\} \colon Av = \lambda v)\} \subseteq \mathbb{R}$,

*(ii) it holds that*

$$\sup_{v \in \mathbb{R}^d \backslash \{0\}} \left[ \frac{\|Av\|_2}{\|v\|_2} \right] = \max_{\lambda \in \mathscr{S}} |\lambda|, \tag{14.77}$$

*and*

*(iii) it holds for all* $v \in \mathbb{R}^d$ *that*

$$\min(\mathscr{S}) \|v\|_2^2 \leq \langle v, Av \rangle \leq \max(\mathscr{S}) \|v\|_2^2. \tag{14.78}$$

*Proof of Lemma 14.1.11.* Throughout this proof let $e_1, e_2, \ldots, e_d \in \mathbb{R}^d$ be the vectors given by

$$e_1 = (1, 0, \ldots, 0), \qquad e_2 = (0, 1, 0, \ldots, 0), \qquad \ldots, \qquad e_d = (0, \ldots, 0, 1). \tag{14.79}$$

Observe that the spectral theorem for symmetric matrices (see, e.g., Petersen [24, Theorem 4.3.4]) proves that there exist $(d \times d)$-matrices $\Lambda = (\Lambda_{i,j})_{i,j \in \{1,2,\ldots,d\}}$, $O = (O_{i,j})_{i,j \in \{1,2,\ldots,d\}} \in \mathbb{R}^{d \times d}$ such that $\mathscr{S} = \{\Lambda_{1,1}, \Lambda_{2,2}, \ldots, \Lambda_{d,d}\}$, $O^*O = OO^* = \mathrm{I}_d$, $A = O\Lambda O^*$, and

$$\Lambda = (\Lambda_{i,j})_{i,j \in \{1,2,\ldots,d\}} = \begin{pmatrix} \Lambda_{1,1} & & 0 \\ & \ddots & \\ 0 & & \Lambda_{d,d} \end{pmatrix} \in \mathbb{R}^{d \times d} \tag{14.80}$$

(cf. Definition 2.2.9). Hence, we obtain that $\mathscr{S} \subseteq \mathbb{R}$. Next note that the assumption that $\mathscr{S} = \{\lambda \in \mathbb{C} \colon (\exists\, v \in \mathbb{C}^d \backslash \{0\} \colon Av = \lambda v)\}$ ensures that for every $\lambda \in \mathscr{S}$ there exists $v \in \mathbb{C}^d \backslash \{0\}$ such that

$$A\mathfrak{Re}(v) + \mathbf{i}A\mathfrak{Im}(v) = Av = \lambda v = \lambda\mathfrak{Re}(v) + \mathbf{i}\lambda\mathfrak{Im}(v). \tag{14.81}$$

The fact that $\mathscr{S} \subseteq \mathbb{R}$ therefore demonstrates that for every $\lambda \in \mathscr{S}$ there exists $v \in \mathbb{R}^d \backslash \{0\}$ such that $Av = \lambda v$. This and the fact that $\mathscr{S} \subseteq \mathbb{R}$ ensure that $\mathscr{S} \subseteq \{\lambda \in \mathbb{R} \colon (\exists\, v \in \mathbb{R}^d \backslash \{0\} \colon Av = \lambda v)\}$. Combining this and the fact that $\{\lambda \in \mathbb{R} \colon (\exists\, v \in \mathbb{R}^d \backslash \{0\} \colon Av =$

$\lambda v)\} \subseteq \mathscr{S}$ proves item (i). Furthermore, note that (14.80) assures that for all $v = (v_1, v_2, \ldots, v_d) \in \mathbb{R}^d$ it holds that

$$
\begin{aligned}
\|\Lambda v\| &= \left[ \sum_{i=1}^{d} |\Lambda_{i,i} v_i|^2 \right]^{1/2} \leq \left[ \sum_{i=1}^{d} \max\{|\Lambda_{1,1}|^2, \ldots, |\Lambda_{d,d}|^2\} |v_i|^2 \right]^{1/2} \\
&= \left[ \max\{|\Lambda_{1,1}|, \ldots, |\Lambda_{d,d}|\}^2 \|v\|_2^2 \right]^{1/2} \\
&= \max\{|\Lambda_{1,1}|, \ldots, |\Lambda_{d,d}|\} \|v\|_2 \\
&= \left( \max_{\lambda \in \mathscr{S}} |\lambda| \right) \|v\|_2.
\end{aligned} \tag{14.82}
$$

The fact that $O$ is an orthogonal matrix and the fact that $A = O\Lambda O^*$ therefore imply that for all $v \in \mathbb{R}^d$ it holds that

$$
\begin{aligned}
\|Av\|_2 &= \|O\Lambda O^* v\|_2 = \|\Lambda O^* v\|_2 \\
&\leq \left( \max_{\lambda \in \mathscr{S}} |\lambda| \right) \|O^* v\|_2 \\
&= \left( \max_{\lambda \in \mathscr{S}} |\lambda| \right) \|v\|_2.
\end{aligned} \tag{14.83}
$$

This implies that

$$
\sup_{v \in \mathbb{R}^d \setminus \{0\}} \left[ \frac{\|Av\|_2}{\|v\|_2} \right] \leq \sup_{v \in \mathbb{R}^d \setminus \{0\}} \left[ \frac{\left( \max_{\lambda \in \mathscr{S}} |\lambda| \right) \|v\|_2}{\|v\|_2} \right] = \max_{\lambda \in \mathscr{S}} |\lambda|. \tag{14.84}
$$

In addition, note that the fact that $\mathscr{S} = \{\Lambda_{1,1}, \Lambda_{2,2} \ldots, \Lambda_{d,d}\}$ ensures that there exists $j \in \{1, 2, \ldots, d\}$ such that

$$
|\Lambda_{j,j}| = \max_{\lambda \in \mathscr{S}} |\lambda|. \tag{14.85}
$$

Next observe that the fact that $A = O\Lambda O^*$, the fact that $O$ is an orthogonal matrix, and (14.85) imply that

$$
\begin{aligned}
\sup_{v \in \mathbb{R}^d \setminus \{0\}} \left[ \frac{\|Av\|_2}{\|v\|_2} \right] &\geq \frac{\|AOe_j\|_2}{\|Oe_j\|_2} = \|O\Lambda O^* Oe_j\|_2 = \|O\Lambda e_j\|_2 \\
&= \|\Lambda e_j\|_2 = \|\Lambda_{j,j} e_j\|_2 = |\Lambda_{j,j}| = \max_{\lambda \in \mathscr{S}} |\lambda|.
\end{aligned} \tag{14.86}
$$

Combining this and (14.84) establishes item (ii). It thus remains to prove item (iii). For this note that (14.80) ensures that for all $v = (v_1, v_2, \ldots, v_d) \in \mathbb{R}^d$ it holds that

$$
\begin{aligned}
\langle v, \Lambda v \rangle &= \sum_{i=1}^{d} \Lambda_{i,i} |v_i|^2 \leq \sum_{i=1}^{d} \max\{\Lambda_{1,1}, \ldots, \Lambda_{d,d}\} |v_i|^2 \\
&= \max\{\Lambda_{1,1}, \ldots, \Lambda_{d,d}\} \|v\|_2^2 = \max(\mathscr{S}) \|v\|_2^2.
\end{aligned} \tag{14.87}
$$

The fact that $O$ is an orthogonal matrix and the fact that $A = O\Lambda O^*$ therefore demonstrate that for all $v \in \mathbb{R}^d$ it holds that

$$
\begin{aligned}
\langle v, Av \rangle &= \langle v, O\Lambda O^* v \rangle = \langle O^* v, \Lambda O^* v \rangle \\
&\leq \max(\mathscr{S}) \|O^* v\|_2^2 = \max(\mathscr{S}) \|v\|_2^2.
\end{aligned} \tag{14.88}
$$

Moreover, observe that (14.80) implies that for all $v = (v_1, v_2, \ldots, v_d) \in \mathbb{R}^d$ it holds that

$$
\begin{aligned}
\langle v, \Lambda v \rangle &= \sum_{i=1}^{d} \Lambda_{i,i} |v_i|^2 \geq \sum_{i=1}^{d} \min\{\Lambda_{1,1}, \ldots, \Lambda_{d,d}\} |v_i|^2 \\
&= \min\{\Lambda_{1,1}, \ldots, \Lambda_{d,d}\} \|v\|_2^2 = \min(\mathscr{S}) \|v\|_2^2.
\end{aligned} \tag{14.89}
$$

The fact that $O$ is an orthogonal matrix and the fact that $A = O\Lambda O^*$ hence demonstrate that for all $v \in \mathbb{R}^d$ it holds that

$$
\begin{aligned}
\langle v, Av \rangle = \langle v, O\Lambda O^* v \rangle &= \langle O^* v, \Lambda O^* v \rangle \\
&\geq \min(\mathscr{S}) \| O^* v \|_2^2 = \min(\mathscr{S}) \| v \|_2^2.
\end{aligned}
\tag{14.90}
$$

Combining this with (14.88) establishes item (iii). The proof of Lemma 14.1.11 is thus complete. $\qquad\square$

We now present the promised Lemma 14.1.12 which discloses suitable conditions (cf. (14.91)–(14.92) below) on the Hessians of the objective function of the considered optimization problem which are sufficient to ensure that (14.29) is satisfied so that we are in the position to apply the error analysis in Subsections 14.1.2.1–14.1.2.4 above.

**Lemma 14.1.12** (Conditions on the spectrum of the Hessian of the objective function at a local minimum). *Let $d \in \mathbb{N}$, let $\|\!\|\cdot\|\!\| : \mathbb{R}^{d \times d} \to [0, \infty)$ satisfy for all $A \in \mathbb{R}^{d \times d}$ that $\|\!\| A \|\!\| = \sup_{v \in \mathbb{R}^d \backslash \{0\}} \frac{\|Av\|_2}{\|v\|_2}$, and let $\lambda, \alpha \in (0, \infty)$, $\beta \in [\alpha, \infty)$, $\vartheta \in \mathbb{R}^d$, $f \in C^2(\mathbb{R}^d, \mathbb{R})$ satisfy for all $v, w \in \mathbb{R}^d$ that*

$$
(\nabla f)(\vartheta) = 0, \qquad \|\!\| (\operatorname{Hess} f)(v) - (\operatorname{Hess} f)(w) \|\!\| \leq \lambda \| v - w \|_2,
\tag{14.91}
$$

$$
and \qquad \{ \mu \in \mathbb{R} : (\exists\, u \in \mathbb{R}^d \backslash \{0\} : [(\operatorname{Hess} f)(\vartheta)]u = \mu u) \} \subseteq [\alpha, \beta].
\tag{14.92}
$$

*Then it holds for all $\theta \in \{ w \in \mathbb{R}^d : \| w - \vartheta \|_2 \leq \frac{\alpha}{\lambda} \}$ that*

$$
\langle \theta - \vartheta, (\nabla f)(\theta) \rangle \geq \tfrac{\alpha}{2} \| \theta - \vartheta \|_2^2 \qquad and \qquad \| (\nabla f)(\theta) \|_2 \leq \tfrac{3\beta}{2} \| \theta - \vartheta \|_2.
\tag{14.93}
$$

*Proof of Lemma 14.1.12.* Throughout this proof let $\mathbb{B} \subseteq \mathbb{R}^d$ be the set given by

$$
\mathbb{B} = \left\{ w \in \mathbb{R}^d : \| w - \vartheta \|_2 \leq \tfrac{\alpha}{\lambda} \right\}
\tag{14.94}
$$

and let $\mathscr{S} \subseteq \mathbb{C}$ be the set given by

$$
\mathscr{S} = \{ \mu \in \mathbb{C} : (\exists\, u \in \mathbb{C}^d \backslash \{0\} : [(\operatorname{Hess} f)(\vartheta)]u = \mu u) \}.
\tag{14.95}
$$

Observe that the fact that $(\operatorname{Hess} f)(\vartheta) \in \mathbb{R}^{d \times d}$ is a symmetric matrix, item (i) in Lemma 14.1.11, and (14.92) imply that

$$
\mathscr{S} = \{ \mu \in \mathbb{R} : (\exists\, u \in \mathbb{R}^d \backslash \{0\} : [(\operatorname{Hess} f)(\vartheta)]u = \mu u) \} \subseteq [\alpha, \beta].
\tag{14.96}
$$

Next note that the assumption that $(\nabla f)(\vartheta) = 0$ and the fundamental theorem of calculus ensure that for all $\theta, w \in \mathbb{R}^d$ it holds that

$$
\begin{aligned}
\langle w, (\nabla f)(\theta) \rangle &= \langle w, (\nabla f)(\theta) - (\nabla f)(\vartheta) \rangle \\
&= \Big\langle w, [(\nabla f)(\vartheta + t(\theta - \vartheta))]_{t=0}^{t=1} \Big\rangle \\
&= \Big\langle w, \textstyle\int_0^1 [(\operatorname{Hess} f)(\vartheta + t(\theta - \vartheta))](\theta - \vartheta) \, dt \Big\rangle \\
&= \int_0^1 \big\langle w, [(\operatorname{Hess} f)(\vartheta + t(\theta - \vartheta))](\theta - \vartheta) \big\rangle dt \\
&= \big\langle w, [(\operatorname{Hess} f)(\vartheta)](\theta - \vartheta) \big\rangle \\
&\quad + \int_0^1 \big\langle w, \big[ (\operatorname{Hess} f)(\vartheta + t(\theta - \vartheta)) - (\operatorname{Hess} f)(\vartheta) \big](\theta - \vartheta) \big\rangle dt.
\end{aligned}
\tag{14.97}
$$

The fact that $(\operatorname{Hess} f)(\vartheta) \in \mathbb{R}^{d \times d}$ is a symmetric matrix, item (iii) in Lemma 14.1.11, the Cauchy-Schwarz inequality, (14.96), (14.91), and (14.92) therefore imply that for all $\theta \in \mathbb{B}$ it holds that

$$
\begin{aligned}
&\langle \theta - \vartheta, (\nabla f)(\theta) \rangle \\
&\geq \langle \theta - \vartheta, [(\operatorname{Hess} f)(\vartheta)](\theta - \vartheta) \rangle \\
&\quad - \left| \int_0^1 \langle \theta - \vartheta, [(\operatorname{Hess} f)(\vartheta + t(\theta - \vartheta)) - (\operatorname{Hess} f)(\vartheta)](\theta - \vartheta) \rangle \, \mathrm{d}t \right| \\
&\geq \min(\mathscr{S}) \|\theta - \vartheta\|_2^2 \\
&\quad - \int_0^1 \|\theta - \vartheta\|_2 \left\| [(\operatorname{Hess} f)(\vartheta + t(\theta - \vartheta)) - (\operatorname{Hess} f)(\vartheta)](\theta - \vartheta) \right\|_2 \mathrm{d}t \\
&\geq \alpha \|\theta - \vartheta\|_2^2 \\
&\quad - \int_0^1 \|\theta - \vartheta\|_2 \|(\operatorname{Hess} f)(\vartheta + t(\theta - \vartheta)) - (\operatorname{Hess} f)(\vartheta)\| \|\theta - \vartheta\|_2 \, \mathrm{d}t \\
&\geq \alpha \|\theta - \vartheta\|_2^2 - \left[ \int_0^1 \lambda \|\vartheta + t(\theta - \vartheta) - \vartheta\|_2 \, \mathrm{d}t \right] \|\theta - \vartheta\|_2^2 \\
&= \left( \alpha - \left[ \int_0^1 t \, \mathrm{d}t \right] \lambda \|\theta - \vartheta\|_2 \right) \|\theta - \vartheta\|_2^2 = \left( \alpha - \tfrac{\lambda}{2} \|\theta - \vartheta\|_2 \right) \|\theta - \vartheta\|_2^2 \\
&\geq \left( \alpha - \tfrac{\lambda \alpha}{2\lambda} \right) \|\theta - \vartheta\|_2^2 = \tfrac{\alpha}{2} \|\theta - \vartheta\|_2^2.
\end{aligned}
\tag{14.98}
$$

Moreover, observe that (14.91), (14.96), (14.97), the fact that $(\operatorname{Hess} f)(\vartheta) \in \mathbb{R}^{d \times d}$ is a symmetric matrix, item (ii) in Lemma 14.1.11, the Cauchy-Schwarz inequality, and the assumption that $\alpha \leq \beta$ ensure that for all $\theta \in \mathbb{B}$, $w \in \mathbb{R}^d$ with $\|w\|_2 = 1$ it holds that

$$
\begin{aligned}
&\langle w, (\nabla f)(\theta) \rangle \\
&\leq \left| \langle w, [(\operatorname{Hess} f)(\vartheta)](\theta - \vartheta) \rangle \right| \\
&\quad + \left| \int_0^1 \langle w, [(\operatorname{Hess} f)(\vartheta + t(\theta - \vartheta)) - (\operatorname{Hess} f)(\vartheta)](\theta - \vartheta) \rangle \, \mathrm{d}t \right| \\
&\leq \|w\|_2 \|[(\operatorname{Hess} f)(\vartheta)](\theta - \vartheta)\|_2 \\
&\quad + \int_0^1 \|w\|_2 \|[(\operatorname{Hess} f)(\vartheta + t(\theta - \vartheta)) - (\operatorname{Hess} f)(\vartheta)](\theta - \vartheta)\|_2 \mathrm{d}t \\
&\leq \left[ \sup_{v \in \mathbb{R}^d \setminus \{0\}} \frac{\|[(\operatorname{Hess} f)(\vartheta)]v\|_2}{\|v\|_2} \right] \|\theta - \vartheta\|_2 \\
&\quad + \int_0^1 \|(\operatorname{Hess} f)(\vartheta + t(\theta - \vartheta)) - (\operatorname{Hess} f)(\vartheta)\| \|\theta - \vartheta\|_2 \, \mathrm{d}t \\
&\leq \max(\mathscr{S}) \|\theta - \vartheta\|_2 + \left[ \int_0^1 \lambda \|\vartheta + t(\theta - \vartheta) - \vartheta\|_2 \, \mathrm{d}t \right] \|\theta - \vartheta\|_2 \\
&\leq \left( \beta + \lambda \left[ \int_0^1 t \, \mathrm{d}t \right] \|\theta - \vartheta\|_2 \right) \|\theta - \vartheta\|_2 = \left( \beta + \tfrac{\lambda}{2} \|\theta - \vartheta\|_2 \right) \|\theta - \vartheta\|_2 \\
&\leq \left( \beta + \tfrac{\lambda \alpha}{2\lambda} \right) \|\theta - \vartheta\|_2 = \left[ \tfrac{2\beta + \alpha}{2} \right] \|\theta - \vartheta\|_2 \leq \tfrac{3\beta}{2} \|\theta - \vartheta\|_2.
\end{aligned}
\tag{14.99}
$$

Therefore, we obtain for all $\theta \in \mathbb{B}$ that

$$
\|(\nabla f)(\theta)\|_2 = \sup_{w \in \mathbb{R}^d, \|w\|_2 = 1} [\langle w, (\nabla f)(\theta) \rangle] \leq \tfrac{3\beta}{2} \|\theta - \vartheta\|_2.
\tag{14.100}
$$

Combining this and (14.98) establishes (14.93). The proof of Lemma 14.1.12 is thus complete. $\qquad\square$

The next result, Corollary 14.1.13 below, combines Lemma 14.1.12 with Proposition 14.1.5 to obtain an error analysis which assumes the conditions in (14.91)–(14.92) in Lemma 14.1.12 above. A result similar to Corollary 14.1.13 can, e.g., be found in Nesterov [23, Theorem 1.2.4].

**Corollary 14.1.13** (Error analysis for the GD optimization method under conditions on the Hessian of the objective function). *Let $d \in \mathbb{N}$, let $\|\|\cdot\|\|\colon \mathbb{R}^{d \times d} \to [0, \infty)$ satisfy for all $A \in \mathbb{R}^{d \times d}$ that $\|\|A\|\| = \sup_{v \in \mathbb{R}^d \backslash \{0\}} \frac{\|Av\|_2}{\|v\|_2}$, and let $\lambda, \alpha \in (0, \infty)$, $\beta \in [\alpha, \infty)$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \frac{4\alpha}{9\beta^2}]$, $\vartheta, \xi \in \mathbb{R}^d$, $f \in C^2(\mathbb{R}^d, \mathbb{R})$ satisfy for all $v, w \in \mathbb{R}^d$ that*

$$(\nabla f)(\vartheta) = 0, \qquad \|\|(\text{Hess } f)(v) - (\text{Hess } f)(w)\|\| \leq \lambda \|v - w\|_2, \tag{14.101}$$

$$\{\mu \in \mathbb{R} \colon (\exists\, u \in \mathbb{R}^d \backslash \{0\} \colon [(\text{Hess } f)(\vartheta)]u = \mu u)\} \subseteq [\alpha, \beta], \tag{14.102}$$

*and $\|\xi - \vartheta\|_2 \leq \frac{\alpha}{\lambda}$, and let $\Theta \colon \mathbb{N}_0 \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \gamma_n (\nabla f)(\Theta_{n-1}). \tag{14.103}$$

*Then it holds for all $n \in \mathbb{N}$ that*

$$\|\Theta_n - \vartheta\|_2 \leq \left[ \prod_{k=1}^n \left[ 1 - \alpha\gamma_k + \frac{9\beta^2(\gamma_k)^2}{4} \right]^{1/2} \right] \|\xi - \vartheta\|_2 \qquad and \tag{14.104}$$

$$0 \leq f(\Theta_n) - f(\vartheta) \leq \frac{3\beta}{4} \left[ \prod_{k=1}^n \left[ 1 - \alpha\gamma_k + \frac{9\beta^2(\gamma_k)^2}{4} \right] \right] \|\xi - \vartheta\|_2^2. \tag{14.105}$$

*Proof of Corollary 14.1.13.* Throughout this proof let $\langle \cdot, \cdot \rangle \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ be the $d$-dimensional Euclidean scalar product. Note that (14.101), (14.102), and Lemma 14.1.12 prove that for all $\theta \in \{w \in \mathbb{R}^d \colon \|w - \vartheta\|_2 \leq \frac{\alpha}{\lambda}\}$ it holds that

$$\langle \theta - \vartheta, (\nabla f)(\theta) \rangle \geq \frac{\alpha}{2} \|\theta - \vartheta\|_2^2 \qquad \text{and} \qquad \|(\nabla f)(\theta)\|_2 \leq \frac{3\beta}{2} \|\theta - \vartheta\|_2. \tag{14.106}$$

Combining this, the assumption that $\|\xi - \vartheta\|_2 \leq \frac{\alpha}{\lambda}$, (14.103), and items (iv)–(v) in Proposition 14.1.5 (applied with $c \curvearrowleft \frac{\alpha}{2}$, $L \curvearrowleft \frac{3\beta}{2}$, $r \curvearrowleft \frac{\alpha}{\lambda}$ in the notation of Proposition 14.1.5) establishes (14.104) and (14.105). The proof of Corollary 14.1.13 is thus complete. $\qquad\square$

### 14.1.2.6 Equivalent conditions on the objective function

**Lemma 14.1.14.** *Let $d \in \mathbb{N}$, let $\langle\!\langle \cdot, \cdot \rangle\!\rangle \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ be a scalar product, let $\|\|\cdot\|\|\colon \mathbb{R}^d \to [0, \infty)$ satisfy for all $v \in \mathbb{R}^d$ that $\|\|v\|\| = \sqrt{\langle\!\langle v, v \rangle\!\rangle}$, let $\gamma \in (0, \infty)$, $\varepsilon \in (0, 1)$, $r \in (0, \infty]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{w \in \mathbb{R}^d \colon \|\|w - \vartheta\|\| \leq r\}$, and let $g \colon \mathbb{R}^d \to \mathbb{R}^d$ satisfy for all $\theta \in \mathbb{B}$ that*

$$\|\|\theta - \gamma g(\theta) - \vartheta\|\| \leq \varepsilon \|\|\theta - \vartheta\|\|. \tag{14.107}$$

