

# Statistical Learning Theory for Neural