*Then it holds for all $\theta \in \mathbb{B}$ that*

$$\langle\!\langle \theta - \vartheta, g(\theta) \rangle\!\rangle \geq \max\left\{ \left[ \frac{1 - \varepsilon^2}{2\gamma} \right] \|\|\theta - \vartheta\|\|^2, \frac{\gamma}{2} \|\|g(\theta)\|\|^2 \right\}$$
$$\geq \min\left\{ \frac{1 - \varepsilon^2}{2\gamma}, \frac{\gamma}{2} \right\} \max\left\{ \|\|\theta - \vartheta\|\|^2, \|\|g(\theta)\|\|^2 \right\}. \tag{14.108}$$

*Proof of Lemma 14.1.14.* First, note that (14.107) ensures that for all $\theta \in \mathbb{B}$ it holds that

$$
\begin{aligned}
\varepsilon^2 |\!|\!| \theta - \vartheta |\!|\!|^2 &\geq |\!|\!| \theta - \gamma g(\theta) - \vartheta |\!|\!|^2 = |\!|\!| (\theta - \vartheta) - \gamma g(\theta) |\!|\!|^2 \\
&= |\!|\!| \theta - \vartheta |\!|\!|^2 - 2\gamma \langle\!\langle \theta - \vartheta, g(\theta) \rangle\!\rangle + \gamma^2 |\!|\!| g(\theta) |\!|\!|^2.
\end{aligned}
\tag{14.109}
$$

Hence, we obtain for all $\theta \in \mathbb{B}$ that

$$
\begin{aligned}
2\gamma \langle\!\langle \theta - \vartheta, g(\theta) \rangle\!\rangle &\geq (1 - \varepsilon^2) |\!|\!| \theta - \vartheta |\!|\!|^2 + \gamma^2 |\!|\!| g(\theta) |\!|\!|^2 \\
&\geq \max\{ (1 - \varepsilon^2) |\!|\!| \theta - \vartheta |\!|\!|^2, \gamma^2 \, |\!|\!| g(\theta) |\!|\!|^2 \} \geq 0.
\end{aligned}
\tag{14.110}
$$

This demonstrates that for all $\theta \in \mathbb{B}$ it holds that

$$
\begin{aligned}
\langle\!\langle \theta - \vartheta, g(\theta) \rangle\!\rangle &\geq \tfrac{1}{2\gamma} \max\{ (1 - \varepsilon^2) |\!|\!| \theta - \vartheta |\!|\!|^2, \gamma^2 \, |\!|\!| g(\theta) |\!|\!|^2 \} \\
&= \max\left\{ \left[ \tfrac{1-\varepsilon^2}{2\gamma} \right] |\!|\!| \theta - \vartheta |\!|\!|^2, \tfrac{\gamma}{2} |\!|\!| g(\theta) |\!|\!|^2 \right\} \\
&\geq \min\left\{ \tfrac{1-\varepsilon^2}{2\gamma}, \tfrac{\gamma}{2} \right\} \max\{ |\!|\!| \theta - \vartheta |\!|\!|^2, |\!|\!| g(\theta) |\!|\!|^2 \}.
\end{aligned}
\tag{14.111}
$$

The proof of Lemma 14.1.14 is thus complete. $\qquad\square$

**Lemma 14.1.15.** *Let $d \in \mathbb{N}$, let $\langle\!\langle \cdot, \cdot \rangle\!\rangle \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ be a scalar product, let $|\!|\!| \cdot |\!|\!| \colon \mathbb{R}^d \to [0, \infty)$ satisfy for all $v \in \mathbb{R}^d$ that $|\!|\!| v |\!|\!| = \sqrt{\langle\!\langle v, v \rangle\!\rangle}$, let $c \in (0, \infty)$, $r \in (0, \infty]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{ w \in \mathbb{R}^d \colon |\!|\!| w - \vartheta |\!|\!| \leq r \}$, and let $g \colon \mathbb{R}^d \to \mathbb{R}^d$ satisfy for all $\theta \in \mathbb{B}$ that*

$$
\langle\!\langle \theta - \vartheta, g(\theta) \rangle\!\rangle \geq c \max\{ |\!|\!| \theta - \vartheta |\!|\!|^2, |\!|\!| g(\theta) |\!|\!|^2 \}.
\tag{14.112}
$$

*Then it holds for all $\theta \in \mathbb{B}$ that*

$$
\langle\!\langle \theta - \vartheta, g(\theta) \rangle\!\rangle \geq c |\!|\!| \theta - \vartheta |\!|\!|^2 \qquad and \qquad |\!|\!| g(\theta) |\!|\!| \leq \tfrac{1}{c} |\!|\!| \theta - \vartheta |\!|\!|.
\tag{14.113}
$$

*Proof of Lemma 14.1.15.* Observe that (14.112) and the Cauchy-Schwarz inequality assure that for all $\theta \in \mathbb{B}$ it holds that

$$
|\!|\!| g(\theta) |\!|\!|^2 \leq \max\{ |\!|\!| \theta - \vartheta |\!|\!|^2, |\!|\!| g(\theta) |\!|\!|^2 \} \leq \tfrac{1}{c} \langle\!\langle \theta - \vartheta, g(\theta) \rangle\!\rangle \leq \tfrac{1}{c} |\!|\!| \theta - \vartheta |\!|\!| \, |\!|\!| g(\theta) |\!|\!|.
\tag{14.114}
$$

Therefore, we obtain for all $\theta \in \mathbb{B}$ that

$$
|\!|\!| g(\theta) |\!|\!| \leq \tfrac{1}{c} |\!|\!| \theta - \vartheta |\!|\!|.
\tag{14.115}
$$

Combining this with (14.112) completes the proof of Lemma 14.1.15. $\qquad\square$

**Lemma 14.1.16.** *Let $d \in \mathbb{N}$, $c \in (0, \infty)$, $r \in (0, \infty]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{ w \in \mathbb{R}^d \colon \| w - \vartheta \|_2 \leq r \}$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ satisfy for all $\theta \in \mathbb{B}$ that*

$$
\langle \theta - \vartheta, (\nabla f)(\theta) \rangle \geq c \| \theta - \vartheta \|_2^2.
\tag{14.116}
$$

*Then it holds for all $v \in \mathbb{R}^d$, $s, t \in [0, 1]$ with $\| v \|_2 \leq r$ and $s \leq t$ that*

$$
f(\vartheta + tv) - f(\vartheta + sv) \geq \tfrac{c}{2}(t^2 - s^2) \| v \|_2^2.
\tag{14.117}
$$

*Proof of Lemma 14.1.16.* First of all, observe that (14.116) implies that for all $v \in \mathbb{R}^d$ with $\|v\|_2 \leq r$ it holds that

$$\langle (\nabla f)(\vartheta + v), v \rangle \geq c\|v\|_2^2. \tag{14.118}$$

The fundamental theorem of calculus hence ensures that for all $v \in \mathbb{R}^d$, $s, t \in [0, 1]$ with $\|v\|_2 \leq r$ and $s \leq t$ it holds that

$$\begin{aligned}
f(\vartheta + tv) - f(\vartheta + sv) &= \left[ f(\vartheta + hv) \right]_{h=s}^{h=t} \\
&= \int_s^t f'(\vartheta + hv)v \, \mathrm{d}h \\
&= \int_s^t \tfrac{1}{h} \langle (\nabla f)(\vartheta + hv), hv \rangle \, \mathrm{d}h \\
&\geq \int_s^t \tfrac{c}{h} \|hv\|_2^2 \, \mathrm{d}h \\
&= c \left[ \int_s^t h \, \mathrm{d}h \right] \|v\|_2^2 = \tfrac{c}{2}(t^2 - s^2)\|v\|_2^2.
\end{aligned} \tag{14.119}$$

The proof of Lemma 14.1.16 is thus complete. $\qquad\square$

**Lemma 14.1.17.** *Let $d \in \mathbb{N}$, $c \in (0, \infty)$, $r \in (0, \infty]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{w \in \mathbb{R}^d \colon \|w - \vartheta\|_2 \leq r\}$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ satisfy for all $v \in \mathbb{R}^d$, $s, t \in [0, 1]$ with $\|v\|_2 \leq r$ and $s \leq t$ that*

$$f(\vartheta + tv) - f(\vartheta + sv) \geq c(t^2 - s^2)\|v\|_2^2. \tag{14.120}$$

*Then it holds for all $\theta \in \mathbb{B}$ that*

$$\langle \theta - \vartheta, (\nabla f)(\theta) \rangle \geq 2c\|\theta - \vartheta\|_2^2. \tag{14.121}$$

*Proof of Lemma 14.1.17.* Observe that (14.120) ensures that for all $s \in (0, r] \cap \mathbb{R}, \theta \in \mathbb{R}^d$ with $\|\theta - \vartheta\|_2 < s$ it holds that

$$\begin{aligned}
\langle \theta - \vartheta, (\nabla f)(\theta) \rangle &= f'(\theta)(\theta - \vartheta) = \lim_{h \searrow 0} \left( \tfrac{1}{h} \left[ f(\theta + h(\theta - \vartheta)) - f(\theta) \right] \right) \\
&= \lim_{h \searrow 0} \left( \frac{1}{h} \left[ f\left( \vartheta + \tfrac{(1+h)\|\theta - \vartheta\|_2}{s} \left( \tfrac{s}{\|\theta - \vartheta\|_2}(\theta - \vartheta) \right) \right) \right. \right. \\
&\qquad \left. \left. - f\left( \vartheta + \tfrac{\|\theta - \vartheta\|_2}{s} \left( \tfrac{s}{\|\theta - \vartheta\|_2}(\theta - \vartheta) \right) \right) \right] \right) \\
&\geq \limsup_{h \searrow 0} \left( \tfrac{c}{h} \left( \left[ \tfrac{(1+h)\|\theta - \vartheta\|_2}{s} \right]^2 - \left[ \tfrac{\|\theta - \vartheta\|_2}{s} \right]^2 \right) \left\| \tfrac{s}{\|\theta - \vartheta\|_2}(\theta - \vartheta) \right\|_2^2 \right) \\
&= c \left[ \limsup_{h \searrow 0} \left( \tfrac{(1+h)^2 - 1}{h} \right) \right] \left[ \tfrac{\|\theta - \vartheta\|_2}{s} \right]^2 \left\| \tfrac{s}{\|\theta - \vartheta\|_2}(\theta - \vartheta) \right\|_2^2 \\
&= c \left[ \limsup_{h \searrow 0} \left( \tfrac{2h + h^2}{h} \right) \right] \|\theta - \vartheta\|_2^2 \\
&= c \left[ \limsup_{h \searrow 0} (2 + h) \right] \|\theta - \vartheta\|_2^2 = 2c\|\theta - \vartheta\|_2^2.
\end{aligned} \tag{14.122}$$

Hence, we obtain that for all $\theta \in \mathbb{R}^d$ with $\|\theta - \vartheta\|_2 < r$ it holds that

$$\langle \theta - \vartheta, (\nabla f)(\theta) \rangle \geq 2c\|\theta - \vartheta\|_2^2. \tag{14.123}$$

Combining this with the fact that the function $\mathbb{R}^d \ni v \mapsto (\nabla f)(v) \in \mathbb{R}^d$ is continuous establishes (14.121). The proof of Lemma 14.1.17 is thus complete. $\qquad\square$

**Lemma 14.1.18.** *Let $d \in \mathbb{N}$, $L \in (0, \infty)$, $r \in (0, \infty]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{w \in \mathbb{R}^d \colon \|w - \vartheta\|_2 \leq r\}$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ satisfy for all $\theta \in \mathbb{B}$ that*

$$\|(\nabla f)(\theta)\|_2 \leq L\|\theta - \vartheta\|_2. \tag{14.124}$$

*Then it holds for all $v, w \in \mathbb{B}$ that*

$$|f(v) - f(w)| \leq L \max\{\|v - \vartheta\|_2, \|w - \vartheta\|_2\}\|v - w\|_2. \tag{14.125}$$

*Proof of Lemma 14.1.18.* Observe that (14.124), the fundamental theorem of calculus, and the Cauchy-Schwarz inequality assure that for all $v, w \in \mathbb{B}$ it holds that

$$
\begin{aligned}
|f(v) - f(w)| &= \left| \left[ f(w + h(v - w)) \right]_{h=0}^{h=1} \right| \\
&= \left| \int_0^1 f'(w + h(v - w))(v - w) \, \mathrm{d}h \right| \\
&= \left| \int_0^1 \left\langle (\nabla f)(w + h(v - w)), v - w \right\rangle \mathrm{d}h \right| \\
&\leq \int_0^1 \|(\nabla f)(hv + (1 - h)w)\|_2 \|v - w\|_2 \, \mathrm{d}h \\
&\leq \int_0^1 L\|hv + (1 - h)w - \vartheta\|_2 \|v - w\|_2 \, \mathrm{d}h \\
&\leq \int_0^1 L(h\|v - \vartheta\|_2 + (1 - h)\|w - \vartheta\|_2)\|v - w\|_2 \, \mathrm{d}h \\
&= L\|v - w\|_2 \left[ \int_0^1 (h\|v - \vartheta\|_2 + h\|w - \vartheta\|_2) \, \mathrm{d}h \right] \\
&= L(\|v - \vartheta\|_2 + \|w - \vartheta\|_2)\|v - w\|_2 \left[ \int_0^1 h \, \mathrm{d}h \right] \\
&\leq L \max\{\|v - \vartheta\|_2, \|w - \vartheta\|_2\}\|v - w\|_2.
\end{aligned}
\tag{14.126}
$$

The proof of Lemma 14.1.18 is thus complete. $\qquad\square$

**Lemma 14.1.19.** *Let $d \in \mathbb{N}$, $L \in (0, \infty)$, $r \in (0, \infty]$, $\vartheta \in \mathbb{R}^d$, $\mathbb{B} = \{w \in \mathbb{R}^d \colon \|w - \vartheta\|_2 \leq r\}$, $f \in C^1(\mathbb{R}^d, \mathbb{R})$ satisfy for all $v, w \in \mathbb{B}$ that*

$$|f(v) - f(w)| \leq L \max\{\|v - \vartheta\|_2, \|w - \vartheta\|_2\}\|v - w\|_2. \tag{14.127}$$

*Then it holds for all $\theta \in \mathbb{B}$ that*

$$\|(\nabla f)(\theta)\|_2 \leq L\|\theta - \vartheta\|_2. \tag{14.128}$$

*Proof of Lemma 14.1.19.* Note that (14.127) implies that for all $\theta \in \mathbb{R}^d$ with $\|\theta - \vartheta\|_2 < r$

it holds that

$$
\begin{aligned}
&\|(\nabla f)(\theta)\|_2 \\
&= \sup_{w\in\mathbb{R}^d,\|w\|_2=1}\Big[f'(\theta)(w)\Big] \\
&= \sup_{w\in\mathbb{R}^d,\|w\|_2=1}\Big[\lim_{h\searrow0}\big[\tfrac{1}{h}(f(\theta+hw)-f(\theta))\big]\Big] \\
&\leq \sup_{w\in\mathbb{R}^d,\|w\|_2=1}\Big[\liminf_{h\searrow0}\big[\tfrac{L}{h}\max\{\|\theta+hw-\vartheta\|_2,\|\theta-\vartheta\|_2\}\|\theta+hw-\theta\|_2\big]\Big] \\
&= \sup_{w\in\mathbb{R}^d,\|w\|_2=1}\Big[\liminf_{h\searrow0}\big[L\max\{\|\theta+hw-\vartheta\|_2,\|\theta-\vartheta\|_2\}\tfrac{1}{h}\|hw\|_2\big]\Big] \\
&= \sup_{w\in\mathbb{R}^d,\|w\|_2=1}\Big[\liminf_{h\searrow0}\big[L\max\{\|\theta+hw-\vartheta\|_2,\|\theta-\vartheta\|_2\}\big]\Big] \\
&= \sup_{w\in\mathbb{R}^d,\|w\|_2=1}\Big[L\|\theta-\vartheta\|_2\Big]=L\|\theta-\vartheta\|_2.
\end{aligned}
\tag{14.129}
$$

The fact that the function $\mathbb{R}^d\ni v\mapsto(\nabla f)(v)\in\mathbb{R}^d$ is continuous therefore establishes (14.128). The proof of Lemma 14.1.19 is thus complete. □

**Corollary 14.1.20.** *Let $d\in\mathbb{N}$, $r\in(0,\infty]$, $\vartheta\in\mathbb{R}^d$, $\mathbb{B}=\{w\in\mathbb{R}^d\colon\|w-\vartheta\|_2\leq r\}$, $f\in C^1(\mathbb{R}^d,\mathbb{R})$. Then the following four statements are equivalent:*

*(i) There exist $c,L\in(0,\infty)$ such that for all $\theta\in\mathbb{B}$ it holds that*

$$
\langle\theta-\vartheta,(\nabla f)(\theta)\rangle\geq c\|\theta-\vartheta\|_2^2 \qquad and \qquad \|(\nabla f)(\theta)\|_2\leq L\|\theta-\vartheta\|_2. \tag{14.130}
$$

*(ii) There exist $\gamma\in(0,\infty)$, $\varepsilon\in(0,1)$ such that for all $\theta\in\mathbb{B}$ it holds that*

$$
\|\theta-\gamma(\nabla f)(\theta)-\vartheta\|_2\leq\varepsilon\|\theta-\vartheta\|_2. \tag{14.131}
$$

*(iii) There exists $c\in(0,\infty)$ such that for all $\theta\in\mathbb{B}$ it holds that*

$$
\langle\theta-\vartheta,(\nabla f)(\theta)\rangle\geq c\max\{\|\theta-\vartheta\|_2^2,\|(\nabla f)(\theta)\|_2^2\}. \tag{14.132}
$$

*(iv) There exist $c,L\in(0,\infty)$ such that for all $v,w\in\mathbb{B}$, $s,t\in[0,1]$ with $s\leq t$ it holds that*

$$
f(\vartheta+t(v-\vartheta))-f(\vartheta+s(v-\vartheta))\geq c(t^2-s^2)\|v-\vartheta\|_2^2 \tag{14.133}
$$

$$
and \qquad |f(v)-f(w)|\leq L\max\{\|v-\vartheta\|_2,\|w-\vartheta\|_2\}\|v-w\|_2. \tag{14.134}
$$

*Proof of Corollary 14.1.20.* First, note that items (ii)–(iii) in Lemma 14.1.4 prove that ((i) ⇒ (ii)). Next observe that Lemma 14.1.14 demonstrates that ((ii) ⇒ (iii)). Moreover, note that Lemma 14.1.15 establishes that ((iii) ⇒ (i)). In addition, observe that Lemma 14.1.16 and Lemma 14.1.18 show that ((i) ⇒ (iv)). Finally, note that Lemma 14.1.17 and Lemma 14.1.19 imply that ((iv) ⇒ (i)). The proof of Corollary 14.1.20 is thus complete. □

## 14.2 The gradient descent optimization method with classical momentum

**Definition 14.2.1** (Momentum gradient descent optimization method). *Let $d \in \mathbb{N}$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \infty)$, $(\alpha_n)_{n \in \mathbb{N}} \subseteq [0, 1]$, $\xi \in \mathbb{R}^d$ and let $f \colon \mathbb{R}^d \to \mathbb{R}$ and $g \colon \mathbb{R}^d \to \mathbb{R}^d$ satisfy for all $\theta \in \{v \in \mathbb{R}^d \colon (f \text{ is differentiable at } v)\}$ that*

$$g(\theta) = (\nabla f)(\theta). \tag{14.135}$$

*Then we say that $\Theta$ is the momentum gradient descent process for the objective function $f$ with generalized gradient $g$, learning rates $(\gamma_n)_{n \in \mathbb{N}}$, momentum decay factors $(\alpha_n)_{n \in \mathbb{N}}$, and initial value $\xi$ (we say that $\Theta$ is the momentum gradient descent process for the objective function $f$ with learning rates $(\gamma_n)_{n \in \mathbb{N}}$, momentum decay factors $(\alpha_n)_{n \in \mathbb{N}}$, and initial value $\xi$) if and only if it holds that $\Theta \colon \mathbb{N}_0 \to \mathbb{R}^d$ is the function from $\mathbb{N}_0$ to $\mathbb{R}^d$ which satisfies that there exists $\mathbf{m} \colon \mathbb{N}_0 \to \mathbb{R}^d$ such that for all $n \in \mathbb{N}$ it holds that*

$$\Theta_0 = \xi, \qquad \mathbf{m}_0 = 0, \tag{14.136}$$

$$\mathbf{m}_n = \alpha_n \mathbf{m}_{n-1} + (1 - \alpha_n) g(\Theta_{n-1}), \tag{14.137}$$

$$\text{and} \qquad \Theta_n = \Theta_{n-1} - \gamma_n \mathbf{m}_n. \tag{14.138}$$

### 14.2.1 A representation of the momentum GD optimization method

In (14.136)–(14.138) the momentum GD optimization method is formulated by means of a one-step recursion. This one-step recursion can efficiently be exploited in an implementation. The following elementary lemma, Lemma 14.2.2 below, provides a suitable full-history recursive representation for the momentum GD optimization method, which enables us to develop a better intuition for the momentum GD optimization method.

**Lemma 14.2.2** (A representation of the momentum GD optimization method). *Let $d \in \mathbb{N}$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq (0, \infty)$, $\alpha \in [0, 1]$, $\xi \in \mathbb{R}^d$, let $f \colon \mathbb{R}^d \to \mathbb{R}$ and $g \colon \mathbb{R}^d \to \mathbb{R}^d$ satisfy for all $\theta \in \{v \in \mathbb{R}^d \colon (f \text{ is differentiable at } v)\}$ that*

$$g(\theta) = (\nabla f)(\theta), \tag{14.139}$$

*and let $\Theta \colon \mathbb{N}_0 \to \mathbb{R}^d$ and $\mathbf{m} \colon \mathbb{N}_0 \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi, \qquad \mathbf{m}_0 = 0, \qquad \Theta_n = \Theta_{n-1} - \gamma_n \mathbf{m}_n, \tag{14.140}$$

$$\text{and} \qquad \mathbf{m}_n = \alpha \mathbf{m}_{n-1} + (1 - \alpha) g(\Theta_{n-1}). \tag{14.141}$$

*Then*

*(i) it holds for all $n \in \mathbb{N}_0$ that*

$$\mathbf{m}_n = (1 - \alpha) \left[ \sum_{k=0}^{n-1} \alpha^k g(\Theta_{n-1-k}) \right] \tag{14.142}$$

*and*

*(ii)* it holds for all $n \in \mathbb{N}$ that

$$\Theta_n = \Theta_{n-1} - \gamma_n (1-\alpha) \left[ \sum_{k=0}^{n-1} \alpha^k g(\Theta_{n-1-k}) \right]. \tag{14.143}$$

*Proof of Lemma 14.2.2.* We prove (14.142) by induction on $n \in \mathbb{N}_0$. For the base case $n = 0$ observe that (14.140) ensures that $\mathbf{m}_0 = (1-\alpha)0$. This establishes (14.142) in the base case $n = 0$. For the induction step observe that (14.141) assures that for all $n \in \mathbb{N}_0$ with

$$\mathbf{m}_n = (1-\alpha) \left[ \sum_{k=0}^{n-1} \alpha^k g(\Theta_{n-1-k}) \right] \tag{14.144}$$

it holds that

$$
\begin{aligned}
\mathbf{m}_{n+1} &= \alpha \mathbf{m}_n + (1-\alpha) g(\Theta_n) \\
&= \alpha \left[ (1-\alpha) \left[ \sum_{k=0}^{n-1} \alpha^k g(\Theta_{n-1-k}) \right] \right] + (1-\alpha) g(\Theta_n) \\
&= (1-\alpha) \left[ \sum_{k=1}^{n} \alpha^k g(\Theta_{n-k}) \right] + (1-\alpha) \alpha^0 g(\Theta_{n-0}) \\
&= (1-\alpha) \left[ \sum_{k=0}^{n} \alpha^k g(\Theta_{n-k}) \right] = (1-\alpha) \left[ \sum_{k=0}^{(n+1)-1} \alpha^k g(\Theta_{(n+1)-1-k}) \right].
\end{aligned}
\tag{14.145}
$$

Induction thus establishes (14.142). The proof of Lemma 14.2.2 is thus complete. □

## 14.2.2 Error analysis for the momentum GD optimization method in the case of quadratic objective functions

In this subsection we provide in Subsection 14.2.2.2 below an error analysis for the momentum GD optimization method in the case of quadratic objective functions (cf. Proposition 14.2.7 in Subsection 14.2.2.2 for the precise statement). In this specific case we also provide in Subsection 14.2.2.3 below a comparison of the convergence speeds of the plain vanilla GD optimization method and the momentum GD optimization method. In particular, we prove, roughly speaking, that the momentum GD optimization method outperfoms the plain vanilla GD optimization method in the specific case of quadratic objective functions; see Corollary 14.2.9 in Subsection 14.2.2.3 for the precise statement. For this comparison between the plain vanilla GD optimization method and the momentum GD optimization method we employ a refined error analysis of the plain vanilla GD optimization method in the case of quadratic objective functions. This refined error analysis is the subject of the next subsection (Subsection 14.2.2.1 below).

### 14.2.2.1 Error analysis for the GD optimization method in the case of quadratic objective functions

**Lemma 14.2.3** (Error analysis for the GD optimization method in the case of quadratic objective functions). *Let $d \in \mathbb{N}$, $\xi \in \mathbb{R}^d$, $\vartheta = (\vartheta_1, \vartheta_2, \ldots, \vartheta_d) \in \mathbb{R}^d$, $\kappa, \mathcal{K}, \lambda_1, \lambda_2, \ldots, \lambda_d \in$*

$(0, \infty)$ *satisfy* $\kappa = \min\{\lambda_1, \lambda_2, \ldots, \lambda_d\}$ *and* $\mathcal{K} = \max\{\lambda_1, \lambda_2, \ldots, \lambda_d\}$, *let* $f: \mathbb{R}^d \to \mathbb{R}$ *satisfy for all* $\theta = (\theta_1, \theta_2, \ldots, \theta_d) \in \mathbb{R}^d$ *that*

$$f(\theta) = \tfrac{1}{2}\left[\sum_{i=1}^{d} \lambda_i |\theta_i - \vartheta_i|^2\right], \tag{14.146}$$

*and let* $\Theta: \mathbb{N}_0 \to \mathbb{R}^d$ *satisfy for all* $n \in \mathbb{N}$ *that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \tfrac{2}{(\mathcal{K}+\kappa)}(\nabla f)(\Theta_{n-1}). \tag{14.147}$$

*Then it holds for all* $n \in \mathbb{N}_0$ *that*

$$\|\Theta_n - \vartheta\|_2 \leq \left[\tfrac{\mathcal{K}-\kappa}{\mathcal{K}+\kappa}\right]^n \|\xi - \vartheta\|_2. \tag{14.148}$$

*Proof of Lemma 14.2.3.* Throughout this proof let $\Theta^{(1)}, \Theta^{(2)}, \ldots, \Theta^{(d)}: \mathbb{N}_0 \to \mathbb{R}$ satisfy for all $n \in \mathbb{N}_0$ that $\Theta_n = (\Theta_n^{(1)}, \Theta_n^{(2)}, \ldots, \Theta_n^{(d)})$. Note that (14.146) implies that for all $\theta = (\theta_1, \theta_2, \ldots, \theta_d) \in \mathbb{R}^d$, $i \in \{1, 2, \ldots, d\}$ it holds that

$$\left(\tfrac{\partial f}{\partial \theta_i}\right)(\theta) = \lambda_i(\theta_i - \vartheta_i). \tag{14.149}$$

Combining this and (14.147) ensures that for all $n \in \mathbb{N}$, $i \in \{1, 2, \ldots, d\}$ it holds that

$$\begin{aligned}
\Theta_n^{(i)} - \vartheta_i &= \Theta_{n-1}^{(i)} - \tfrac{2}{(\mathcal{K}+\kappa)}\left(\tfrac{\partial f}{\partial \theta_i}\right)(\Theta_{n-1}) - \vartheta_i \\
&= \Theta_{n-1}^{(i)} - \vartheta_i - \tfrac{2}{(\mathcal{K}+\kappa)}\left[\lambda_i(\Theta_{n-1}^{(i)} - \vartheta_i)\right] \\
&= \left(1 - \tfrac{2\lambda_i}{(\mathcal{K}+\kappa)}\right)(\Theta_{n-1}^{(i)} - \vartheta_i).
\end{aligned} \tag{14.150}$$

Hence, we obtain that for all $n \in \mathbb{N}$ it holds that

$$\begin{aligned}
\|\Theta_n - \vartheta\|_2^2 &= \sum_{i=1}^{d} |\Theta_n^{(i)} - \vartheta_i|^2 \\
&= \sum_{i=1}^{d} \left[\left|1 - \tfrac{2\lambda_i}{(\mathcal{K}+\kappa)}\right|^2 |\Theta_{n-1}^{(i)} - \vartheta_i|^2\right] \\
&\leq \left[\max\left\{\left|1 - \tfrac{2\lambda_1}{(\mathcal{K}+\kappa)}\right|^2, \ldots, \left|1 - \tfrac{2\lambda_d}{(\mathcal{K}+\kappa)}\right|^2\right\}\right]\left[\sum_{i=1}^{d} |\Theta_{n-1}^{(i)} - \vartheta_i|^2\right] \\
&= \left[\max\left\{\left|1 - \tfrac{2\lambda_1}{(\mathcal{K}+\kappa)}\right|, \ldots, \left|1 - \tfrac{2\lambda_d}{(\mathcal{K}+\kappa)}\right|\right\}\right]^2 \|\Theta_{n-1} - \vartheta\|_2^2.
\end{aligned} \tag{14.151}$$

Moreover, note that the fact that for all $i \in \{1, 2, \ldots, d\}$ it holds that $\lambda_i \geq \kappa$ implies that for all $i \in \{1, 2, \ldots, d\}$ it holds that

$$1 - \tfrac{2\lambda_i}{(\mathcal{K}+\kappa)} \leq 1 - \tfrac{2\kappa}{(\mathcal{K}+\kappa)} = \tfrac{\mathcal{K}+\kappa-2\kappa}{\mathcal{K}+\kappa} = \tfrac{\mathcal{K}-\kappa}{\mathcal{K}+\kappa} \geq 0. \tag{14.152}$$

In addition, observe that the fact that for all $i \in \{1, 2, \ldots, d\}$ it holds that $\lambda_i \leq \mathcal{K}$ implies that for all $i \in \{1, 2, \ldots, d\}$ it holds that

$$1 - \tfrac{2\lambda_i}{(\mathcal{K}+\kappa)} \geq 1 - \tfrac{2\mathcal{K}}{(\mathcal{K}+\kappa)} = \tfrac{\mathcal{K}+\kappa-2\mathcal{K}}{(\mathcal{K}+\kappa)} = -\left[\tfrac{\mathcal{K}-\kappa}{\mathcal{K}+\kappa}\right] \leq 0. \tag{14.153}$$

This and (14.152) ensure that for all $i \in \{1, 2, \ldots, d\}$ it holds that

$$\left| 1 - \tfrac{2\lambda_i}{(\mathcal{K}+\kappa)} \right| \leq \tfrac{\mathcal{K}-\kappa}{\mathcal{K}+\kappa}. \tag{14.154}$$

Combining this with (14.151) demonstrates that for all $n \in \mathbb{N}$ it holds that

$$\begin{aligned}
\|\Theta_n - \vartheta\|_2 &\leq \left[ \max\left\{ \left| 1 - \tfrac{2\lambda_1}{\mathcal{K}+\kappa} \right|, \ldots, \left| 1 - \tfrac{2\lambda_d}{\mathcal{K}+\kappa} \right| \right\} \right] \|\Theta_{n-1} - \vartheta\|_2 \\
&\leq \left[ \tfrac{\mathcal{K}-\kappa}{\mathcal{K}+\kappa} \right] \|\Theta_{n-1} - \vartheta\|_2.
\end{aligned} \tag{14.155}$$

Induction therefore establishes that for all $n \in \mathbb{N}_0$ it holds that

$$\|\Theta_n - \vartheta\|_2 \leq \left[ \tfrac{\mathcal{K}-\kappa}{\mathcal{K}+\kappa} \right]^n \|\Theta_0 - \vartheta\|_2 = \left[ \tfrac{\mathcal{K}-\kappa}{\mathcal{K}+\kappa} \right]^n \|\xi - \vartheta\|_2. \tag{14.156}$$

The proof of Lemma 14.2.3 is thus complete. $\qquad\square$

Lemma 14.2.3 above establishes, roughly speaking, the convergence rate $\tfrac{\mathcal{K}-\kappa}{\mathcal{K}+\kappa}$ (see (14.148) above for the precise statement) for the GD optimization method in the case of the objective function (14.146). The next result, Lemma 14.2.4 below, essentially proves in the situation of Lemma 14.2.3 that this convergence rate cannot be improved by means of a difference choice of the learning rate.

**Lemma 14.2.4** (Lower bound for the convergence rate of gradient descent for quadratic objective functions)**.** *Let $d \in \mathbb{N}$, $\xi = (\xi_1, \xi_2, \ldots, \xi_d)$, $\vartheta = (\vartheta_1, \vartheta_2, \ldots, \vartheta_d) \in \mathbb{R}^d$, $\gamma, \kappa, \mathcal{K}, \lambda_1, \lambda_2 \ldots, \lambda_d \in (0, \infty)$ satisfy $\kappa = \min\{\lambda_1, \lambda_2, \ldots, \lambda_d\}$ and $\mathcal{K} = \max\{\lambda_1, \lambda_2, \ldots, \lambda_d\}$, let $f \colon \mathbb{R}^d \to \mathbb{R}$ satisfy for all $\theta = (\theta_1, \theta_2, \ldots, \theta_d) \in \mathbb{R}^d$ that*

$$f(\theta) = \tfrac{1}{2} \left[ \sum_{i=1}^{d} \lambda_i |\theta_i - \vartheta_i|^2 \right], \tag{14.157}$$

*and let $\Theta \colon \mathbb{N}_0 \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \gamma(\nabla f)(\Theta_{n-1}). \tag{14.158}$$

*Then it holds for all $n \in \mathbb{N}_0$ that*

$$\begin{aligned}
\|\Theta_n - \vartheta\|_2 &\geq \left[ \max\{\gamma\mathcal{K} - 1, 1 - \gamma\kappa\} \right]^n \left[ \min\{|\xi_1 - \vartheta_1|, \ldots, |\xi_d - \vartheta_d|\} \right] \\
&\geq \left[ \tfrac{\mathcal{K}-\kappa}{\mathcal{K}+\kappa} \right]^n \left[ \min\{|\xi_1 - \vartheta_1|, \ldots, |\xi_d - \vartheta_d|\} \right].
\end{aligned} \tag{14.159}$$

*Proof of Lemma 14.2.4.* Throughout this proof let $\Theta^{(1)}, \Theta^{(2)}, \ldots, \Theta^{(d)} \colon \mathbb{N}_0 \to \mathbb{R}$ satisfy for all $n \in \mathbb{N}_0$ that $\Theta_n = (\Theta_n^{(1)}, \Theta_n^{(2)}, \ldots, \Theta_n^{(d)})$ and let $\iota, \mathcal{I} \in \{1, 2, \ldots, d\}$ satisfy $\lambda_\iota = \kappa$ and $\lambda_{\mathcal{I}} = \mathcal{K}$. Observe that (14.157) implies that for all $\theta = (\theta_1, \theta_2, \ldots, \theta_d) \in \mathbb{R}^d$, $i \in \{1, 2, \ldots, d\}$ it holds that

$$\left( \tfrac{\partial f}{\partial \theta_i} \right)(\theta) = \lambda_i(\theta_i - \vartheta_i). \tag{14.160}$$

Combining this with (14.158) implies that for all $n \in \mathbb{N}$, $i \in \{1, 2, \ldots, d\}$ it holds that

$$\begin{aligned}
\Theta_n^{(i)} - \vartheta_i &= \Theta_{n-1}^{(i)} - \gamma \left( \tfrac{\partial f}{\partial \theta_i} \right)(\Theta_{n-1}) - \vartheta_i \\
&= \Theta_{n-1}^{(i)} - \vartheta_i - \gamma\lambda_i(\Theta_{n-1}^{(i)} - \vartheta_i) \\
&= (1 - \gamma\lambda_i)(\Theta_{n-1}^{(i)} - \vartheta_i).
\end{aligned} \tag{14.161}$$

Induction hence proves that for all $n \in \mathbb{N}_0$, $i \in \{1, 2, \ldots, d\}$ it holds that

$$\Theta_n^{(i)} - \vartheta_i = (1 - \gamma\lambda_i)^n(\Theta_0^{(i)} - \vartheta_i) = (1 - \gamma\lambda_i)^n(\xi_i - \vartheta_i). \tag{14.162}$$

This shows that for all $n \in \mathbb{N}_0$ it holds that

$$
\begin{aligned}
\|\Theta_n - \vartheta\|_2^2 &= \sum_{i=1}^d |\Theta_n^{(i)} - \vartheta_i|^2 = \sum_{i=1}^d \Big[|1 - \gamma\lambda_i|^{2n}|\xi_i - \vartheta_i|^2\Big] \\
&\geq \big[\min\{|\xi_1 - \vartheta_1|^2, \ldots, |\xi_d - \vartheta_d|^2\}\big] \Big[\sum_{i=1}^d |1 - \gamma\lambda_i|^{2n}\Big] \\
&\geq \big[\min\{|\xi_1 - \vartheta_1|^2, \ldots, |\xi_d - \vartheta_d|^2\}\big] \big[\max\{|1 - \gamma\lambda_1|^{2n}, \ldots, |1 - \gamma\lambda_d|^{2n}\}\big] \\
&= \big[\min\{|\xi_1 - \vartheta_1|, \ldots, |\xi_d - \vartheta_d|\}\big]^2 \big[\max\{|1 - \gamma\lambda_1|, \ldots, |1 - \gamma\lambda_d|\}\big]^{2n}.
\end{aligned}
\tag{14.163}
$$

Furthermore, note that

$$
\begin{aligned}
\max\{|1 - \gamma\lambda_1|, \ldots, |1 - \gamma\lambda_d|\} &\geq \max\{|1 - \gamma\lambda_{\mathcal{I}}|, |1 - \gamma\lambda_\iota|\} \\
&= \max\{|1 - \gamma\mathcal{K}|, |1 - \gamma\kappa|\} \geq \max\{\gamma\mathcal{K} - 1, 1 - \gamma\kappa\}.
\end{aligned}
\tag{14.164}
$$

In addition, observe that for all $\alpha \in (-\infty, \frac{2}{\mathcal{K}+\kappa}]$ it holds that

$$\max\{\alpha\mathcal{K} - 1, 1 - \alpha\kappa\} \geq 1 - \alpha\kappa \geq 1 - \big[\tfrac{2}{\mathcal{K}+\kappa}\big]\kappa = \tfrac{\mathcal{K}+\kappa-2\kappa}{\mathcal{K}+\kappa} = \tfrac{\mathcal{K}-\kappa}{\mathcal{K}+\kappa}. \tag{14.165}$$

Moreover, note that for all $\alpha \in [\frac{2}{\mathcal{K}+\kappa}, \infty)$ it holds that

$$\max\{\alpha\mathcal{K} - 1, 1 - \alpha\kappa\} \geq \alpha\mathcal{K} - 1 \geq \big[\tfrac{2}{\mathcal{K}+\kappa}\big]\mathcal{K} - 1 = \tfrac{2\mathcal{K}-(\mathcal{K}+\kappa)}{\mathcal{K}+\kappa} = \tfrac{\mathcal{K}-\kappa}{\mathcal{K}+\kappa}. \tag{14.166}$$

Combining this, (14.164), and (14.165) proves that

$$\max\{|1 - \gamma\lambda_1|, \ldots, |1 - \gamma\lambda_d|\} \geq \max\{\gamma\mathcal{K} - 1, 1 - \gamma\kappa\} \geq \tfrac{\mathcal{K}-\kappa}{\mathcal{K}+\kappa} \geq 0. \tag{14.167}$$

This and (14.163) demonstrate that for all $n \in \mathbb{N}_0$ it holds that

$$
\begin{aligned}
&\|\Theta_n - \vartheta\|_2 \\
&\geq \big[\max\{|1 - \gamma\lambda_1|, \ldots, |1 - \gamma\lambda_d|\}\big]^n \big[\min\{|\xi_1 - \vartheta_1|, \ldots, |\xi_d - \vartheta_d|\}\big] \\
&\geq \big[\max\{\gamma\mathcal{K} - 1, 1 - \gamma\kappa\}\big]^n \big[\min\{|\xi_1 - \vartheta_1|, \ldots, |\xi_d - \vartheta_d|\}\big] \\
&\geq \big[\tfrac{\mathcal{K}-\kappa}{\mathcal{K}+\kappa}\big]^n \big[\min\{|\xi_1 - \vartheta_1|, \ldots, |\xi_d - \vartheta_d|\}\big].
\end{aligned}
\tag{14.168}
$$

The proof of Lemma 14.2.4 is thus complete. $\qquad\square$

### 14.2.2.2   Error analysis for the momentum GD optimization method in the case of quadratic objective functions

In this subsection we provide in Proposition 14.2.7 below an error analysis for the momentum GD optimization method in the case of quadratic objective functions. Our proof of Proposition 14.2.7 employs the two auxiliary results on quadratic matrices in Lemma 14.2.5 and Lemma 14.2.6 below. Lemma 14.2.5 is a special case of the so-called Gelfand spectral radius formula in the literature. The proof of Lemma 14.2.5 can, e.g., be found in Tropp [29] and Einsiedler & Ward [10, Theorem 11.6]. Lemma 14.2.6 establishes a formula for the determinants of quadratic block matrices (see (14.170) below for the precise statement). Lemma 14.2.6 and its proof can, e.g., be found in Silvester [27, Theorem 3].

**Lemma 14.2.5** (A special case of Gelfand's spectral radius formula for real matrices)**.** *Let $d \in \mathbb{N}$, $A \in \mathbb{R}^{d \times d}$, $\mathscr{S} = \{\lambda \in \mathbb{C} \colon (\exists\, v \in \mathbb{C}^d \backslash \{0\} \colon Av = \lambda v)\}$ and let $\|\cdot\| \colon \mathbb{R}^d \to [0, \infty)$ be a norm. Then*

$$\liminf_{n \to \infty} \left( \left[ \sup_{v \in \mathbb{R}^d \backslash \{0\}} \frac{\|A^n v\|}{\|v\|} \right]^{1/n} \right) = \limsup_{n \to \infty} \left( \left[ \sup_{v \in \mathbb{R}^d \backslash \{0\}} \frac{\|A^n v\|}{\|v\|} \right]^{1/n} \right) \tag{14.169}$$
$$= \max_{\lambda \in \mathscr{S} \cup \{0\}} |\lambda|.$$

**Lemma 14.2.6** (Determinants for block matrices)**.** *Let $d \in \mathbb{N}$, $A, B, C, D \in \mathbb{R}^{d \times d}$ satisfy $CD = DC$. Then*

$$\det \underbrace{\begin{pmatrix} A & B \\ C & D \end{pmatrix}}_{\in \mathbb{R}^{2d \times 2d}} = \det(AD - BC) \tag{14.170}$$

*Proof of Lemma 14.2.6.* Throughout this proof let $D_x \in \mathbb{R}^{d \times d}$, $x \in \mathbb{R}$, satisfy for all $x \in \mathbb{R}$ that

$$D_x = D - x \, \mathrm{I}_d \tag{14.171}$$

(cf. Definition 2.2.9). Observe that the fact that for all $x \in \mathbb{R}$ it holds that $CD_x = D_x C$ and the fact that for all $X, Y, Z \in \mathbb{R}^{d \times d}$ it holds that

$$\det \begin{pmatrix} X & Y \\ 0 & Z \end{pmatrix} = \det(X) \det(Z) = \det \begin{pmatrix} X & 0 \\ Y & Z \end{pmatrix} \tag{14.172}$$

(cf., e.g., Petersen [24, Proposition 5.5.3 and Proposition 5.5.4]) imply that for all $x \in \mathbb{R}$ it holds that

$$\det \left( \begin{pmatrix} A & B \\ C & D_x \end{pmatrix} \begin{pmatrix} D_x & 0 \\ -C & \mathrm{I}_d \end{pmatrix} \right) = \det \begin{pmatrix} (AD_x - BC) & B \\ (CD_x - D_x C) & D_x \end{pmatrix}$$
$$= \det \begin{pmatrix} (AD_x - BC) & B \\ 0 & D_x \end{pmatrix} \tag{14.173}$$
$$= \det(AD_x - BC) \det(D_x).$$

Moreover, note that the multiplicative property of the determinant (see, e.g., Petersen [24, (1) in Proposition 5.5.2]) implies that for all $x \in \mathbb{R}$ it holds that

$$\det \left( \begin{pmatrix} A & B \\ C & D_x \end{pmatrix} \begin{pmatrix} D_x & 0 \\ -C & \mathrm{I}_d \end{pmatrix} \right) = \det \begin{pmatrix} A & B \\ C & D_x \end{pmatrix} \det \begin{pmatrix} D_x & 0 \\ -C & \mathrm{I}_d \end{pmatrix}$$
$$= \det \begin{pmatrix} A & B \\ C & D_x \end{pmatrix} \det(D_x) \det(\mathrm{I}_d) \tag{14.174}$$
$$= \det \begin{pmatrix} A & B \\ C & D_x \end{pmatrix} \det(D_x).$$

Combining this and (14.173) demonstrates that for all $x \in \mathbb{R}$ it holds that

$$\det \begin{pmatrix} A & B \\ C & D_x \end{pmatrix} \det(D_x) = \det(AD_x - BC) \det(D_x). \tag{14.175}$$

Hence, we obtain for all $x \in \mathbb{R}$ that

$$\left( \det \begin{pmatrix} A & B \\ C & D_x \end{pmatrix} - \det(AD_x - BC) \right) \det(D_x) = 0. \tag{14.176}$$

This implies that for all $x \in \mathbb{R}$ with $\det(D_x) \neq 0$ it holds that

$$\det \begin{pmatrix} A & B \\ C & D_x \end{pmatrix} - \det(AD_x - BC) = 0. \tag{14.177}$$

Moreover, note that the fact that $\{x \in \mathbb{R} \colon \det(D_x) = 0\} = \{x \in \mathbb{R} \colon \det(D - x\,\mathrm{I}_d) = 0\}$ is a finite set, the fact that the function

$$\mathbb{R} \ni x \mapsto \det \begin{pmatrix} A & B \\ C & D_x \end{pmatrix} - \det(AD_x - BC) \in \mathbb{R} \tag{14.178}$$

is continuous and (14.177) ensure that for all $x \in \mathbb{R}$ it holds that

$$\det \begin{pmatrix} A & B \\ C & D_x \end{pmatrix} - \det(AD_x - BC) = 0. \tag{14.179}$$

Hence, we obtain for all $x \in \mathbb{R}$ that

$$\det \begin{pmatrix} A & B \\ C & D_x \end{pmatrix} = \det(AD_x - BC). \tag{14.180}$$

This establishes that

$$\det \begin{pmatrix} A & B \\ C & D \end{pmatrix} = \det \begin{pmatrix} A & B \\ C & D_0 \end{pmatrix} = \det(AD_0 - BC) = \det(AD - BC). \tag{14.181}$$

The proof of Lemma 14.2.6 is thus completed. $\qquad \square$

We are now in the position to formulate and prove the promised error analysis for the momentum GD optimization method in the case of quadratic objective functions; see Proposition 14.2.7 below.

**Proposition 14.2.7** (Error analysis for the momentum GD optimization method in the case of quadratic objective functions). *Let $d \in \mathbb{N}$, $\xi \in \mathbb{R}^d$, $\vartheta = (\vartheta_1, \vartheta_2 \ldots, \vartheta_d) \in \mathbb{R}^d$, $\kappa, \mathcal{K}, \lambda_1, \lambda_2, \ldots, \lambda_d \in (0, \infty)$ satisfy $\kappa = \min\{\lambda_1, \lambda_2, \ldots, \lambda_d\}$ and $\mathcal{K} = \max\{\lambda_1, \lambda_2, \ldots, \lambda_d\}$, let $f \colon \mathbb{R}^d \to \mathbb{R}$ satisfy for all $\theta = (\theta_1, \theta_2, \ldots, \theta_d) \in \mathbb{R}^d$ that*

$$f(\theta) = \tfrac{1}{2} \left[ \sum_{i=1}^{d} \lambda_i |\theta_i - \vartheta_i|^2 \right], \tag{14.182}$$

*and let $\Theta \colon \mathbb{N}_0 \cup \{-1\} \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$ that $\Theta_{-1} = \Theta_0 = \xi$ and*

$$\Theta_n = \Theta_{n-1} - \tfrac{4}{(\sqrt{\mathcal{K}} + \sqrt{\kappa})^2} (\nabla f)(\Theta_{n-1}) + \left[ \tfrac{\sqrt{\mathcal{K}} - \sqrt{\kappa}}{\sqrt{\mathcal{K}} + \sqrt{\kappa}} \right]^2 (\Theta_{n-1} - \Theta_{n-2}). \tag{14.183}$$

*Then*

*(i) it holds that $\Theta|_{\mathbb{N}_0} \colon \mathbb{N}_0 \to \mathbb{R}^d$ is the momentum gradient descent process for the objective function $f$ with learning rates $\mathbb{N} \ni n \mapsto \frac{1}{\sqrt{\mathcal{K}\kappa}} \in [0, \infty)$, momentum decay factors $\mathbb{N} \ni n \mapsto \left[ \frac{\mathcal{K}^{1/2} - \kappa^{1/2}}{\mathcal{K}^{1/2} + \kappa^{1/2}} \right]^2 \in [0, 1]$, and initial value $\xi$ (cf. Definition 14.2.1) and*

*(ii) for every $\varepsilon \in (0, \infty)$ there exists $C \in (0, \infty)$ such that for all $n \in \mathbb{N}_0$ it holds that*

$$\|\Theta_n - \vartheta\|_2 \leq C \left[ \tfrac{\sqrt{\mathcal{K}} - \sqrt{\kappa}}{\sqrt{\mathcal{K}} + \sqrt{\kappa}} + \varepsilon \right]^n. \tag{14.184}$$

*Proof of Proposition 14.2.7.* Throughout this proof let $\varepsilon \in (0, \infty)$, let $\|\!|\!|\cdot|\!|\!|\colon \mathbb{R}^{2d \times 2d} \to [0, \infty)$ satisfy for all $B \in \mathbb{R}^{2d \times 2d}$ that

$$\|\!|\!|B|\!|\!| = \sup_{v \in \mathbb{R}^{2d} \setminus \{0\}} \left[ \frac{\|Bv\|_2}{\|v\|_2} \right], \tag{14.185}$$

let $\Theta^{(1)}, \Theta^{(2)}, \dots, \Theta^{(d)} \colon \mathbb{N}_0 \to \mathbb{R}$ satisfy for all $n \in \mathbb{N}_0$ that $\Theta_n = (\Theta_n^{(1)}, \Theta_n^{(2)}, \dots, \Theta_n^{(d)})$, let $\mathbf{m} \colon \mathbb{N}_0 \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}_0$ that

$$\mathbf{m}_n = -\sqrt{\mathcal{K}\kappa}(\Theta_n - \Theta_{n-1}), \tag{14.186}$$

let $\varrho \in (0, \infty)$, $\alpha \in [0, 1)$ be given by

$$\varrho = \frac{4}{(\sqrt{\mathcal{K}} + \sqrt{\kappa})^2} \qquad \text{and} \qquad \alpha = \left[ \frac{\sqrt{\mathcal{K}} - \sqrt{\kappa}}{\sqrt{\mathcal{K}} + \sqrt{\kappa}} \right]^2, \tag{14.187}$$

let $M \in \mathbb{R}^{d \times d}$ be the diagonal $(d \times d)$-matrix given by

$$M = \begin{pmatrix} (1 - \varrho\lambda_1 + \alpha) & & 0 \\ & \ddots & \\ 0 & & (1 - \varrho\lambda_d + \alpha) \end{pmatrix}, \tag{14.188}$$

let $A \in \mathbb{R}^{2d \times 2d}$ be the $((2d) \times (2d))$-matrix given by

$$A = \begin{pmatrix} M & (-\alpha \, \mathrm{I}_d) \\ \mathrm{I}_d & 0 \end{pmatrix}, \tag{14.189}$$

and let $\mathscr{S} \subseteq \mathbb{C}$ be the set given by

$$\mathscr{S} = \{ \mu \in \mathbb{C} \colon (\exists\, v \in \mathbb{C}^{2d} \setminus \{0\} \colon Av = \mu v) \} = \{ \mu \in \mathbb{C} \colon \det(A - \mu \, \mathrm{I}_{2d}) = 0 \} \tag{14.190}$$

(cf. Definition 2.2.9). Observe that (14.183), (14.186), and the fact that

$$\begin{aligned}
&\frac{(\sqrt{\mathcal{K}} + \sqrt{\kappa})^2 - (\sqrt{\mathcal{K}} - \sqrt{\kappa})^2}{4} \\
&= \tfrac{1}{4} \Big[ (\sqrt{\mathcal{K}} + \sqrt{\kappa} + \sqrt{\mathcal{K}} - \sqrt{\kappa})(\sqrt{\mathcal{K}} + \sqrt{\kappa} - [\sqrt{\mathcal{K}} - \sqrt{\kappa}]) \Big] \\
&= \tfrac{1}{4} \Big[ (2\sqrt{\mathcal{K}})(2\sqrt{\kappa}) \Big] = \sqrt{\mathcal{K}\kappa}
\end{aligned} \tag{14.191}$$

assure that for all $n \in \mathbb{N}$ it holds that

$$\begin{aligned}
\mathbf{m}_n &\\
&= -\sqrt{\mathcal{K}\kappa}(\Theta_n - \Theta_{n-1}) \\
&= -\sqrt{\mathcal{K}\kappa}\left( \Theta_{n-1} - \left[ \tfrac{4}{(\sqrt{\mathcal{K}} + \sqrt{\kappa})^2} \right](\nabla f)(\Theta_{n-1}) + \left[ \tfrac{\sqrt{\mathcal{K}} - \sqrt{\kappa}}{\sqrt{\mathcal{K}} + \sqrt{\kappa}} \right]^2 (\Theta_{n-1} - \Theta_{n-2}) - \Theta_{n-1} \right) \\
&= \sqrt{\mathcal{K}\kappa}\left( \left[ \tfrac{4}{(\sqrt{\mathcal{K}} + \sqrt{\kappa})^2} \right](\nabla f)(\Theta_{n-1}) - \left[ \tfrac{\sqrt{\mathcal{K}} - \sqrt{\kappa}}{\sqrt{\mathcal{K}} + \sqrt{\kappa}} \right]^2 (\Theta_{n-1} - \Theta_{n-2}) \right) \\
&= \frac{(\sqrt{\mathcal{K}} + \sqrt{\kappa})^2 - (\sqrt{\mathcal{K}} - \sqrt{\kappa})^2}{4} \left[ \tfrac{4}{(\sqrt{\mathcal{K}} + \sqrt{\kappa})^2} \right](\nabla f)(\Theta_{n-1}) \\
&\quad - \sqrt{\mathcal{K}\kappa}\left[ \tfrac{\sqrt{\mathcal{K}} - \sqrt{\kappa}}{\sqrt{\mathcal{K}} + \sqrt{\kappa}} \right]^2 (\Theta_{n-1} - \Theta_{n-2}) \\
&= \left[ 1 - \tfrac{(\sqrt{\mathcal{K}} - \sqrt{\kappa})^2}{(\sqrt{\mathcal{K}} + \sqrt{\kappa})^2} \right](\nabla f)(\Theta_{n-1}) + \left[ \tfrac{\sqrt{\mathcal{K}} - \sqrt{\kappa}}{\sqrt{\mathcal{K}} + \sqrt{\kappa}} \right]^2 \left[ -\sqrt{\mathcal{K}\kappa}(\Theta_{n-1} - \Theta_{n-2}) \right] \\
&= \left[ 1 - \left[ \tfrac{\sqrt{\mathcal{K}} - \sqrt{\kappa}}{\sqrt{\mathcal{K}} + \sqrt{\kappa}} \right]^2 \right](\nabla f)(\Theta_{n-1}) + \left[ \tfrac{\sqrt{\mathcal{K}} - \sqrt{\kappa}}{\sqrt{\mathcal{K}} + \sqrt{\kappa}} \right]^2 \mathbf{m}_{n-1}.
\end{aligned} \tag{14.192}$$

Moreover, note that (14.186) implies that for all $n \in \mathbb{N}_0$ it holds that

$$
\begin{aligned}
\Theta_n &= \Theta_{n-1} + (\Theta_n - \Theta_{n-1}) \\
&= \Theta_{n-1} - \tfrac{1}{\sqrt{\mathcal{K}\kappa}} \left( \left[ -\sqrt{\mathcal{K}\kappa} \right] (\Theta_n - \Theta_{n-1}) \right) = \Theta_{n-1} - \tfrac{1}{\sqrt{\mathcal{K}\kappa}} \mathbf{m}_n.
\end{aligned}
\tag{14.193}
$$

In addition, observe that the assumption that $\Theta_{-1} = \Theta_0 = \xi$ and (14.186) ensure that

$$
\mathbf{m}_0 = -\sqrt{\mathcal{K}\kappa} \left( \Theta_0 - \Theta_{-1} \right) = 0.
\tag{14.194}
$$

Combining this and the assumption that $\Theta_0 = \xi$ with (14.192) and (14.193) proves item (i). It thus remains to prove item (ii). For this observe that (14.182) implies that for all $\theta = (\theta_1, \theta_2, \ldots, \theta_d) \in \mathbb{R}^d$, $i \in \{1, 2, \ldots, d\}$ it holds that

$$
\left( \tfrac{\partial f}{\partial \theta_i} \right)(\theta) = \lambda_i (\theta_i - \vartheta_i).
\tag{14.195}
$$

This, (14.183), and (14.187) imply that for all $n \in \mathbb{N}$, $i \in \{1, 2, \ldots, d\}$ it holds that

$$
\begin{aligned}
\Theta_n^{(i)} &- \vartheta_i \\
&= \Theta_{n-1}^{(i)} - \varrho \left( \tfrac{\partial f}{\partial \theta_i} \right)(\Theta_{n-1}) + \alpha(\Theta_{n-1}^{(i)} - \Theta_{n-2}^{(i)}) - \vartheta_i \\
&= (\Theta_{n-1}^{(i)} - \vartheta_i) - \varrho \lambda_i (\Theta_{n-1}^{(i)} - \vartheta_i) + \alpha \left( (\Theta_{n-1}^{(i)} - \vartheta_i) - (\Theta_{n-2}^{(i)} - \vartheta_i) \right) \\
&= (1 - \varrho \lambda_i + \alpha)(\Theta_{n-1}^{(i)} - \vartheta_i) - \alpha(\Theta_{n-2}^{(i)} - \vartheta_i).
\end{aligned}
\tag{14.196}
$$

Combining this with (14.188) demonstrates that for all $n \in \mathbb{N}$ it holds that

$$
\begin{aligned}
\mathbb{R}^d \ni (\Theta_n - \vartheta) &= M(\Theta_{n-1} - \vartheta) - \alpha(\Theta_{n-2} - \vartheta) \\
&= \underbrace{\left( M \quad (-\alpha \, \mathrm{I}_d) \right)}_{\in \mathbb{R}^{d \times 2d}} \underbrace{\begin{pmatrix} \Theta_{n-1} - \vartheta \\ \Theta_{n-2} - \vartheta \end{pmatrix}}_{\in \mathbb{R}^{2d}}.
\end{aligned}
\tag{14.197}
$$

This and (14.189) assure that for all $n \in \mathbb{N}$ it holds that

$$
\mathbb{R}^{2d} \ni \begin{pmatrix} \Theta_n - \vartheta \\ \Theta_{n-1} - \vartheta \end{pmatrix} = \begin{pmatrix} M & (-\alpha \, \mathrm{I}_d) \\ \mathrm{I}_d & 0 \end{pmatrix} \begin{pmatrix} \Theta_{n-1} - \vartheta \\ \Theta_{n-2} - \vartheta \end{pmatrix} = A \begin{pmatrix} \Theta_{n-1} - \vartheta \\ \Theta_{n-2} - \vartheta \end{pmatrix}.
\tag{14.198}
$$

Induction hence proves that for all $n \in \mathbb{N}_0$ it holds that

$$
\mathbb{R}^{2d} \ni \begin{pmatrix} \Theta_n - \vartheta \\ \Theta_{n-1} - \vartheta \end{pmatrix} = A^n \begin{pmatrix} \Theta_0 - \vartheta \\ \Theta_{-1} - \vartheta \end{pmatrix} = A^n \begin{pmatrix} \xi - \vartheta \\ \xi - \vartheta \end{pmatrix}.
\tag{14.199}
$$

This, in turn, implies that for all $n \in \mathbb{N}_0$ it holds that

$$
\begin{aligned}
\|\Theta_n - \vartheta\|_2 &\leq \sqrt{\|\Theta_n - \vartheta\|_2^2 + \|\Theta_{n-1} - \vartheta\|_2^2} \\
&= \left\| \begin{pmatrix} \Theta_n - \vartheta \\ \Theta_{n-1} - \vartheta \end{pmatrix} \right\|_2 \\
&= \left\| A^n \begin{pmatrix} \xi - \vartheta \\ \xi - \vartheta \end{pmatrix} \right\|_2 \\
&\leq \|\!|A^n|\!\| \left\| \begin{pmatrix} \xi - \vartheta \\ \xi - \vartheta \end{pmatrix} \right\|_2 \\
&= \|\!|A^n|\!\| \sqrt{\|\xi - \vartheta\|_2^2 + \|\xi - \vartheta\|_2^2} \\
&= \|\!|A^n|\!\| \sqrt{2} \|\xi - \vartheta\|_2.
\end{aligned}
\tag{14.200}
$$

Next note that Lemma 14.2.5 demonstrates that

$$\limsup_{n\to\infty}\left(\left[\|\!|A^n|\!\|\right]^{1/n}\right)=\liminf_{n\to\infty}\left(\left[\|\!|A^n|\!\|\right]^{1/n}\right)=\max_{\mu\in\mathscr{S}\cup\{0\}}|\mu|. \tag{14.201}$$

This implies that there exists $m\in\mathbb{N}$ such that for all $n\in\mathbb{N}_0\cap[m,\infty)$ it holds that

$$\left[\|\!|A^n|\!\|\right]^{1/n}\leq\varepsilon+\max_{\mu\in\mathscr{S}\cup\{0\}}|\mu|. \tag{14.202}$$

Therefore, we obtain for all $n\in\mathbb{N}_0\cap[m,\infty)$ that

$$\|\!|A^n|\!\|\leq\left[\varepsilon+\max_{\mu\in\mathscr{S}\cup\{0\}}|\mu|\right]^n. \tag{14.203}$$

Furthermore, note that for all $n\in\mathbb{N}_0\cap[0,m)$ it holds that

$$
\begin{aligned}
&\|\!|A^n|\!\|\\
&=\left[\varepsilon+\max_{\mu\in\mathscr{S}\cup\{0\}}|\mu|\right]^n\left[\frac{\|\!|A^n|\!\|}{(\varepsilon+\max_{\mu\in\mathscr{S}\cup\{0\}}|\mu|)^n}\right]\\
&\leq\left[\varepsilon+\max_{\mu\in\mathscr{S}\cup\{0\}}|\mu|\right]^n\left[\max\left(\left\{\frac{\|\!|A^k|\!\|}{(\varepsilon+\max_{\mu\in\mathscr{S}\cup\{0\}}|\mu|)^k}:k\in\mathbb{N}_0\cap[0,m)\right\}\cup\{1\}\right)\right].
\end{aligned} \tag{14.204}
$$

Combining this and (14.203) proves that for all $n\in\mathbb{N}_0$ it holds that

$$
\begin{aligned}
&\|\!|A^n|\!\|\\
&\leq\left[\varepsilon+\max_{\mu\in\mathscr{S}\cup\{0\}}|\mu|\right]^n\left[\max\left(\left\{\frac{\|\!|A^k|\!\|}{(\varepsilon+\max_{\mu\in\mathscr{S}\cup\{0\}}|\mu|)^k}:k\in\mathbb{N}_0\cap[0,m)\right\}\cup\{1\}\right)\right].
\end{aligned} \tag{14.205}
$$

Next observe that Lemma 14.2.6, (14.189), and the fact that for all $\mu\in\mathbb{C}$ it holds that $\mathrm{I}_d(-\mu\,\mathrm{I}_d)=-\mu\,\mathrm{I}_d=(-\mu\,\mathrm{I}_d)\,\mathrm{I}_d$ ensure that for all $\mu\in\mathbb{C}$ it holds that

$$
\begin{aligned}
\det(A-\mu\,\mathrm{I}_{2d})&=\det\begin{pmatrix}(M-\mu\,\mathrm{I}_d)&(-\alpha\,\mathrm{I}_d)\\\mathrm{I}_d&-\mu\,\mathrm{I}_d\end{pmatrix}\\
&=\det\big((M-\mu\,\mathrm{I}_d)(-\mu\,\mathrm{I}_d)-(-\alpha\,\mathrm{I}_d)\,\mathrm{I}_d\big)\\
&=\det\big((M-\mu\,\mathrm{I}_d)(-\mu\,\mathrm{I}_d)+\alpha\,\mathrm{I}_d\big).
\end{aligned} \tag{14.206}
$$

This and (14.188) demonstrate that for all $\mu\in\mathbb{C}$ it holds that

$$
\begin{aligned}
&\det(A-\mu\,\mathrm{I}_{2d})\\
&=\det\begin{pmatrix}\big((1-\varrho\lambda_1+\alpha-\mu)(-\mu)+\alpha\big)&&0\\&\ddots&\\0&&\big((1-\varrho\lambda_d+\alpha-\mu)(-\mu)+\alpha\big)\end{pmatrix}\\
&=\prod_{i=1}^d\big((1-\varrho\lambda_i+\alpha-\mu)(-\mu)+\alpha\big)\\
&=\prod_{i=1}^d\big(\mu^2-(1-\varrho\lambda_i+\alpha)\mu+\alpha\big).
\end{aligned} \tag{14.207}
$$

Moreover, note that for all $\mu \in \mathbb{C}$, $i \in \{1, 2, \ldots, d\}$ it holds that

$$
\begin{aligned}
&\mu^2 - (1 - \varrho\lambda_i + \alpha)\mu + \alpha \\
&= \mu^2 - 2\mu\left[\frac{(1-\varrho\lambda_i+\alpha)}{2}\right] + \left[\frac{(1-\varrho\lambda_i+\alpha)}{2}\right]^2 + \alpha - \left[\frac{(1-\varrho\lambda_i+\alpha)}{2}\right]^2 \\
&= \left[\mu - \frac{(1-\varrho\lambda_i+\alpha)}{2}\right]^2 + \alpha - \tfrac{1}{4}[1 - \varrho\lambda_i + \alpha]^2 \\
&= \left[\mu - \frac{(1-\varrho\lambda_i+\alpha)}{2}\right]^2 - \tfrac{1}{4}\Big[[1 - \varrho\lambda_i + \alpha]^2 - 4\alpha\Big].
\end{aligned}
\tag{14.208}
$$

Hence, we obtain that for all $i \in \{1, 2, \ldots, d\}$ it holds that

$$
\begin{aligned}
&\left\{\mu \in \mathbb{C}\colon \mu^2 - (1 - \varrho\lambda_i + \alpha)\mu + \alpha = 0\right\} \\
&= \left\{\mu \in \mathbb{C}\colon \left[\mu - \frac{(1-\varrho\lambda_i+\alpha)}{2}\right]^2 = \tfrac{1}{4}\Big[[1 - \varrho\lambda_i + \alpha]^2 - 4\alpha\Big]\right\} \\
&= \left\{\frac{(1-\varrho\lambda_i+\alpha)+\sqrt{[1-\varrho\lambda_i+\alpha]^2-4\alpha}}{2}, \frac{(1-\varrho\lambda_i+\alpha)-\sqrt{[1-\varrho\lambda_i+\alpha]^2-4\alpha}}{2},\right\} \\
&= \bigcup_{s\in\{-1,1\}}\left\{\tfrac{1}{2}\left[1 - \varrho\lambda_i + \alpha + s\sqrt{(1 - \varrho\lambda_i + \alpha)^2 - 4\alpha}\right]\right\}.
\end{aligned}
\tag{14.209}
$$

Combining this, (14.190), and (14.207) demonstrates that

$$
\begin{aligned}
\mathscr{S} &= \{\mu \in \mathbb{C}\colon \det(A - \mu\,\mathrm{I}_{2d}) = 0\} \\
&= \left\{\mu \in \mathbb{C}\colon \left[\prod_{i=1}^{d}(\mu^2 - (1 - \varrho\lambda_i + \alpha)\mu + \alpha) = 0\right]\right\} \\
&= \bigcup_{i=1}^{d}\{\mu \in \mathbb{C}\colon \mu^2 - (1 - \varrho\lambda_i + \alpha)\mu + \alpha = 0\} \\
&= \bigcup_{i=1}^{d}\bigcup_{s\in\{-1,1\}}\left\{\tfrac{1}{2}\left[1 - \varrho\lambda_i + \alpha + s\sqrt{(1 - \varrho\lambda_i + \alpha)^2 - 4\alpha}\right]\right\}.
\end{aligned}
\tag{14.210}
$$

Moreover, observe that the fact that for all $i \in \{1, 2, \ldots, d\}$ it holds that $\lambda_i \geq \kappa$ and (14.187) ensure that for all $i \in \{1, 2, \ldots, d\}$ it holds that

$$
\begin{aligned}
1 - \varrho\lambda_i + \alpha &\leq 1 - \varrho\kappa + \alpha = 1 - \left[\frac{4}{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2}\right]\kappa + \frac{(\sqrt{\mathcal{K}}-\sqrt{\kappa})^2}{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2} \\
&= \frac{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2 - 4\kappa + (\sqrt{\mathcal{K}}-\sqrt{\kappa})^2}{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2} = \frac{\mathcal{K}+2\sqrt{\mathcal{K}}\sqrt{\kappa}+\kappa-4\kappa+\mathcal{K}-2\sqrt{\mathcal{K}}\sqrt{\kappa}+\kappa}{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2} \\
&= \frac{2\mathcal{K}-2\kappa}{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2} = \frac{2(\sqrt{\mathcal{K}}-\sqrt{\kappa})(\sqrt{\mathcal{K}}+\sqrt{\kappa})}{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2} = 2\left[\frac{\sqrt{\mathcal{K}}-\sqrt{\kappa}}{\sqrt{\mathcal{K}}+\sqrt{\kappa}}\right] \geq 0.
\end{aligned}
\tag{14.211}
$$

In addition, note that the fact that for all $i \in \{1, 2, \ldots, d\}$ it holds that $\lambda_i \leq \mathcal{K}$ and (14.187) assure that for all $i \in \{1, 2, \ldots, d\}$ it holds that

$$
\begin{aligned}
1 - \varrho\lambda_i + \alpha &\geq 1 - \varrho\mathcal{K} + \alpha = 1 - \left[\frac{4}{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2}\right]\mathcal{K} + \frac{(\sqrt{\mathcal{K}}-\sqrt{\kappa})^2}{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2} \\
&= \frac{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2 - 4\mathcal{K} + (\sqrt{\mathcal{K}}-\sqrt{\kappa})^2}{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2} = \frac{\mathcal{K}+2\sqrt{\mathcal{K}}\sqrt{\kappa}+\kappa-4\mathcal{K}+\mathcal{K}-2\sqrt{\mathcal{K}}\sqrt{\kappa}+\kappa}{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2} \\
&= \frac{-2\mathcal{K}+2\kappa}{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2} = -2\left[\frac{\mathcal{K}-\kappa}{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2}\right] = -2\left[\frac{(\sqrt{\mathcal{K}}-\sqrt{\kappa})(\sqrt{\mathcal{K}}+\sqrt{\kappa})}{(\sqrt{\mathcal{K}}+\sqrt{\kappa})^2}\right] \\
&= -2\left[\frac{\sqrt{\mathcal{K}}-\sqrt{\kappa}}{\sqrt{\mathcal{K}}+\sqrt{\kappa}}\right] \leq 0.
\end{aligned}
\tag{14.212}
$$

Combining this and (14.211) implies that for all $i \in \{1, 2, \ldots, d\}$ it holds that

$$(1 - \varrho\lambda_i + \alpha)^2 \le \left[2\left(\tfrac{\sqrt{\mathcal{K}}-\sqrt{\kappa}}{\sqrt{\mathcal{K}}+\sqrt{\kappa}}\right)\right]^2 = 4\left[\tfrac{\sqrt{\mathcal{K}}-\sqrt{\kappa}}{\sqrt{\mathcal{K}}+\sqrt{\kappa}}\right]^2 = 4\alpha. \tag{14.213}$$

This and (14.210) demonstrate that

$$\begin{aligned}
&\max_{\mu \in \mathscr{S} \cup \{0\}} |\mu| = \max_{\mu \in \mathscr{S}} |\mu| \\
&= \max_{i \in \{1,2,\ldots,d\}} \max_{s \in \{-1,1\}} \left| \frac{1}{2}\left[1 - \varrho\lambda_i + \alpha + s\sqrt{(1 - \varrho\lambda_i + \alpha)^2 - 4\alpha}\right]\right| \\
&= \frac{1}{2}\left[\max_{i \in \{1,2,\ldots,d\}} \max_{s \in \{-1,1\}} \left|\left[1 - \varrho\lambda_i + \alpha + s\sqrt{(-1)(4\alpha - [1 - \varrho\lambda_i + \alpha]^2)}\right]\right|\right] \\
&= \frac{1}{2}\left[\max_{i \in \{1,2,\ldots,d\}} \max_{s \in \{-1,1\}} \left|\left[1 - \varrho\lambda_i + \alpha + s\mathbf{i}\sqrt{4\alpha - (1 - \varrho\lambda_i + \alpha)^2}\right]\right|^2\right]^{1/2}.
\end{aligned} \tag{14.214}$$

Combining this with (14.213) proves that

$$\begin{aligned}
&\max_{\mu \in \mathscr{S} \cup \{0\}} |\mu| \\
&= \tfrac{1}{2}\left[\max_{i \in \{1,2,\ldots,d\}} \max_{s \in \{-1,1\}} \left(\left|1 - \varrho\lambda_i + \alpha\right|^2 + \left|s\sqrt{4\alpha - (1 - \varrho\lambda_i + \alpha)^2}\right|^2\right)\right]^{1/2} \\
&= \tfrac{1}{2}\left[\max_{i \in \{1,2,\ldots,d\}} \max_{s \in \{-1,1\}} \left((1 - \varrho\lambda_i + \alpha)^2 + 4\alpha - (1 - \varrho\lambda_i + \alpha)^2\right)\right]^{1/2} \\
&= \tfrac{1}{2}[4\alpha]^{1/2} = \sqrt{\alpha}.
\end{aligned} \tag{14.215}$$

Combining (14.200) and (14.205) hence ensures that for all $n \in \mathbb{N}_0$ it holds that

$$\begin{aligned}
&\left\|\Theta_n - \vartheta\right\|_2 \\
&\le \sqrt{2}\left\|\xi - \vartheta\right\|_2 \|A^n\| \\
&\le \sqrt{2}\left\|\xi - \vartheta\right\|_2 \left[\varepsilon + \max_{\mu \in \mathscr{S} \cup \{0\}} |\mu|\right]^n \\
&\quad \cdot \left[\max\left(\left\{\tfrac{\|A^k\|}{(\varepsilon + \max_{\mu \in \mathscr{S} \cup \{0\}} |\mu|)^k} \in \mathbb{R} \colon k \in \mathbb{N}_0 \cap [0, m)\right\} \cup \{1\}\right)\right] \\
&= \sqrt{2}\left\|\xi - \vartheta\right\|_2 \left[\varepsilon + \alpha^{1/2}\right]^n \left[\max\left(\left\{\tfrac{\|A^k\|}{(\varepsilon + \alpha^{1/2})^k} \in \mathbb{R} \colon k \in \mathbb{N}_0 \cap [0, m)\right\} \cup \{1\}\right)\right] \\
&= \sqrt{2}\left\|\xi - \vartheta\right\|_2 \left[\varepsilon + \tfrac{\sqrt{\mathcal{K}}-\sqrt{\kappa}}{\sqrt{\mathcal{K}}+\sqrt{\kappa}}\right]^n \left[\max\left(\left\{\tfrac{\|A^k\|}{(\varepsilon + \alpha^{1/2})^k} \in \mathbb{R} \colon k \in \mathbb{N}_0 \cap [0, m)\right\} \cup \{1\}\right)\right].
\end{aligned} \tag{14.216}$$

This establishes item (ii). The proof of Proposition 14.2.7 it thus completed. □

### 14.2.2.3 Comparison of the convergence speeds of the GD optimization method with and without classical momentum

In this subsection we provide in Corollary 14.2.9 below a comparison between the convergence speeds of the plain vanilla GD optimization method and the momentum GD optimization method. Our proof of Corollary 14.2.9 employs the auxiliary and elementary estimate in Lemma 14.2.8 below, the refined error analysis for the plain vanilla GD optimization in Subsection 14.2.2.1 above (see Lemma 14.2.3 and Lemma 14.2.4 in Subsection 14.2.2.1), as well as the error analysis for the momentum GD optimization method in Subsection 14.2.2.2 above (see Proposition 14.2.7 in Subsection 14.2.2.2).

**Lemma 14.2.8** (Comparison of the convergence rates of the GD optimization method and the momentum GD optimization method)**.** *Let $\mathcal{K}, \kappa \in (0, \infty)$ satisfy $\kappa < \mathcal{K}$. Then*

$$\frac{\sqrt{\mathcal{K}} - \sqrt{\kappa}}{\sqrt{\mathcal{K}} + \sqrt{\kappa}} < \frac{\mathcal{K} - \kappa}{\mathcal{K} + \kappa}. \tag{14.217}$$

*Proof of Lemma 14.2.8.* Note that the fact that $\mathcal{K} - \kappa > 0 < 2\sqrt{\mathcal{K}}\sqrt{\kappa}$ ensures that

$$\frac{\sqrt{\mathcal{K}} - \sqrt{\kappa}}{\sqrt{\mathcal{K}} + \sqrt{\kappa}} = \frac{(\sqrt{\mathcal{K}} - \sqrt{\kappa})(\sqrt{\mathcal{K}} + \sqrt{\kappa})}{(\sqrt{\mathcal{K}} + \sqrt{\kappa})^2} = \frac{\mathcal{K} - \kappa}{\mathcal{K} + 2\sqrt{\mathcal{K}}\sqrt{\kappa} + \kappa} < \frac{\mathcal{K} - \kappa}{\mathcal{K} + \kappa}. \tag{14.218}$$

The proof of Lemma 14.2.8 it thus completed. □

**Corollary 14.2.9** (Convergence rate comparisons between the GD optimization method and the momentum GD optimization method)**.** *Let $d \in \mathbb{N}$, $\kappa, \mathcal{K}, \lambda_1, \lambda_2, \ldots, \lambda_d \in (0, \infty)$, $\xi = (\xi_1, \xi_2, \ldots, \xi_d)$, $\vartheta = (\vartheta_1, \vartheta_2, \ldots, \vartheta_d) \in \mathbb{R}^d$ satisfy*

$$\kappa = \min\{\lambda_1, \lambda_2, \ldots, \lambda_d\} < \max\{\lambda_1, \lambda_2, \ldots, \lambda_d\} = \mathcal{K}, \tag{14.219}$$

*let $f \colon \mathbb{R}^d \to \mathbb{R}$ satisfy for all $\theta = (\theta_1, \theta_2, \ldots, \theta_d) \in \mathbb{R}^d$ that*

$$f(\theta) = \tfrac{1}{2}\left[\sum_{i=1}^{d} \lambda_i |\theta_i - \vartheta_i|^2\right], \tag{14.220}$$

*let $\Theta^\gamma \colon \mathbb{N}_0 \to \mathbb{R}^d$, $\gamma \in (0, \infty)$, satisfy for all $\gamma \in (0, \infty)$, $n \in \mathbb{N}$ that*

$$\Theta_0^\gamma = \xi \qquad and \qquad \Theta_n^\gamma = \Theta_{n-1}^\gamma - \gamma(\nabla f)(\Theta_{n-1}^\gamma), \tag{14.221}$$

*and let $\mathcal{M} \colon \mathbb{N}_0 \cup \{-1\} \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$ that $\mathcal{M}_{-1} = \mathcal{M}_0 = \xi$ and*

$$\mathcal{M}_n = \mathcal{M}_{n-1} - \tfrac{4}{(\sqrt{\mathcal{K}} + \sqrt{\kappa})^2} (\nabla f)(\mathcal{M}_{n-1}) + \left[\tfrac{\sqrt{\mathcal{K}} - \sqrt{\kappa}}{\sqrt{\mathcal{K}} + \sqrt{\kappa}}\right]^2 (\mathcal{M}_{n-1} - \mathcal{M}_{n-2}). \tag{14.222}$$

*Then*

*(i) there exist $\gamma, C \in (0, \infty)$ such that for all $n \in \mathbb{N}_0$ it holds that*

$$\|\Theta_n^\gamma - \vartheta\|_2 \le C\left[\tfrac{\mathcal{K} - \kappa}{\mathcal{K} + \kappa}\right]^n, \tag{14.223}$$

*(ii) it holds for all $\gamma \in (0, \infty), n \in \mathbb{N}_0$ that*

$$\|\Theta_n^\gamma - \vartheta\|_2 \ge \left[\min\{|\xi_1 - \vartheta_1|, \ldots, |\xi_d - \vartheta_d|\}\right]\left[\tfrac{\mathcal{K} - \kappa}{\mathcal{K} + \kappa}\right]^n, \tag{14.224}$$

*(iii) for every $\varepsilon \in (0, \infty)$ there exists $C \in (0, \infty)$ such that for all $n \in \mathbb{N}_0$ it holds that*

$$\|\mathcal{M}_n - \vartheta\|_2 \le C\left[\tfrac{\sqrt{\mathcal{K}} - \sqrt{\kappa}}{\sqrt{\mathcal{K}} + \sqrt{\kappa}} + \varepsilon\right]^n, \tag{14.225}$$

*and*

*(iv) it holds that $\tfrac{\sqrt{\mathcal{K}} - \sqrt{\kappa}}{\sqrt{\mathcal{K}} + \sqrt{\kappa}} < \tfrac{\mathcal{K} - \kappa}{\mathcal{K} + \kappa}$.*

*Proof of Corollary 14.2.9.* First, note that Lemma 14.2.3 proves item (i). Next observe that Lemma 14.2.4 establishes item (ii). In addition, note that Proposition 14.2.7 proves item (iii). Finally, observe that Lemma 14.2.8 establishes item (iv). The proof of Corollary 14.2.9 is thus complete. □

Corollary 14.2.9 above, roughly speaking, shows in the case of quadratic objective functions that the momentum GD optimization method in (14.222) outperforms the classical plain vanilla GD optimization method (and, in particular, the classical plain vanilla GD optimization method in (14.147) in Lemma 14.2.3 above) provided that the parameters $\lambda_1, \lambda_2, \ldots, \lambda_d \in (0, \infty)$ in the objective function in (14.220) satisfy the assumption that $\min\{\lambda_1, \ldots, \lambda_d\} < \max\{\lambda_1, \ldots, \lambda_d\}$. The next elementary result, Lemma 14.2.10 below, demonstrates that the momentum GD optimization method in (14.222) and the plain vanilla GD optimization method in (14.147) in Lemma 14.2.3 above coincide in the case where $\min\{\lambda_1, \ldots, \lambda_d\} = \max\{\lambda_1, \ldots, \lambda_d\}$.

**Lemma 14.2.10** (Concurrence of the GD optimization method and the momentum GD optimization method). *Let $d \in \mathbb{N}$, $\xi, \vartheta \in \mathbb{R}^d$, $\alpha \in (0, \infty)$, let $f \colon \mathbb{R}^d \to \mathbb{R}$ satisfy for all $\theta \in \mathbb{R}^d$ that*

$$f(\theta) = \tfrac{\alpha}{2} \|\theta - \vartheta\|_2^2, \tag{14.226}$$

*let $\Theta \colon \mathbb{N}_0 \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \tfrac{2}{(\alpha + \alpha)} (\nabla f)(\Theta_{n-1}), \tag{14.227}$$

*and let $\mathcal{M} \colon \mathbb{N}_0 \cup \{-1\} \to \mathbb{R}^d$ satisfy for all $n \in \mathbb{N}$ that $\mathcal{M}_{-1} = \mathcal{M}_0 = \xi$ and*

$$\mathcal{M}_n = \mathcal{M}_{n-1} - \tfrac{4}{(\sqrt{\alpha} + \sqrt{\alpha})^2} (\nabla f)(\mathcal{M}_{n-1}) + \left[ \tfrac{\sqrt{\alpha} - \sqrt{\alpha}}{\sqrt{\alpha} + \sqrt{\alpha}} \right]^2 (\mathcal{M}_{n-1} - \mathcal{M}_{n-2}). \tag{14.228}$$

*Then*

(i) *it holds that $\mathcal{M}|_{\mathbb{N}_0} \colon \mathbb{N}_0 \to \mathbb{R}^d$ is the momentum gradient descent process for the objective function $f$ with learning rates $\mathbb{N} \ni n \mapsto 1/\alpha \in [0, \infty)$, momentum decay factors $\mathbb{N} \ni n \mapsto 0 \in [0, 1]$, and initial value $\xi$ (cf. Definition 14.2.1),*

(ii) *it holds for all $n \in \mathbb{N}_0$ that $\mathcal{M}_n = \Theta_n$, and*

(iii) *it holds for all $n \in \mathbb{N}$ that $\Theta_n = \vartheta = \mathcal{M}_n$.*

*Proof of Lemma 14.2.10.* First, note that (14.228) implies that for all $n \in \mathbb{N}$ it holds that

$$\mathcal{M}_n = \mathcal{M}_{n-1} - \tfrac{4}{(2\sqrt{\alpha})^2}(\nabla f)(\mathcal{M}_{n-1}) = \mathcal{M}_{n-1} - \tfrac{1}{\alpha}(\nabla f)(\mathcal{M}_{n-1}). \tag{14.229}$$

Combining this with the assumption that $\mathcal{M}_0 = \xi$ establishes item (i). Next note that (14.227) ensures that for all $n \in \mathbb{N}$ it holds that

$$\Theta_n = \Theta_{n-1} - \tfrac{1}{\alpha}(\nabla f)(\Theta_{n-1}). \tag{14.230}$$

Combining this with (14.229) and the assumption that $\Theta_0 = \xi = \mathcal{M}_0$ proves item (ii). Furthermore, observe that Lemma 13.2.4 assures that for all $\theta \in \mathbb{R}^d$ it holds that

$$(\nabla f)(\theta) = \tfrac{\alpha}{2}(2(\theta - \vartheta)) = \alpha(\theta - \vartheta). \tag{14.231}$$

Next we claim that for all $n \in \mathbb{N}$ it holds that

$$\Theta_n = \vartheta. \tag{14.232}$$

We now prove (14.232) by induction on $n \in \mathbb{N}$. For the base case $n = 1$ note that (14.230) and (14.231) imply that

$$\Theta_1 = \Theta_0 - \tfrac{1}{\alpha}(\nabla f)(\Theta_0) = \xi - \tfrac{1}{\alpha}(\alpha(\xi - \vartheta)) = \xi - (\xi - \vartheta) = \vartheta. \tag{14.233}$$

This establishes (14.232) in the base case $n = 1$. For the induction step observe that (14.230) and (14.231) assure that for all $n \in \mathbb{N}$ with $\Theta_n = \vartheta$ it holds that

$$\Theta_{n+1} = \Theta_n - \tfrac{1}{\alpha}(\nabla f)(\Theta_n) = \vartheta - \tfrac{1}{\alpha}(\alpha(\vartheta - \vartheta)) = \vartheta. \tag{14.234}$$

Induction thus proves (14.232). Combining (14.232) and item (ii) establishes item (iii). The proof of Lemma 14.2.10 is thus complete. $\qquad\square$

## 14.2.3 Comparison of the GD optimization method with and without momentum in the case of a numerical example

In this subsection we provide a numerical comparison of the plain vanilla GD optimization method and the momentum GD optimization method in the case of the specific quadratic optimization problem in (14.235)–(14.236) below; see Illustration 14.2.11 below, PYTHON code 14.1, and Figure 14.1 below.

**Illustration 14.2.11.** *Let* $\mathcal{K} = 10$, $\kappa = 1$, $\vartheta = (\vartheta_1, \vartheta_2) \in \mathbb{R}^2$, $\xi = (\xi_1, \xi_2) \in \mathbb{R}^2$ *satisfy*

$$\vartheta = \begin{pmatrix} \vartheta_1 \\ \vartheta_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \qquad and \qquad \xi = \begin{pmatrix} \xi_1 \\ \xi_2 \end{pmatrix} = \begin{pmatrix} 5 \\ 3 \end{pmatrix}, \tag{14.235}$$

*let* $f\colon \mathbb{R}^2 \to \mathbb{R}$ *satisfy for all* $\theta = (\theta_1, \theta_2) \in \mathbb{R}^2$ *that*

$$f(\theta) = \left(\tfrac{\kappa}{2}\right)|\theta_1 - \vartheta_1|^2 + \left(\tfrac{\mathcal{K}}{2}\right)|\theta_2 - \vartheta_2|^2, \tag{14.236}$$

*let* $\Theta\colon \mathbb{N}_0 \to \mathbb{R}^d$ *satisfy for all* $n \in \mathbb{N}$ *that* $\Theta_0 = \xi$ *and*

$$\begin{aligned} \Theta_n &= \Theta_{n-1} - \tfrac{2}{(\mathcal{K}+\kappa)}(\nabla f)(\Theta_{n-1}) = \Theta_{n-1} - \tfrac{2}{11}(\nabla f)(\Theta_{n-1}) \\ &= \Theta_{n-1} - 0.\overline{18}\,(\nabla f)(\Theta_{n-1}) \approx \Theta_{n-1} - 0.18\,(\nabla f)(\Theta_{n-1}), \end{aligned} \tag{14.237}$$

*and let* $\mathcal{M}, \mathbf{m}\colon \mathbb{N}_0 \to \mathbb{R}^d$ *satisfy for all* $n \in \mathbb{N}$ *that* $\mathcal{M}_0 = \xi$, $\mathbf{m}_0 = 0$, $\mathcal{M}_n = \mathcal{M}_{n-1} - 0.3\,\mathbf{m}_n$, *and*

$$\begin{aligned} \mathbf{m}_n &= 0.5\,\mathbf{m}_{n-1} + (1 - 0.5)\,(\nabla f)(\mathcal{M}_{n-1}) \\ &= 0.5\,(\mathbf{m}_{n-1} + (\nabla f)(\mathcal{M}_{n-1})). \end{aligned} \tag{14.238}$$

*Then*

(i) *it holds for all* $\theta = (\theta_1, \theta_2) \in \mathbb{R}^2$ *that*

$$(\nabla f)(\theta) = \begin{pmatrix} \kappa(\theta_1 - \vartheta_1) \\ \mathcal{K}(\theta_2 - \vartheta_2) \end{pmatrix} = \begin{pmatrix} \theta_1 - 1 \\ 10\,(\theta_2 - 1) \end{pmatrix}, \tag{14.239}$$

*(ii) it holds that*

$$\Theta_0 = \begin{pmatrix} 5 \\ 3 \end{pmatrix}, \tag{14.240}$$

$$
\begin{aligned}
\Theta_1 &= \Theta_0 - \tfrac{2}{11}(\nabla f)(\Theta_0) \approx \Theta_0 - 0.18(\nabla f)(\Theta_0) \\
&= \begin{pmatrix} 5 \\ 3 \end{pmatrix} - 0.18 \begin{pmatrix} 5-1 \\ 10(3-1) \end{pmatrix} = \begin{pmatrix} 5 - 0.18 \cdot 4 \\ 3 - 0.18 \cdot 10 \cdot 2 \end{pmatrix} \\
&= \begin{pmatrix} 5 - 0.72 \\ 3 - 3.6 \end{pmatrix} = \begin{pmatrix} 4.28 \\ -0.6 \end{pmatrix},
\end{aligned}
\tag{14.241}
$$

$$
\begin{aligned}
\Theta_2 &\approx \Theta_1 - 0.18(\nabla f)(\Theta_1) = \begin{pmatrix} 4.28 \\ -0.6 \end{pmatrix} - 0.18 \begin{pmatrix} 4.28 - 1 \\ 10(-0.6 - 1) \end{pmatrix} \\
&= \begin{pmatrix} 4.28 - 0.18 \cdot 3.28 \\ -0.6 - 0.18 \cdot 10 \cdot (-1.6) \end{pmatrix} = \begin{pmatrix} 4.10 - 0.18 \cdot 2 - 0.18 \cdot 0.28 \\ -0.6 + 1.8 \cdot 1.6 \end{pmatrix} \\
&= \begin{pmatrix} 4.10 - 0.36 - 2 \cdot 9 \cdot 4 \cdot 7 \cdot 10^{-4} \\ -0.6 + 1.6 \cdot 1.6 + 0.2 \cdot 1.6 \end{pmatrix} = \begin{pmatrix} 3.74 - 9 \cdot 56 \cdot 10^{-4} \\ -0.6 + 2.56 + 0.32 \end{pmatrix} \\
&= \begin{pmatrix} 3.74 - 504 \cdot 10^{-4} \\ 2.88 - 0.6 \end{pmatrix} = \begin{pmatrix} 3.6896 \\ 2.28 \end{pmatrix} \approx \begin{pmatrix} 3.69 \\ 2.28 \end{pmatrix},
\end{aligned}
\tag{14.242}
$$

$$
\begin{aligned}
\Theta_3 &\approx \Theta_2 - 0.18(\nabla f)(\Theta_2) \approx \begin{pmatrix} 3.69 \\ 2.28 \end{pmatrix} - 0.18 \begin{pmatrix} 3.69 - 1 \\ 10(2.28 - 1) \end{pmatrix} \\
&= \begin{pmatrix} 3.69 - 0.18 \cdot 2.69 \\ 2.28 - 0.18 \cdot 10 \cdot 1.28 \end{pmatrix} = \begin{pmatrix} 3.69 - 0.2 \cdot 2.69 + 0.02 \cdot 2.69 \\ 2.28 - 1.8 \cdot 1.28 \end{pmatrix} \\
&= \begin{pmatrix} 3.69 - 0.538 + 0.0538 \\ 2.28 - 1.28 - 0.8 \cdot 1.28 \end{pmatrix} = \begin{pmatrix} 3.7438 - 0.538 \\ 1 - 1.28 + 0.2 \cdot 1.28 \end{pmatrix} \\
&= \begin{pmatrix} 3.2058 \\ 0.256 - 0.280 \end{pmatrix} = \begin{pmatrix} 3.2058 \\ -0.024 \end{pmatrix} \approx \begin{pmatrix} 3.21 \\ -0.02 \end{pmatrix},
\end{aligned}
\tag{14.243}
$$

$$\vdots$$

*and*

*(iii) it holds that*

$$\mathcal{M}_0 = \begin{pmatrix} 5 \\ 3 \end{pmatrix}, \tag{14.244}$$

$$
\begin{aligned}
\mathbf{m}_1 &= 0.5 \, (\mathbf{m}_0 + (\nabla f)(\mathcal{M}_0)) = 0.5 \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix} + \begin{pmatrix} 5-1 \\ 10(3-1) \end{pmatrix} \right) \\
&= \begin{pmatrix} 0.5\,(0+4) \\ 0.5\,(0 + 10 \cdot 2) \end{pmatrix} = \begin{pmatrix} 2 \\ 10 \end{pmatrix},
\end{aligned}
\tag{14.245}
$$

$$\mathcal{M}_1 = \mathcal{M}_0 - 0.3\,\mathbf{m}_1 = \begin{pmatrix} 5 \\ 3 \end{pmatrix} - 0.3 \begin{pmatrix} 2 \\ 10 \end{pmatrix} = \begin{pmatrix} 4.4 \\ 0 \end{pmatrix}, \qquad (14.246)$$

$$\begin{aligned} \mathbf{m}_2 &= 0.5\,(\mathbf{m}_1 + (\nabla f)(\mathcal{M}_1)) = 0.5\left(\begin{pmatrix} 2 \\ 10 \end{pmatrix} + \begin{pmatrix} 4.4 - 1 \\ 10(0 - 1) \end{pmatrix}\right) \\ &= \begin{pmatrix} 0.5\,(2 + 3.4) \\ 0.5\,(10 - 10) \end{pmatrix} = \begin{pmatrix} 2.7 \\ 0 \end{pmatrix}, \end{aligned} \qquad (14.247)$$

$$\mathcal{M}_2 = \mathcal{M}_1 - 0.3\,\mathbf{m}_2 = \begin{pmatrix} 4.4 \\ 0 \end{pmatrix} - 0.3 \begin{pmatrix} 2.7 \\ 0 \end{pmatrix} = \begin{pmatrix} 4.4 - 0.81 \\ 0 \end{pmatrix} = \begin{pmatrix} 3.59 \\ 0 \end{pmatrix}, \quad (14.248)$$

$$\begin{aligned} \mathbf{m}_3 &= 0.5\,(\mathbf{m}_2 + (\nabla f)(\mathcal{M}_2)) = 0.5\left(\begin{pmatrix} 2.7 \\ 0 \end{pmatrix} + \begin{pmatrix} 3.59 - 1 \\ 10(0 - 1) \end{pmatrix}\right) \\ &= \begin{pmatrix} 0.5\,(2.7 + 2.59) \\ 0.5\,(0 - 10) \end{pmatrix} = \begin{pmatrix} 0.5 \cdot 5.29 \\ 0.5(-10) \end{pmatrix} \\ &= \begin{pmatrix} 2.5 + 0.145 \\ -5 \end{pmatrix} = \begin{pmatrix} 2.645 \\ -5 \end{pmatrix} \approx \begin{pmatrix} 2.65 \\ -5 \end{pmatrix}, \end{aligned} \qquad (14.249)$$

$$\begin{aligned} \mathcal{M}_3 &= \mathcal{M}_2 - 0.3\,\mathbf{m}_3 \approx \begin{pmatrix} 3.59 \\ 0 \end{pmatrix} - 0.3 \begin{pmatrix} 2.65 \\ -5 \end{pmatrix} \\ &= \begin{pmatrix} 3.59 - 0.795 \\ 1.5 \end{pmatrix} = \begin{pmatrix} 3 - 0.205 \\ 1.5 \end{pmatrix} = \begin{pmatrix} 2.795 \\ 1.5 \end{pmatrix} \approx \begin{pmatrix} 2.8 \\ 1.5 \end{pmatrix}, \end{aligned} \qquad (14.250)$$

$$\vdots$$

.

```python
# Example for GD and momentum GD

import numpy as np
import matplotlib.pyplot as plt

# Number of steps for the schemes
N = 8

# Problem setting
d = 2
K = [1., 10.]

vartheta = np.array([1., 1.])
xi = np.array([5., 3.])

def f(x, y):
    result = K[0] / 2. * np.abs(x - vartheta[0]) ** 2 \
    + K[1] / 2. * np.abs(y - vartheta[1]) ** 2
    return result

```

```python
21 def nabla_f(x):
22     return K * (x - vartheta)
23
24 # Coefficients for GD
25 gamma_GD = 2 /(K[0] + K[1])
26
27 # Coefficients for momentum
28 gamma_momentum = 0.3
29 alpha = 0.5
30
31 # Placeholder for processes
32 Theta = np.zeros((N+1, d))
33 M = np.zeros((N+1, d))
34 m = np.zeros((N+1, d))
35
36 Theta[0] = xi
37 M[0] = xi
38
39 # Perform gradient descent
40 for i in range(N):
41     Theta[i+1] = Theta[i] - gamma_GD * nabla_f(Theta[i])
42
43 # Perform momentum GD
44 for i in range(N):
45     m[i+1] = alpha * m[i] + (1 - alpha) * nabla_f(M[i])
46     M[i+1] = M[i] - gamma_momentum * m[i+1]
47
48
49 ### Plot ###
50 plt.figure()
51
52 # Plot the gradient descent process
53 plt.plot(Theta[:, 0], Theta[:, 1],
54          label = "GD", color = "c",
55          linestyle = "-", marker = "*")
56
57 # Plot the momentum gradient descent process
58 plt.plot(M[:, 0], M[:, 1],
59          label = "Momentum", color = "orange", marker = "*")
60
61 # Target value
62 plt.scatter(vartheta[0],vartheta[1],
63             label = "vartheta", color = "red", marker = "x")
64
65 # Plot contour lines of f
66 x = np.linspace(-3., 7., 100)
67 y = np.linspace(-2., 4., 100)
68 X, Y = np.meshgrid(x, y)
69 Z = f(X, Y)
70 cp = plt.contour(X, Y, Z, colors="black",
71                  levels = [0.5,2,4,8,16],
72                  linestyles=":")
73
74 plt.legend()
75 plt.savefig("GD_momentum_plots.pdf")
76 plt.show()
```
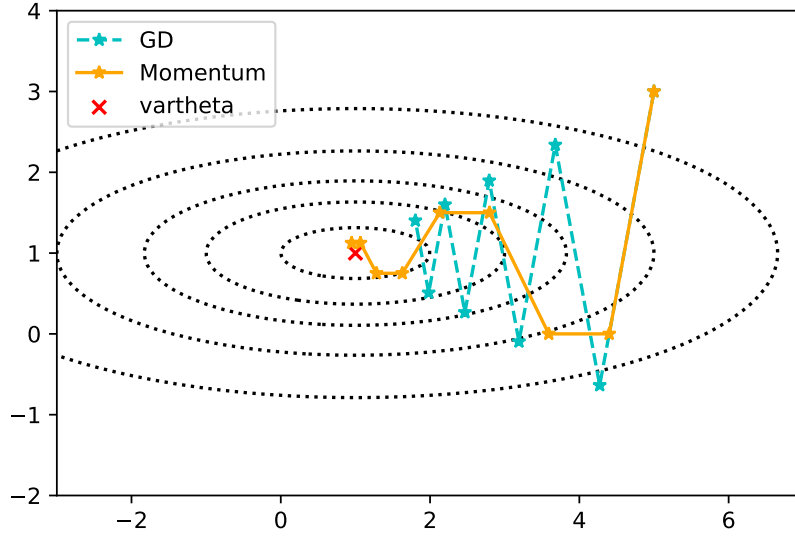
Source code 14.1: PYTHON code for Figure 14.1

Figure 14.1: Result of a call of PYTHON code 14.1

## 14.3 The gradient descent optimization method with Nesterov momentum

**Definition 14.3.1** (Nesterov accelerated gradient descent optimization method)**.** *Let $d \in \mathbb{N}$, $(\gamma_n)_{n\in\mathbb{N}} \subseteq [0,\infty)$, $(\alpha_n)_{n\in\mathbb{N}} \subseteq [0,1]$, $\xi \in \mathbb{R}^d$ and let $f\colon \mathbb{R}^d \to \mathbb{R}$ and $g\colon \mathbb{R}^d \to \mathbb{R}^d$ satisfy for all $\theta \in \{v \in \mathbb{R}^d\colon (f \text{ is differentiable at } v)\}$ that*

$$g(\theta) = (\nabla f)(\theta). \tag{14.251}$$

*Then we say that $\Theta$ is the Nesterov accelerated gradient descent process for the objective function $f$ with generalized gradient $g$, learning rates $(\gamma_n)_{n\in\mathbb{N}}$, momentum decay factors $(\alpha_n)_{n\in\mathbb{N}}$, and initial value $\xi$ (we say that $\Theta$ is the Nesterov accelerated gradient descent process for the objective function $f$ with learning rates $(\gamma_n)_{n\in\mathbb{N}}$, momentum decay factors $(\alpha_n)_{n\in\mathbb{N}}$, and initial value $\xi$) if and only if it holds that $\Theta\colon \mathbb{N}_0 \to \mathbb{R}^d$ is the function from $\mathbb{N}_0$ to $\mathbb{R}^d$ which satisfies that there exists a function $\mathbf{m}\colon \mathbb{N}_0 \to \mathbb{R}^d$ such that for all $n \in \mathbb{N}$ it holds that*

$$\Theta_0 = \xi, \qquad \mathbf{m}_0 = 0, \tag{14.252}$$

$$\mathbf{m}_n = \alpha_n \mathbf{m}_{n-1} + (1-\alpha_n)\, g(\Theta_{n-1} - \gamma_n \alpha_n \mathbf{m}_{n-1}), \tag{14.253}$$

$$\text{and} \qquad \Theta_n = \Theta_{n-1} - \gamma_n \mathbf{m}_n. \tag{14.254}$$

## 14.4 The adaptive gradient descent optimization method (Adagrad optimization method)

**Definition 14.4.1** (Adagrad optimization method)**.** *Let $d \in \mathbb{N}$, $(\gamma_n)_{n\in\mathbb{N}} \subseteq [0,\infty)$, $\varepsilon \in (0,\infty)$, $\xi \in \mathbb{R}^d$ and let $f\colon \mathbb{R}^d \to \mathbb{R}$ and $g = (g_1,\ldots,g_d)\colon \mathbb{R}^d \to \mathbb{R}^d$ satisfy for all $\theta \in \{v \in$*

$\mathbb{R}^d$: (*f is differentiable at v*)} *that*

$$g(\theta) = (\nabla f)(\theta). \tag{14.255}$$

*Then we say that* $\Theta$ *is the Adagrad gradient descent process for the objective function f with generalized gradient g, learning rates* $(\gamma_n)_{n\in\mathbb{N}}$, *regularizing factor* $\varepsilon$, *and initial value* $\xi$ *(we say that* $\Theta$ *is the Adagrad gradient descent process for the objective function f with learning rates* $(\gamma_n)_{n\in\mathbb{N}}$, *regularizing factor* $\varepsilon$, *and initial value* $\xi$*) if and only if it holds that* $\Theta = (\Theta^{(1)}, \ldots, \Theta^{(d)}): \mathbb{N}_0 \to \mathbb{R}^d$ *is the function from* $\mathbb{N}_0$ *to* $\mathbb{R}^d$ *which satisfies for all* $n \in \mathbb{N}$, $i \in \{1, 2, \ldots, d\}$ *that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n^{(i)} = \Theta_{n-1}^{(i)} - \gamma_n \left[ \varepsilon + \sum_{k=0}^{n-1} |g_i(\Theta_k)|^2 \right]^{-1/2} g_i(\Theta_{n-1}). \tag{14.256}$$

## 14.5 The root mean square propagation gradient descent optimization method (RMSprop gradient descent optimization method)

**Definition 14.5.1** (RMSprop gradient descent optimization method). *Let* $d \in \mathbb{N}$, $(\gamma_n)_{n\in\mathbb{N}} \subseteq [0, \infty)$, $(\beta_n)_{n\in\mathbb{N}} \subseteq [0, 1]$, $\varepsilon \in (0, \infty)$, $\xi \in \mathbb{R}^d$ *and let* $f: \mathbb{R}^d \to \mathbb{R}$ *and* $g = (g_1, \ldots, g_d): \mathbb{R}^d \to \mathbb{R}^d$ *satisfy for all* $\theta \in \{v \in \mathbb{R}^d: (f \text{ is differentiable at } v)\}$ *that*

$$g(\theta) = (\nabla f)(\theta). \tag{14.257}$$

*Then we say that* $\Theta$ *is the RMSprop gradient descent process for the objective function f with generalized gradient g, learning rates* $(\gamma_n)_{n\in\mathbb{N}}$, *second moment decay factors* $(\beta_n)_{n\in\mathbb{N}}$, *regularizing factor* $\varepsilon$, *and initial value* $\xi$ *(we say that* $\Theta$ *is the RMSprop gradient descent process for the objective function f with learning rates* $(\gamma_n)_{n\in\mathbb{N}}$, *second moment decay factors* $(\beta_n)_{n\in\mathbb{N}}$, *regularizing factor* $\varepsilon$, *and initial value* $\xi$*) if and only if it holds that* $\Theta = (\Theta^{(1)}, \ldots, \Theta^{(d)}): \mathbb{N}_0 \to \mathbb{R}^d$ *is the function from* $\mathbb{N}_0$ *to* $\mathbb{R}^d$ *which satisfies that there exists a function* $\mathbb{M} = (\mathbb{M}^{(1)}, \ldots, \mathbb{M}^{(d)}): \mathbb{N}_0 \to \mathbb{R}^d$ *such that for all* $n \in \mathbb{N}$, $i \in \{1, 2, \ldots, d\}$ *it holds that*

$$\Theta_0 = \xi, \qquad \mathbb{M}_0 = 0, \qquad \mathbb{M}_n^{(i)} = \beta_n \, \mathbb{M}_{n-1}^{(i)} + (1 - \beta_n)|g_i(\Theta_{n-1})|^2, \tag{14.258}$$

$$and \qquad \Theta_n^{(i)} = \Theta_{n-1}^{(i)} - \gamma_n \left[ \varepsilon + \mathbb{M}_n^{(i)} \right]^{-1/2} g_i(\Theta_{n-1}). \tag{14.259}$$

## 14.6 The Adadelta gradient descent optimization method

**Definition 14.6.1** (Adadelta gradient descent optimization method). *Let* $d \in \mathbb{N}$, $(\beta_n)_{n\in\mathbb{N}}$, $(\delta_n)_{n\in\mathbb{N}} \subseteq [0, 1]$, $\varepsilon \in (0, \infty)$, $\xi \in \mathbb{R}^d$ *and let* $f: \mathbb{R}^d \to \mathbb{R}$ *and* $g = (g_1, \ldots, g_d): \mathbb{R}^d \to \mathbb{R}^d$ *satisfy for all* $\theta \in \{v \in \mathbb{R}^d: (f \text{ is differentiable at } v)\}$ *that*

$$g(\theta) = (\nabla f)(\theta). \tag{14.260}$$

*Then we say that* $\Theta$ *is the Adadelta gradient descent process for the objective function f with generalized gradient g, second moment decay factors* $(\beta_n)_{n\in\mathbb{N}}$, *delta decay factors* $(\delta_n)_{n\in\mathbb{N}}$, *regularizing factor* $\varepsilon$, *and initial value* $\xi$ *(we say that* $\Theta$ *is the Adadelta gradient*

descent process for the objective function $f$ with second moment decay factors $(\beta_n)_{n\in\mathbb{N}}$, delta decay factors $(\delta_n)_{n\in\mathbb{N}}$, regularizing factor $\varepsilon$, and initial value $\xi$) if and only if it holds that $\Theta = (\Theta^{(1)}, \ldots, \Theta^{(d)})\colon \mathbb{N}_0 \to \mathbb{R}^d$ is the function from $\mathbb{N}_0$ to $\mathbb{R}^d$ which satisfies that there exist functions $\mathbb{M} = (\mathbb{M}^{(1)}, \ldots, \mathbb{M}^{(d)})$, $\Delta = (\Delta^{(1)}, \ldots, \Delta^{(d)})\colon \mathbb{N}_0 \to \mathbb{R}^d$ such that for all $n \in \mathbb{N}$, $i \in \{1, 2, \ldots, d\}$ it holds that

$$\Theta_0 = \xi, \qquad \mathbb{M}_0 = 0, \qquad \Delta_0 = 0, \tag{14.261}$$

$$\mathbb{M}_n^{(i)} = \beta_n\,\mathbb{M}_{n-1}^{(i)} + (1 - \beta_n)|g_i(\Theta_{n-1})|^2, \tag{14.262}$$

$$\Theta_n^{(i)} = \Theta_{n-1}^{(i)} - \left[\frac{\varepsilon + \Delta_{n-1}^{(i)}}{\varepsilon + \mathbb{M}_n^{(i)}}\right]^{1/2} g_i(\Theta_{n-1}), \tag{14.263}$$

$$\text{and} \qquad \Delta_n^{(i)} = \delta_n\,\Delta_{n-1}^{(i)} + (1 - \delta_n)\,|\Theta_n^{(i)} - \Theta_{n-1}^{(i)}|^2. \tag{14.264}$$

## 14.7 The adaptive moment estimation gradient descent optimization method (Adam gradient descent optimization method)

**Definition 14.7.1** (Adam gradient descent optimization method). *Let $d \in \mathbb{N}$, $(\gamma_n)_{n\in\mathbb{N}} \subseteq [0, \infty)$, $(\alpha_n)_{n\in\mathbb{N}}$, $(\beta_n)_{n\in\mathbb{N}} \subseteq [0, 1)$, $\xi \in \mathbb{R}^d$ and let $f\colon \mathbb{R}^d \to \mathbb{R}$ and $g = (g_1, \ldots, g_d)\colon \mathbb{R}^d \to \mathbb{R}^d$ satisfy for all $\theta \in \{v \in \mathbb{R}^d\colon (f \text{ is differentiable at } v)\}$ that*

$$g(\theta) = (\nabla f)(\theta). \tag{14.265}$$

*Then we say that $\Theta$ is the Adam gradient descent process for the objective function $f$ with generalized gradient $g$, learning rates $(\gamma_n)_{n\in\mathbb{N}}$, momentum decay factors $(\alpha_n)_{n\in\mathbb{N}}$, second moment decay factors $(\beta_n)_{n\in\mathbb{N}}$, and initial value $\xi$ (we say that $\Theta$ is the Adam gradient descent process for the objective function $f$ with learning rates $(\gamma_n)_{n\in\mathbb{N}}$, momentum decay factors $(\alpha_n)_{n\in\mathbb{N}}$, second moment decay factors $(\beta_n)_{n\in\mathbb{N}}$, and initial value $\xi$) if and only if it holds that $\Theta = (\Theta^{(1)}, \ldots, \Theta^{(d)})\colon \mathbb{N}_0 \to \mathbb{R}^d$ is the function from $\mathbb{N}_0$ to $\mathbb{R}^d$ which satisfies that there exist functions $\mathbf{m} = (\mathbf{m}^{(1)}, \ldots, \mathbf{m}^{(d)})$, $\mathbb{M} = (\mathbb{M}^{(1)}, \ldots, \mathbb{M}^{(d)})\colon \mathbb{N}_0 \to \mathbb{R}^d$ such that for all $n \in \mathbb{N}$, $i \in \{1, 2, \ldots, d\}$ it holds that*

$$\Theta_0 = \xi, \qquad \mathbf{m}_0 = 0, \qquad \mathbb{M}_0 = 0, \tag{14.266}$$

$$\mathbf{m}_n = \alpha_n\,\mathbf{m}_{n-1} + (1 - \alpha_n)\,g(\Theta_{n-1}), \tag{14.267}$$

$$\mathbb{M}_n^{(i)} = \beta_n\,\mathbb{M}_{n-1}^{(i)} + (1 - \beta_n)|g_i(\Theta_{n-1})|^2, \tag{14.268}$$

$$\text{and} \qquad \Theta_n^{(i)} = \Theta_{n-1}^{(i)} - \gamma_n \left[\varepsilon + \left[\frac{\mathbb{M}_n^{(i)}}{(1 - \prod_{l=1}^n \beta_l)}\right]^{1/2}\right]^{-1} \left[\frac{\mathbf{m}_n^{(i)}}{(1 - \prod_{l=1}^n \alpha_l)}\right]. \tag{14.269}$$

# Chapter 15

# Optimization through gradient descent processes

## 15.1 The deterministic gradient descent optimization method

**Definition 15.1.1** (Gradient descent optimization method). *Let $d \in \mathbb{N}$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \infty)$, $\xi \in \mathbb{R}^d$ and let $f \colon \mathbb{R}^d \to \mathbb{R}$ and $g \colon \mathbb{R}^d \to \mathbb{R}^d$ satisfy for all $\theta \in \{v \in \mathbb{R}^d \colon (f$ is differentiable at $v)\}$ that*

$$g(\theta) = (\nabla f)(\theta). \tag{15.1}$$

*Then we say that $\Theta$ is the gradient descent process for the objective function $f$ with generalized gradient $g$, learning rates $(\gamma_n)_{n \in \mathbb{N}}$, and initial value $\xi$ (we say that $\Theta$ is the gradient descent process for the objective function $f$ with learning rates $(\gamma_n)_{n \in \mathbb{N}}$ and initial value $\xi$) if and only if it holds that $\Theta \colon \mathbb{N}_0 \to \mathbb{R}^d$ is the function from $\mathbb{N}_0$ to $\mathbb{R}^d$ which satisfies for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \gamma_n g(\Theta_{n-1}). \tag{15.2}$$

## 15.2 The stochastic gradient descent optimization method

**Definition 15.2.1** (Stochastic gradient descent optimization method). *Let $d \in \mathbb{N}$, $(\gamma_n)_{n \in \mathbb{N}} \subseteq [0, \infty)$, $(J_n)_{n \in \mathbb{N}} \subseteq \mathbb{N}$, let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space, let $(S, \mathcal{S})$ be a measurable space, let $\xi \colon \Omega \to \mathbb{R}^d$ and $X_{n,j} \colon \Omega \to S$, $j \in \{1, 2, \ldots, J_n\}$, $n \in \mathbb{N}$, be random variables, and let $F = (F(\theta, x))_{(\theta, x) \in \mathbb{R}^d \times S} \colon \mathbb{R}^d \times S \to \mathbb{R}$ and $G \colon \mathbb{R}^d \times S \to \mathbb{R}^d$ satisfy for all $x \in S$, $\theta \in \{v \in \mathbb{R}^d \colon F(\cdot, x)$ is differentiable at $v\}$ that*

$$G(\theta, x) = (\nabla_\theta F)(\theta, x). \tag{15.3}$$

*Then we say that $\Theta$ is the stochastic gradient descent process on $((\Omega, \mathcal{F}, \mathbb{P}), (S, \mathcal{S}))$ for the loss function $F$ with generalized gradient $G$, learning rates $(\gamma_n)_{n \in \mathbb{N}}$, batch sizes $(J_n)_{n \in \mathbb{N}}$, initial value $\xi$, and data $(X_{n,j})_{j \in \{1,2,\ldots,J_n\}, n \in \mathbb{N}}$ (we say that $\Theta$ is the stochastic gradient descent process for the loss function $F$ with learning rates $(\gamma_n)_{n \in \mathbb{N}}$, batch sizes $(J_n)_{n \in \mathbb{N}}$, initial value $\xi$, and data $(X_{n,j})_{j \in \{1,2,\ldots,J_n\}, n \in \mathbb{N}}$) if and only if it holds that $\Theta \colon \mathbb{N}_0 \times \Omega \to \mathbb{R}^d$ is the function from $\mathbb{N}_0 \times \Omega$ to $\mathbb{R}^d$ which satisfies for all $n \in \mathbb{N}$ that*

$$\Theta_0 = \xi \qquad and \qquad \Theta_n = \Theta_{n-1} - \gamma_n \left[ \frac{1}{J_n} \sum_{j=1}^{J_n} G(\Theta_{n-1}, X_{n,j}) \right]. \tag{15.4}$$

# Chapter 16

# Additional material

## 16.1 Compositions of ANNs and affine linear transformations

**Corollary 16.1.1.** *Let $\Phi \in \mathbf{N}$ (cf. Definition 2.2.1). Then*

(i) *it holds for all $\mathbb{A} \in \mathbf{N}$ with $\mathcal{L}(\mathbb{A}) = 1$ and $\mathcal{I}(\mathbb{A}) = \mathcal{O}(\Phi)$ that*

$$\mathcal{P}(\mathbb{A} \bullet \Phi) \leq \left[\max\left\{1, \tfrac{\mathcal{O}(\mathbb{A})}{\mathcal{O}(\Phi)}\right\}\right]\mathcal{P}(\Phi) \tag{16.1}$$

*and*

(ii) *it holds for all $\mathbb{A} \in \mathbf{N}$ with $\mathcal{L}(\mathbb{A}) = 1$ and $\mathcal{I}(\Phi) = \mathcal{O}(\mathbb{A})$ that*

$$\mathcal{P}(\Phi \bullet \mathbb{A}) \leq \left[\max\left\{1, \tfrac{\mathcal{I}(\mathbb{A})+1}{\mathcal{I}(\Phi)+1}\right\}\right]\mathcal{P}(\Phi) \tag{16.2}$$

*(cf. Definition 2.2.5).*

*Proof of Corollary 16.1.1.* Throughout this proof let $L \in \mathbb{N}$, $l_0, l_1, \ldots, l_L \in \mathbb{N}$, $\mathbb{A}_1, \mathbb{A}_2 \in \mathbf{N}$ satisfy $\mathcal{L}(\mathbb{A}_1) = \mathcal{L}(\mathbb{A}_2) = 1$, $\mathcal{I}(\mathbb{A}_1) = \mathcal{O}(\Phi)$, $\mathcal{I}(\Phi) = \mathcal{O}(\mathbb{A}_2)$, and $\mathcal{D}(\Phi) = (l_0, l_1, \ldots, l_L)$. Observe that item (iv) in Proposition 2.2.7, the fact that $\mathcal{O}(\Phi) = l_L$, the fact that $\mathcal{I}(\Phi) = l_0$, and the fact that for all $k \in \{1, 2\}$ it holds that $\mathcal{D}(\mathbb{A}_k) = (\mathcal{I}(\mathbb{A}_k), \mathcal{O}(\mathbb{A}_k))$ ensure that

$$
\begin{aligned}
\mathcal{P}(\mathbb{A}_1 \bullet \Phi) &= \left[\sum_{m=1}^{L-1} l_m(l_{m-1}+1)\right] + \left[\mathcal{O}(\mathbb{A}_1)\right](l_{L-1}+1) \\
&= \left[\sum_{m=1}^{L-1} l_m(l_{m-1}+1)\right] + \left[\tfrac{\mathcal{O}(\mathbb{A}_1)}{l_L}\right]l_L(l_{L-1}+1) \\
&\leq \left[\max\left\{1, \tfrac{\mathcal{O}(\mathbb{A}_1)}{l_L}\right\}\right]\left[\sum_{m=1}^{L-1} l_m(l_{m-1}+1)\right] + \left[\max\left\{1, \tfrac{\mathcal{O}(\mathbb{A}_1)}{l_L}\right\}\right]l_L(l_{L-1}+1) \\
&= \left[\max\left\{1, \tfrac{\mathcal{O}(\mathbb{A}_1)}{l_L}\right\}\right]\left[\sum_{m=1}^{L} l_m(l_{m-1}+1)\right] = \left[\max\left\{1, \tfrac{\mathcal{O}(\mathbb{A}_1)}{\mathcal{O}(\Phi)}\right\}\right]\mathcal{P}(\Phi)
\end{aligned}
\tag{16.3}
$$

and

$$
\begin{aligned}
\mathcal{P}(\Phi \bullet \mathbb{A}_2) &= \left[ \sum_{m=2}^{L} l_m(l_{m-1}+1) \right] + l_1 \big[ \mathcal{I}(\mathbb{A}_2) + 1 \big] \\
&= \left[ \sum_{m=2}^{L} l_m(l_{m-1}+1) \right] + \left[ \tfrac{\mathcal{I}(\mathbb{A}_2)+1}{l_0+1} \right] l_1(l_0+1) \\
&\leq \left[ \max\!\left\{ 1, \tfrac{\mathcal{I}(\mathbb{A}_2)+1}{l_0+1} \right\} \right] \left[ \sum_{m=2}^{L} l_m(l_{m-1}+1) \right] + \left[ \max\!\left\{ 1, \tfrac{\mathcal{I}(\mathbb{A}_2)+1}{l_0+1} \right\} \right] l_1(l_0+1) \\
&= \left[ \max\!\left\{ 1, \tfrac{\mathcal{I}(\mathbb{A}_2)+1}{l_0+1} \right\} \right] \left[ \sum_{m=1}^{L} l_m(l_{m-1}+1) \right] = \left[ \max\!\left\{ 1, \tfrac{\mathcal{I}(\mathbb{A}_2)+1}{\mathcal{I}(\Phi)+1} \right\} \right] \mathcal{P}(\Phi).
\end{aligned}
\tag{16.4}
$$

This establishes items (i)–(ii). The proof of Corollary 16.1.1 is thus complete.  □

## 16.2   Powers and extensions of ANNs

**Definition 16.2.1** (Extension of ANNs). *Let $L \in \mathbb{N}$, $\Psi \in \mathbf{N}$ satisfy $\mathcal{I}(\Psi) = \mathcal{O}(\Psi)$. Then we denote by $\mathcal{E}_{L,\Psi} \colon \{\Phi \in \mathbf{N} \colon (\mathcal{L}(\Phi) \leq L \text{ and } \mathcal{O}(\Phi) = \mathcal{I}(\Psi))\} \to \mathbf{N}$ the function which satisfies for all $\Phi \in \mathbf{N}$ with $\mathcal{L}(\Phi) \leq L$ and $\mathcal{O}(\Phi) = \mathcal{I}(\Psi)$ that*

$$
\mathcal{E}_{L,\Psi}(\Phi) = (\Psi^{\bullet(L-\mathcal{L}(\Phi))}) \bullet \Phi
\tag{16.5}
$$

*(cf. Definitions 2.2.1, 2.2.5, and 2.2.10).*

**Lemma 16.2.2.** *Let $d, \mathfrak{i} \in \mathbb{N}$, $\Psi \in \mathbf{N}$ satisfy that $\mathcal{D}(\Psi) = (d, \mathfrak{i}, d)$ (cf. Definition 2.2.1). Then*

*(i)  it holds for all $n \in \mathbb{N}_0$ that $\mathcal{L}(\Psi^{\bullet n}) = n+1$, $\mathcal{D}(\Psi^{\bullet n}) \in \mathbb{N}^{n+2}$, and*

$$
\mathcal{D}(\Psi^{\bullet n}) = \begin{cases} (d, d) & : n = 0 \\ (d, \mathfrak{i}, \mathfrak{i}, \ldots, \mathfrak{i}, d) & : n \in \mathbb{N} \end{cases}
\tag{16.6}
$$

*and*

*(ii)  it holds for all $\Phi \in \mathbf{N}$, $L \in \mathbb{N} \cap [\mathcal{L}(\Phi), \infty)$ with $\mathcal{O}(\Phi) = d$ that $\mathcal{L}\big(\mathcal{E}_{L,\Psi}(\Phi)\big) = L$*

*(cf. Definitions 2.2.10 and 16.2.1).*

*Proof of Lemma 16.2.2.* Throughout this proof let $\Phi \in \mathbf{N}$ satisfy $\mathcal{O}(\Phi) = d$. We claim that for all $n \in \mathbb{N}_0$ it holds that

$$
\mathcal{L}(\Psi^{\bullet n}) = n+1 \qquad \text{and} \qquad \mathbb{N}^{n+2} \ni \mathcal{D}(\Psi^{\bullet n}) = \begin{cases} (d, d) & : n = 0 \\ (d, \mathfrak{i}, \mathfrak{i}, \ldots, \mathfrak{i}, d) & : n \in \mathbb{N}. \end{cases}
\tag{16.7}
$$

We now prove (16.7) by induction on $n \in \mathbb{N}_0$. Note that the fact that $\Psi^{\bullet 0} = (\mathrm{I}_d, 0) \in \mathbb{R}^{d \times d} \times \mathbb{R}^d$ (cf. Definition 2.2.9) establishes (16.7) in the base case $n = 0$. For the induction step assume that there exists $n \in \mathbb{N}_0$ such that

$$
\mathcal{L}(\Psi^{\bullet n}) = n+1 \qquad \text{and} \qquad \mathbb{N}^{n+2} \ni \mathcal{D}(\Psi^{\bullet n}) = \begin{cases} (d, d) & : n = 0 \\ (d, \mathfrak{i}, \mathfrak{i}, \ldots, \mathfrak{i}, d) & : n \in \mathbb{N}. \end{cases}
\tag{16.8}
$$

Observe that Lemma 2.2.4, (2.109), items (i)–(ii) in Proposition 2.2.7, (16.8), and the assumption that $\mathcal{D}(\Psi) = (d, \mathfrak{i}, d)$ imply that

$$
\begin{aligned}
\mathcal{L}(\Psi^{\bullet(n+1)}) &= \mathcal{L}(\Psi \bullet (\Psi^{\bullet n})) = \mathcal{L}(\Psi) + \mathcal{L}(\Psi^{\bullet n}) - 1 = 2 + (n+1) - 1 = (n+1) + 1 \\
&\text{and} \qquad \mathcal{D}(\Psi^{\bullet(n+1)}) = \mathcal{D}(\Psi \bullet (\Psi^{\bullet n})) = (d, \mathfrak{i}, \mathfrak{i}, \dots, \mathfrak{i}, d) \in \mathbb{N}^{n+3}.
\end{aligned}
$$
$$(16.9)$$

Induction thus proves (16.7). Next note that (16.7) establishes item (i). In addition, observe that item (ii) in Proposition 2.2.7, item (i), and (16.5) ensure that for all $L \in \mathbb{N} \cap [\mathcal{L}(\Phi), \infty)$ it holds that

$$
\begin{aligned}
\mathcal{L}\big(\mathcal{E}_{L,\Psi}(\Phi)\big) &= \mathcal{L}\big((\Psi^{\bullet(L-\mathcal{L}(\Phi))}) \bullet \Phi\big) = \mathcal{L}\big(\Psi^{\bullet(L-\mathcal{L}(\Phi))}\big) + \mathcal{L}(\Phi) - 1 \\
&= (L - \mathcal{L}(\Phi) + 1) + \mathcal{L}(\Phi) - 1 = L.
\end{aligned}
$$
$$(16.10)$$

This establishes item (ii). The proof of Lemma 16.2.2 is thus complete. $\qquad \square$

**Lemma 16.2.3.** *Let $a \in C(\mathbb{R}, \mathbb{R})$, $\mathbb{I} \in \mathbf{N}$ satisfy for all $x \in \mathbb{R}^{\mathcal{I}(\mathbb{I})}$ that $\mathcal{I}(\mathbb{I}) = \mathcal{O}(\mathbb{I})$ and $(\mathcal{R}_a(\mathbb{I}))(x) = x$ (cf. Definitions 2.2.1 and 2.2.3). Then*

*(i) it holds for all $n \in \mathbb{N}_0$, $x \in \mathbb{R}^{\mathcal{I}(\mathbb{I})}$ that*

$$
\mathcal{R}_a(\mathbb{I}^{\bullet n}) \in C(\mathbb{R}^{\mathcal{I}(\mathbb{I})}, \mathbb{R}^{\mathcal{I}(\mathbb{I})}) \qquad and \qquad (\mathcal{R}_a(\mathbb{I}^{\bullet n}))(x) = x \qquad (16.11)
$$

*and*

*(ii) it holds for all $\Phi \in \mathbf{N}$, $L \in \mathbb{N} \cap [\mathcal{L}(\Phi), \infty)$, $x \in \mathbb{R}^{\mathcal{I}(\Phi)}$ with $\mathcal{O}(\Phi) = \mathcal{I}(\mathbb{I})$ that*

$$
\mathcal{R}_a(\mathcal{E}_{L,\mathbb{I}}(\Phi)) \in C(\mathbb{R}^{\mathcal{I}(\Phi)}, \mathbb{R}^{\mathcal{O}(\Phi)}) \quad and \quad \big(\mathcal{R}_a(\mathcal{E}_{L,\mathbb{I}}(\Phi))\big)(x) = \big(\mathcal{R}_a(\Phi)\big)(x) \quad (16.12)
$$

*(cf. Definitions 2.2.10 and 16.2.1).*

*Proof of Lemma 16.2.3.* Throughout this proof let $\Phi \in \mathbf{N}$, $L, d \in \mathbb{N}$ satisfy $\mathcal{L}(\Phi) \le L$ and $\mathcal{I}(\mathbb{I}) = \mathcal{O}(\Phi) = d$. We claim that for all $n \in \mathbb{N}_0$ it holds that

$$
\mathcal{R}_a(\mathbb{I}^{\bullet n}) \in C(\mathbb{R}^d, \mathbb{R}^d) \qquad and \qquad \forall\, x \in \mathbb{R}^d \colon (\mathcal{R}_a(\mathbb{I}^{\bullet n}))(x) = x. \qquad (16.13)
$$

We now prove (16.13) by induction on $n \in \mathbb{N}_0$. Note that (2.109) and the fact that $\mathcal{O}(\mathbb{I}) = d$ demonstrate that $\mathcal{R}_a(\mathbb{I}^{\bullet 0}) \in C(\mathbb{R}^d, \mathbb{R}^d)$ and $\forall\, x \in \mathbb{R}^d \colon (\mathcal{R}_a(\mathbb{I}^{\bullet 0}))(x) = x$. This establishes (16.13) in the base case $n = 0$. For the induction step observe that for all $n \in \mathbb{N}_0$ with $\mathcal{R}_a(\mathbb{I}^{\bullet n}) \in C(\mathbb{R}^d, \mathbb{R}^d)$ and $\forall\, x \in \mathbb{R}^d \colon (\mathcal{R}_a(\mathbb{I}^{\bullet n}))(x) = x$ it holds that

$$
\mathcal{R}_a(\mathbb{I}^{\bullet(n+1)}) = \mathcal{R}_a(\mathbb{I} \bullet (\mathbb{I}^{\bullet n})) = (\mathcal{R}_a(\mathbb{I})) \circ (\mathcal{R}_a(\mathbb{I}^{\bullet n})) \in C(\mathbb{R}^d, \mathbb{R}^d) \qquad (16.14)
$$

and

$$
\begin{aligned}
\forall\, x \in \mathbb{R}^d \colon \big(\mathcal{R}_a(\mathbb{I}^{\bullet(n+1)})\big)(x) &= \big([\mathcal{R}_a(\mathbb{I})] \circ [\mathcal{R}_a(\mathbb{I}^{\bullet n})]\big)(x) \\
&= (\mathcal{R}_a(\mathbb{I}))\big((\mathcal{R}_a(\mathbb{I}^{\bullet n}))(x)\big) = (\mathcal{R}_a(\mathbb{I}))(x) = x.
\end{aligned}
$$
$$(16.15)$$

Induction thus proves (16.13). Next observe that (16.13) establishes item (i). Moreover, note that (16.5), item (v) in Proposition 2.2.7, item (i), and the fact that $\mathcal{I}(\mathbb{I}) = \mathcal{O}(\Phi)$ ensure that

$$
\begin{aligned}
\mathcal{R}_a(\mathcal{E}_{L,\mathbb{I}}(\Phi)) &= \mathcal{R}_a((\mathbb{I}^{\bullet(L-\mathcal{L}(\Phi))}) \bullet \Phi) \\
&\in C(\mathbb{R}^{\mathcal{I}(\Phi)}, \mathbb{R}^{\mathcal{O}(\mathbb{I})}) = C(\mathbb{R}^{\mathcal{I}(\Phi)}, \mathbb{R}^{\mathcal{I}(\mathbb{I})}) = C(\mathbb{R}^{\mathcal{I}(\Phi)}, \mathbb{R}^{\mathcal{O}(\Phi)})
\end{aligned}
$$
$$(16.16)$$

and

$$\forall\, x \in \mathbb{R}^{\mathcal{I}(\Phi)} \colon \; \big(\mathcal{R}_a(\mathcal{E}_{L,\mathbb{I}}(\Phi))\big)(x) = \big(\mathcal{R}_a(\mathbb{I}^{\bullet(L-\mathcal{L}(\Phi))})\big)\big((\mathcal{R}_a(\Phi))(x)\big)$$
$$= (\mathcal{R}_a(\Phi))(x). \tag{16.17}$$

This establishes item (ii). The proof of Lemma 16.2.3 is thus complete. □

**Lemma 16.2.4.** *Let* $d, \mathfrak{i}, L, \mathfrak{L} \in \mathbb{N}$, $l_0, l_1, \ldots, l_{L-1} \in \mathbb{N}$, $\Phi, \Psi \in \mathbf{N}$ *satisfy* $\mathfrak{L} \geq L$, $\mathcal{D}(\Phi) = (l_0, l_1, \ldots, l_{L-1}, d)$ *and* $\mathcal{D}(\Psi) = (d, \mathfrak{i}, d)$ *(cf. Definition 2.2.1). Then it holds that* $\mathcal{D}(\mathcal{E}_{\mathfrak{L},\Psi}(\Phi)) \in \mathbb{N}^{\mathfrak{L}+1}$ *and*

$$\mathcal{D}(\mathcal{E}_{\mathfrak{L},\Psi}(\Phi)) = \begin{cases} (l_0, l_1, \ldots, l_{L-1}, d) & : \mathfrak{L} = L \\ (l_0, l_1, \ldots, l_{L-1}, \mathfrak{i}, \mathfrak{i}, \ldots, \mathfrak{i}, d) & : \mathfrak{L} > L \end{cases} \tag{16.18}$$

*(cf. Definition 16.2.1).*

*Proof of Lemma 16.2.4.* Observe that item (i) in Lemma 16.2.2 ensures that $\mathcal{L}(\Psi^{\bullet(\mathfrak{L}-L)}) = \mathfrak{L} - L + 1$, $\mathcal{D}(\Psi^{\bullet(\mathfrak{L}-L)}) \in \mathbb{N}^{\mathfrak{L}-L+2}$, and

$$\mathcal{D}(\Psi^{\bullet(\mathfrak{L}-L)}) = \begin{cases} (d, d) & : \mathfrak{L} = L \\ (d, \mathfrak{i}, \mathfrak{i}, \ldots, \mathfrak{i}, d) & : \mathfrak{L} > L \end{cases} \tag{16.19}$$

*(cf. Definition 2.2.10).* Combining this with Proposition 2.2.7 shows that $\mathcal{L}((\Psi^{\bullet(\mathfrak{L}-L)}) \bullet \Phi) = \mathcal{L}(\Psi^{\bullet(\mathfrak{L}-L)}) + \mathcal{L}(\Phi) - 1 = \mathfrak{L}$, $\mathcal{D}((\Psi^{\bullet(\mathfrak{L}-L)}) \bullet \Phi) \in \mathbb{N}^{\mathfrak{L}+1}$, and

$$\mathcal{D}((\Psi^{\bullet(\mathfrak{L}-L)}) \bullet \Phi) = \begin{cases} (l_0, l_1, \ldots, l_{L-1}, d) & : \mathfrak{L} = L \\ (l_0, l_1, \ldots, l_{L-1}, \mathfrak{i}, \mathfrak{i}, \ldots, \mathfrak{i}, d) & : \mathfrak{L} > L. \end{cases} \tag{16.20}$$

This and (16.5) establish (16.18). The proof of Lemma 16.2.4 is thus complete. □

**Lemma 16.2.5.** *Let* $d, \mathfrak{i} \in \mathbb{N}$, $\Psi \in \mathbf{N}$ *satisfy that* $\mathcal{D}(\Psi) = (d, \mathfrak{i}, d)$ *(cf. Definition 2.2.1). Then*

(i) *it holds for all* $n \in \mathbb{N}_0$ *that* $\mathcal{L}(\Psi^{\bullet n}) = n + 1$, $\mathcal{D}(\Psi^{\bullet n}) \in \mathbb{N}^{n+2}$, *and*

$$\mathcal{D}(\Psi^{\bullet n}) = \begin{cases} (d, d) & : n = 0 \\ (d, \mathfrak{i}, \mathfrak{i}, \ldots, \mathfrak{i}, d) & : n \in \mathbb{N} \end{cases} \tag{16.21}$$

*and*

(ii) *it holds for all* $\Phi \in \mathbf{N}$, $L \in \mathbb{N} \cap [\mathcal{L}(\Phi), \infty)$ *with* $\mathcal{O}(\Phi) = d$ *that* $\mathcal{L}\big(\mathcal{E}_{L,\Psi}(\Phi)\big) = L$ *and*

$$\mathcal{P}(\mathcal{E}_{L,\Psi}(\Phi))$$
$$\leq \begin{cases} \mathcal{P}(\Phi) & : \mathcal{L}(\Phi) = L \\ \big[(\max\{1, \frac{\mathfrak{i}}{d}\})\mathcal{P}(\Phi) + \big((L - \mathcal{L}(\Phi) - 1)\,\mathfrak{i} + d\big)(\mathfrak{i}+1)\big] & : \mathcal{L}(\Phi) < L \end{cases} \tag{16.22}$$

*(cf. Definitions 2.2.10, 16.2.1, and 16.2.1).*

Chapter 16. Additional material

*Proof of Lemma 16.2.2.* Throughout this proof let $\Phi \in \mathbf{N}$, $l_0, l_1, \ldots, l_{\mathcal{L}(\Phi)} \in \mathbb{N}$ satisfy $\mathcal{O}(\Phi) = d$ and $\mathcal{D}(\Phi) = (l_0, l_1, \ldots, l_{\mathcal{L}(\Phi)}) \in \mathbb{N}^{\mathcal{L}(\Phi)+1}$ and let $a_{L,k} \in \mathbb{N}$, $k \in \mathbb{N}_0 \cap [0, L]$, $L \in \mathbb{N} \cap [\mathcal{L}(\Phi), \infty)$, satisfy for all $L \in \mathbb{N} \cap [\mathcal{L}(\Phi), \infty)$, $k \in \mathbb{N}_0 \cap [0, L]$ that

$$a_{L,k} = \begin{cases} l_k & : k < \mathcal{L}(\Phi) \\ \mathfrak{i} & : \mathcal{L}(\Phi) \leq k < L \\ d & : k = L \end{cases}. \tag{16.23}$$

We claim that for all $n \in \mathbb{N}_0$ it holds that

$$\mathcal{L}(\Psi^{\bullet n}) = n + 1 \qquad \text{and} \qquad \mathbb{N}^{n+2} \ni \mathcal{D}(\Psi^{\bullet n}) = \begin{cases} (d, d) & : n = 0 \\ (d, \mathfrak{i}, \mathfrak{i}, \ldots, \mathfrak{i}, d) & : n \in \mathbb{N} \end{cases}. \tag{16.24}$$

We now prove (16.7) by induction on $n \in \mathbb{N}_0$. Note that the fact that $\Psi^{\bullet 0} = (\mathrm{I}_d, 0) \in \mathbb{R}^{d \times d} \times \mathbb{R}^d$ (cf. Definition 2.2.9) establishes (16.6) in the base case $n = 0$. For the induction step assume that there exists $n \in \mathbb{N}_0$ such that

$$\mathcal{L}(\Psi^{\bullet n}) = n + 1 \qquad \text{and} \qquad \mathbb{N}^{n+2} \ni \mathcal{D}(\Psi^{\bullet n}) = \begin{cases} (d, d) & : n = 0 \\ (d, \mathfrak{i}, \mathfrak{i}, \ldots, \mathfrak{i}, d) & : n \in \mathbb{N} \end{cases}. \tag{16.25}$$

Observe that Lemma 2.2.4, (2.109), items (i)–(ii) in Proposition 2.2.7, (16.8), and the assumption that $\mathcal{D}(\Psi) = (d, \mathfrak{i}, d)$ imply that

$$\mathcal{L}(\Psi^{\bullet(n+1)}) = \mathcal{L}(\Psi \bullet (\Psi^{\bullet n})) = \mathcal{L}(\Psi) + \mathcal{L}(\Psi^{\bullet n}) - 1 = 2 + (n+1) - 1 = (n+1) + 1$$
$$\text{and} \qquad \mathcal{D}(\Psi^{\bullet(n+1)}) = \mathcal{D}(\Psi \bullet (\Psi^{\bullet n})) = (d, \mathfrak{i}, \mathfrak{i}, \ldots, \mathfrak{i}, d) \in \mathbb{N}^{n+3}. \tag{16.26}$$

Induction thus proves (16.7). Next note that (16.7) establishes item (i). In addition, observe that items (i)–(ii) in Proposition 2.2.7, item (i), (16.5), and (16.23) ensure that for all $L \in \mathbb{N} \cap [\mathcal{L}(\Phi), \infty)$ it holds that

$$\mathcal{L}(\mathcal{E}_{L,\Psi}(\Phi)) = \mathcal{L}((\Psi^{\bullet(L-\mathcal{L}(\Phi))}) \bullet \Phi) = \mathcal{L}(\Psi^{\bullet(L-\mathcal{L}(\Phi))}) + \mathcal{L}(\Phi) - 1$$
$$= (L - \mathcal{L}(\Phi) + 1) + \mathcal{L}(\Phi) - 1 = L \tag{16.27}$$

and

$$\mathcal{D}(\mathcal{E}_{L,\Psi}(\Phi)) = \mathcal{D}((\Psi^{\bullet(L-\mathcal{L}(\Phi))}) \bullet \Phi) = (a_{L,0}, a_{L,1}, \ldots, a_{L,L}). \tag{16.28}$$

Combining this with (16.23) demonstrates that

$$\mathcal{L}(\mathcal{E}_{\mathcal{L}(\Phi),\Psi}(\Phi)) = \mathcal{L}(\Phi) \tag{16.29}$$

and

$$\mathcal{D}(\mathcal{E}_{\mathcal{L}(\Phi),\Psi}(\Phi)) = (a_{\mathcal{L}(\Phi),0}, a_{\mathcal{L}(\Phi),1}, \ldots, a_{\mathcal{L}(\Phi),\mathcal{L}(\Phi)})$$
$$= (l_0, l_1, \ldots, l_{\mathcal{L}(\Phi)}) = \mathcal{D}(\Phi). \tag{16.30}$$

Hence, we obtain that

$$\mathcal{P}(\mathcal{E}_{\mathcal{L}(\Phi),\Psi}(\Phi)) = \mathcal{P}(\Phi). \tag{16.31}$$

Dissemination prohibited. July 29, 2021261

Next note that (16.23), (16.28), and the fact that $l_{\mathcal{L}(\Phi)} = \mathcal{O}(\Phi) = d$ imply that for all $L \in \mathbb{N} \cap (\mathcal{L}(\Phi), \infty)$ it holds that

$$
\begin{aligned}
\mathcal{P}\big(\mathcal{E}_{L,\Psi}(\Phi)\big) &= \sum_{k=1}^{L} a_{L,k}(a_{L,k-1} + 1) \\
&= \left[ \sum_{k=1}^{\mathcal{L}(\Phi)-1} a_{L,k}(a_{L,k-1} + 1) \right] + \left[ \sum_{k=\mathcal{L}(\Phi)}^{L} a_{L,k}(a_{L,k-1} + 1) \right] \\
&= \left[ \sum_{k=1}^{\mathcal{L}(\Phi)-1} l_k(l_{k-1} + 1) \right] + \left[ \sum_{k=\mathcal{L}(\Phi)}^{\mathcal{L}(\Phi)} a_{L,k}(a_{L,k-1} + 1) \right] \\
&\quad + \left[ \sum_{k=\mathcal{L}(\Phi)+1}^{L} a_{L,k}(a_{L,k-1} + 1) \right] \\
&= \left[ \sum_{k=1}^{\mathcal{L}(\Phi)-1} l_k(l_{k-1} + 1) \right] + a_{L,\mathcal{L}(\Phi)}(a_{L,\mathcal{L}(\Phi)-1} + 1) \\
&\quad + \left[ \sum_{k=\mathcal{L}(\Phi)+1}^{L-1} a_{L,k}(a_{L,k-1} + 1) \right] + \left[ \sum_{k=L}^{L} a_{L,k}(a_{L,k-1} + 1) \right] \\
&= \left[ \sum_{k=1}^{\mathcal{L}(\Phi)-1} l_k(l_{k-1} + 1) \right] + \mathfrak{i}(l_{\mathcal{L}(\Phi)-1} + 1) \\
&\quad + \big(L - 1 - (\mathcal{L}(\Phi) + 1) + 1\big)\mathfrak{i}(\mathfrak{i} + 1) + a_{L,L}(a_{L,L-1} + 1) \\
&= \left[ \sum_{k=1}^{\mathcal{L}(\Phi)-1} l_k(l_{k-1} + 1) \right] + \tfrac{\mathfrak{i}}{d}\big[l_{\mathcal{L}(\Phi)}(l_{\mathcal{L}(\Phi)-1} + 1)\big] \\
&\quad + \big(L - \mathcal{L}(\Phi) - 1\big)\mathfrak{i}(\mathfrak{i} + 1) + d(\mathfrak{i} + 1) \\
&\leq \big[\max\{1, \tfrac{\mathfrak{i}}{d}\}\big]\left[ \sum_{k=1}^{\mathcal{L}(\Phi)} l_k(l_{k-1} + 1) \right] + \big(L - \mathcal{L}(\Phi) - 1\big)\mathfrak{i}(\mathfrak{i} + 1) + d(\mathfrak{i} + 1) \\
&= \big[\max\{1, \tfrac{\mathfrak{i}}{d}\}\big]\mathcal{P}(\Phi) + \big(L - \mathcal{L}(\Phi) - 1\big)\mathfrak{i}(\mathfrak{i} + 1) + d(\mathfrak{i} + 1).
\end{aligned}
\tag{16.32}
$$

Combining this with (16.31) establishes (16.22). The proof of Lemma 16.2.2 is thus complete. □

## 16.3 Compositions of ANNs involving artificial identities

**Definition 16.3.1** (Composition of ANNs involving artificial identities)**.** *Let $\Psi \in \mathbf{N}$. Then we denote by*

$$
(\cdot) \odot_{\Psi} (\cdot) \colon \{(\Phi_1, \Phi_2) \in \mathbf{N} \times \mathbf{N} \colon \mathcal{I}(\Phi_1) = \mathcal{O}(\Psi) \text{ and } \mathcal{O}(\Phi_2) = \mathcal{I}(\Psi)\} \to \mathbf{N}
\tag{16.33}
$$

*the function which satisfies for all $\Phi_1, \Phi_2 \in \mathbf{N}$ with $\mathcal{I}(\Phi_1) = \mathcal{O}(\Psi)$ and $\mathcal{O}(\Phi_2) = \mathcal{I}(\Psi)$ that*

$$
\Phi_1 \odot_{\Psi} \Phi_2 = \Phi_1 \bullet (\Psi \bullet \Phi_2) = (\Phi_1 \bullet \Psi) \bullet \Phi_2
\tag{16.34}
$$

*(cf. Definitions 2.2.1 and 2.2.5 and Lemma 2.2.8).*

**Proposition 16.3.2.** *Let* $\Psi, \Phi_1, \Phi_2 \in \mathbf{N}$ *satisfy that* $\mathcal{H}(\Psi) = 1$, $\mathcal{I}(\Phi_1) = \mathcal{O}(\Psi)$, *and* $\mathcal{O}(\Phi_2) = \mathcal{I}(\Psi)$ *(cf. Definition 2.2.1). Then*

(i) *it holds that*

$$\mathcal{D}(\Phi_1 \odot_\Psi \Phi_2) = (\mathbb{D}_0(\Phi_2), \mathbb{D}_1(\Phi_2), \dots, \mathbb{D}_{\mathcal{L}(\Phi_2)-1}(\Phi_2), \mathbb{D}_1(\Psi), \mathbb{D}_1(\Phi_1), \mathbb{D}_2(\Phi_1), \dots, \mathbb{D}_{\mathcal{L}(\Phi_1)}(\Phi_1)),$$
(16.35)

(ii) *it holds that*

$$\mathcal{L}(\Phi_1 \odot_\Psi \Phi_2) = \mathcal{L}(\Phi_1) + \mathcal{L}(\Phi_2),$$
(16.36)

(iii) *it holds that*

$$\mathcal{P}(\Phi_1 \odot_\Psi \Phi_2) \le \left[\max\left\{1, \tfrac{\mathbb{D}_1(\Psi)}{\mathcal{I}(\Psi)}, \tfrac{\mathbb{D}_1(\Psi)}{\mathcal{O}(\Psi)}\right\}\right] (\mathcal{P}(\Phi_1) + \mathcal{P}(\Phi_2)),$$
(16.37)

*and*

(iv) *it holds for all* $a \in C(\mathbb{R}, \mathbb{R})$ *that* $\mathcal{R}_a(\Phi_1 \odot_\Psi \Phi_2) \in C(\mathbb{R}^{\mathcal{I}(\Phi_2)}, \mathbb{R}^{\mathcal{O}(\Phi_1)})$ *and*

$$\mathcal{R}_a(\Phi_1 \odot_\Psi \Phi_2) = [\mathcal{R}_a(\Phi_1)] \circ [\mathcal{R}_a(\Psi)] \circ [\mathcal{R}_a(\Phi_2)]$$
(16.38)

*(cf. Definitions 2.2.3 and 16.3.1).*

*Proof of Propositions 16.3.2.* Throughout this proof let $a \in C(\mathbb{R}, \mathbb{R})$, $L_1, L_2, l_{1,0}, l_{1,1}, \dots,$
$l_{1,\mathcal{L}(\Phi_1)}$,
$l_{2,0}, l_{2,1}, \dots, l_{2,\mathcal{L}(\Phi_2)}, \mathfrak{i} \in \mathbb{N}$ satisfy for all $k \in \{1, 2\}$ that $L_k = \mathcal{L}(\Phi_k)$, $\mathcal{D}(\Phi_k) = (l_{k,0}, l_{k,1}, \dots, l_{k,\mathcal{L}(\Phi_k)})$, and $\mathfrak{i} = \mathbb{D}_1(\Psi)$. Note that item (i) in Proposition 2.2.7, the fact that $\mathcal{D}(\Phi_2) = (l_{2,0}, l_{2,1}, \dots, l_{2,L_2})$, the fact that $\mathcal{L}(\Psi) = 2$, and the assumption that $\mathcal{I}(\Psi) = \mathcal{O}(\Phi_2)$ show that

$$\mathcal{D}(\Psi \bullet \Phi_2) = (l_{2,0}, l_{2,1}, \dots, l_{2,L_2-1}, \mathfrak{i}, \mathcal{O}(\Psi))$$
(16.39)

(cf. Definition 2.2.5). Combining this with item (i) in Proposition 2.2.7, the fact that $\mathcal{D}(\Phi_1) = (l_{1,0}, l_{1,1}, \dots, l_{1,L_1})$, and the assumption that $\mathcal{I}(\Phi_1) = \mathcal{O}(\Psi)$ proves that

$$\mathcal{D}(\Phi_1 \odot_\Psi \Phi_2) = \mathcal{D}\big(\Phi_1 \bullet (\Psi \bullet \Phi_2)\big) = (l_{2,0}, l_{2,1}, \dots, l_{2,L_2-1}, \mathfrak{i}, l_{1,1}, l_{1,2}, \dots, l_{1,L_1}). \quad (16.40)$$

This establishes item (i). Moreover, observe that item (ii) in Proposition 2.2.7 and the fact that $\mathcal{L}(\Psi) = 2$ ensure that

$$\begin{aligned}
\mathcal{L}(\Phi_1 \odot_\Psi \Phi_2) &= \mathcal{L}\big(\Phi_1 \bullet (\Psi \bullet \Phi_2)\big) = \mathcal{L}(\Phi_1) + \mathcal{L}(\Psi \bullet \Phi_2) - 1 \\
&= \mathcal{L}(\Phi_1) + \mathcal{L}(\Psi) + \mathcal{L}(\Phi_2) - 2 = \mathcal{L}(\Phi_1) + \mathcal{L}(\Phi_2).
\end{aligned}$$
(16.41)

This establishes item (ii). In addition, observe that (16.40), the fact that $\mathcal{I}(\Psi) = \mathcal{O}(\Phi_2) =$

$l_{2,L_2}$, and the fact that $\mathcal{O}(\Psi) = \mathcal{I}(\Phi_1) = l_{1,0}$ demonstrate that

$$
\begin{aligned}
\mathcal{P}(\Phi_1 \odot_\Psi \Phi_2) &= \left[\sum_{m=1}^{L_2-1} l_{2,m}(l_{2,m-1}+1)\right] + \left[\sum_{m=2}^{L_1} l_{1,m}(l_{1,m-1}+1)\right] \\
&\quad + \mathfrak{i}\big(l_{2,L_2-1}+1\big) + l_{1,1}(\mathfrak{i}+1) \\
&= \left[\sum_{m=1}^{L_2-1} l_{2,m}(l_{2,m-1}+1)\right] + \left[\sum_{m=2}^{L_1} l_{1,m}(l_{1,m-1}+1)\right] \\
&\quad + \tfrac{\mathfrak{i}}{\mathcal{I}(\Psi)}\, l_{2,L_2}\big(l_{2,L_2-1}+1\big) + l_{1,1}\big(\tfrac{\mathfrak{i}}{\mathcal{O}(\Psi)}\, l_{1,0}+1\big) \\
&\leq \left[\max\big\{1, \tfrac{\mathfrak{i}}{\mathcal{I}(\Psi)}\big\}\right]\left[\sum_{m=1}^{L_2} l_{2,m}(l_{2,m-1}+1)\right] \\
&\quad + \left[\max\big\{1, \tfrac{\mathfrak{i}}{\mathcal{O}(\Psi)}\big\}\right]\left[\sum_{m=1}^{L_1} l_{1,m}(l_{1,m-1}+1)\right] \\
&\leq \left[\max\big\{1, \tfrac{\mathfrak{i}}{\mathcal{I}(\Psi)}, \tfrac{\mathfrak{i}}{\mathcal{O}(\Psi)}\big\}\right]\big(\mathcal{P}(\Phi_1)+\mathcal{P}(\Phi_2)\big).
\end{aligned}
\tag{16.42}
$$

This establishes item (iii). Next note that item (v) in Proposition 2.2.7 implies that $\mathcal{R}_a(\Phi_1 \odot_\Psi \Phi_2) \in C(\mathbb{R}^{\mathcal{I}(\Phi_2)}, \mathbb{R}^{\mathcal{O}(\Phi_1)})$ and

$$
\begin{aligned}
\mathcal{R}_a(\Phi_1 \odot_\Psi \Phi_2) &= \mathcal{R}_a\big(\Phi_1 \bullet (\Psi \bullet \Phi_2)\big) \\
&= \big[\mathcal{R}_a(\Phi_1)\big] \circ \big[\mathcal{R}_a(\Psi \bullet \Phi_2)\big] \\
&= \big(\big[\mathcal{R}_a(\Phi_1)\big] \circ \big[\mathcal{R}_a(\Psi)\big] \circ \big[\mathcal{R}_a(\Phi_2)\big]\big) \in C(\mathbb{R}^{\mathcal{I}(\Phi_2)}, \mathbb{R}^{\mathcal{O}(\Phi_1)}).
\end{aligned}
\tag{16.43}
$$

This establishes item (iv). The proof of Proposition 16.3.2 is thus complete. $\qquad\square$

## 16.4   Parallelization of ANNs with different lengths

**Corollary 16.4.1.** *Let* $n, L \in \mathbb{N}$, $\mathfrak{i}_1, \mathfrak{i}_2, \ldots, \mathfrak{i}_n \in \mathbb{N}$, $\Psi = (\Psi_1, \Psi_2, \ldots, \Psi_n)$, $\Phi = (\Phi_1, \Phi_2, \ldots, \Phi_n) \in \mathbf{N}^n$ *satisfy for all* $j \in \{1, 2, \ldots, n\}$ *that* $\mathcal{D}(\Psi_j) = (\mathcal{O}(\Phi_j), \mathfrak{i}_j, \mathcal{O}(\Phi_j))$ *and* $L = \max_{k\in\{1,2,\ldots,n\}} \mathcal{L}(\Phi_k)$ *(cf. Definition 2.2.1). Then it holds that*

$$
\begin{aligned}
&\mathcal{P}\big(\mathrm{P}_{n,\Psi}(\Phi)\big) \\
&\leq \tfrac{1}{2}\Bigg(\left[\sum_{j=1}^n \big[\max\{1, \tfrac{\mathfrak{i}_j}{\mathcal{O}(\Phi_j)}\}\big]\, \mathcal{P}(\Phi_j)\, \mathbb{1}_{(\mathcal{L}(\Phi_j),\infty)}(L)\right] \\
&\quad + \left[\sum_{j=1}^n \big((L-\mathcal{L}(\Phi_j)-1)\,\mathfrak{i}_j\,(\mathfrak{i}_j+1) + \mathcal{O}(\Phi_j)\,(\mathfrak{i}_j+1)\big)\, \mathbb{1}_{(\mathcal{L}(\Phi_j),\infty)}(L)\right] \\
&\quad + \left[\sum_{j=1}^n \mathcal{P}(\Phi_j)\, \mathbb{1}_{\{\mathcal{L}(\Phi_j)\}}(L)\right]\Bigg)^2
\end{aligned}
\tag{16.44}
$$

*(cf. Definition 2.2.16).*

*Proof of Corollary 16.4.1.* Observe that (2.128), item (iii) in Proposition 2.2.14, and

item (ii) in Lemma 16.2.2 assure that

$$
\begin{aligned}
&\mathcal{P}\big(\mathrm{P}_{n,\Psi}(\Phi)\big)\\
&= \mathcal{P}\big(\mathbf{P}_n\big(\mathcal{E}_{L,\Psi_1}(\Phi_1), \mathcal{E}_{L,\Psi_2}(\Phi_2), \ldots, \mathcal{E}_{L,\Psi_n}(\Phi_n)\big)\big)\\
&\leq \tfrac{1}{2}\Big[\textstyle\sum_{j=1}^n \mathcal{P}(\mathcal{E}_{L,\Psi_j}(\Phi_j))\Big]^2\\
&\leq \tfrac{1}{2}\bigg(\Big[\textstyle\sum_{j=1}^n \big[\max\{1, \tfrac{\mathbf{i}_j}{\mathcal{O}(\Phi_j)}\}\big]\, \mathcal{P}(\Phi_j)\, \mathbb{1}_{(\mathcal{L}(\Phi_j),\infty)}(L)\Big]\\
&\qquad + \Big[\textstyle\sum_{j=1}^n \big((L - \mathcal{L}(\Phi_j) - 1)\,\mathbf{i}_j\,(\mathbf{i}_j + 1) + \mathcal{O}(\Phi_j)\,(\mathbf{i}_j + 1)\big)\, \mathbb{1}_{(\mathcal{L}(\Phi_j),\infty)}(L)\Big]\\
&\qquad + \Big[\textstyle\sum_{j=1}^n \mathcal{P}(\Phi_j)\, \mathbb{1}_{\{\mathcal{L}(\Phi_j)\}}(L)\Big]\bigg)^2
\end{aligned}
\tag{16.45}
$$

(cf. Definitions 2.2.11 and 16.2.1). The proof of Corollary 16.4.1 is thus complete. □

# Bibliography

[1]   BECK, C., BECKER, S., GROHS, P., JAAFARI, N., AND JENTZEN, A.  Solving stochastic differential equations and Kolmogorov equations by means of deep learning. *arXiv:1806.00421* (2018), 56 pages.

[2]   BOYD, S., AND VANDENBERGHE, L.  *Convex optimization.* Cambridge University Press, 2004.

[3]   BUBECK, S.  Convex Optimization: Algorithms and Complexity. *Foundations and Trends® in Machine Learning* 8, 3-4 (2015), 231–357. ISSN: 1935-8237. URL: http://dx.doi.org/10.1561/2200000050.

[4]   CARL, B., AND STEPHANI, I.  *Entropy, compactness and the approximation of operators.* Vol. 98. Cambridge Tracts in Mathematics. Cambridge University Press, Cambridge, 1990, x+277. ISBN: 0-521-33011-4. URL: https://doi.org/10.1017/CBO9780511897467.

[5]   COLEMAN, R.  *Calculus on normed vector spaces.* Springer Science & Business Media, 2012.

[6]   CUCKER, F., AND SMALE, S.  On the mathematical foundations of learning. *Bull. Amer. Math. Soc. (N.S.)* 39, 1 (2002), 1–49. ISSN: 0273-0979. URL: https://doi.org/10.1090/S0273-0979-01-00923-5.

[7]   DEREICH, S., AND MUELLER-GRONBACH, T.  General multilevel adaptations for stochastic approximation algorithms. *arXiv:1506.05482* (2017), 33 pages.

[8]   DUCHI, J.  *Probability Bounds.* https://stanford.edu/~jduchi/projects/probability_bounds.pdf.

[9]   E, W., HAN, J., AND JENTZEN, A.  Deep learning-based numerical methods for high-dimensional parabolic partial differential equations and backward stochastic differential equations. *Communications in Mathematics and Statistics* 5, 4 (2017), 349–380. ISSN: 2194-671X. URL: https://doi.org/10.1007/s40304-017-0117-6.

[10]  EINSIEDLER, M., AND WARD, T.  *Functional analysis, spectral theory, and applications.* Vol. 276. Graduate Texts in Mathematics. Springer, Cham, 2017, xiv+614. ISBN: 978-3-319-58539-0; 978-3-319-58540-6.

[11]  GROHS, P., HORNUNG, F., JENTZEN, A., AND ZIMMERMANN, P.  Space-time error estimates for deep neural network approximations for differential equations. *arXiv:1908.03833* (2019).

[12]  HAN, J., JENTZEN, A., AND E, W.  Solving high-dimensional partial differential equations using deep learning. *Proc. Natl. Acad. Sci. USA* 115, 34 (2018), 8505–8510.

[13]  HENRY, D.  *Geometric theory of semilinear parabolic equations.* Vol. 840. Lecture Notes in Mathematics. 348 pages. Springer-Verlag, Berlin, 1981, iv+348. ISBN: 3-540-10557-3.

[14]  HINTON, G., SRIVASTAVA, N., AND SWERSKY, K.  *Lecture 6e: Rmsprop: Divide the gradient by a running average of its recent magnitude.* http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf. 26-31, Accessed: 01.12.2017.

[15]  HUTZENTHALER, M., JENTZEN, A., KRUSE, T., AND NGUYEN, T. A.  A proof that rectified deep neural networks overcome the curse of dimensionality in the numerical approximation of semilinear heat equations. *SN Partial Differ. Equ. Appl.* 10, 1 (2020). URL: https://doi.org/10.1007/s42985-019-0006-9.

[16]  JENTZEN, A., KUCKUCK, B., NEUFELD, A., AND von WURSTEMBERGER, P.  Strong error analysis for stochastic gradient descent optimization algorithms. *arXiv:1801.09324* (2018), 75 pages.

[17]  JENTZEN, A., AND von WURSTEMBERGER, P.  Lower error bounds for the stochastic gradient descent optimization algorithm: Sharp convergence rates for slowly and fast decaying learning rates. *J. Complexity* 57 (2020), 101438. ISSN: 0885-064X. URL: https://doi.org/10.1016/j.jco.2019.101438.

[18]  KINGMA, D. P., AND BA, J.  Adam: A Method for Stochastic Optimization. *arXiv:1412.6980* (2014), 15 pages. arXiv: 1412.6980. URL: http://arxiv.org/abs/1412.6980.

[19]  KLENKE, A.  *Probability Theory.* 2nd ed. Universitext. Springer-Verlag London Ltd., 2014.

[20]  *Lasagne: Updates.* http://lasagne.readthedocs.io/en/latest/modules/updates.html#lasagne.updates.rmsprop. [Accessed 6-December-2017].

[21]  LeCun, Y., BENGIO, Y., AND HINTON, G.  Deep learning. *Nature* 521 (2015), 436–444.

[22]  NESTEROV, Y.  A method of solving a convex programming problem with convergence rate O (1/k2). In: *Soviet Mathematics Doklady.* Vol. 27. 1983, 372–376.

[23]  NESTEROV, Y.  *Introductory lectures on convex optimization: A basic course.* Vol. 87. Springer Science & Business Media, 2013.

[24]  PETERSEN, P.  *Linear Algebra.* Undergraduate Texts in Mathematics. Springer New York, 2012. ISBN: 9781461436119. URL: https://books.google.ch/books?id=3klgLwEACAAJ.

[25]  POLYAK, B. T.  Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics* 4, 5 (1964), 1–17.

[26]  RUDER, S.  An overview of gradient descent optimization algorithms. *arXiv:1609.04747* (2016), 12 pages.

[27]  SILVESTER, J. R.  Determinants of block matrices. *The Mathematical Gazette* 84, 501 (2000), 460–467.

[28]  *Tensorflow: RMSProp.* https://www.tensorflow.org/api_docs/python/tf/train/RMSPropOptimizer. [Accessed 6-December-2017].

[29]  TROPP, J. A.  An Elementary Proof of the Spectral Radius Formula for Matrices. In: http://users.cms.caltech.edu/~jtropp/notes/Tro01-Spectral-Radius.pdf. [Accessed 16-February-2018]. 2001.

[30]  ZEILER, M. D.  ADADELTA: An Adaptive Learning Rate Method. *arXiv:1212.5701* (2012), 6 pages. arXiv: 1212.5701. URL: http://arxiv.org/abs/1212.5701.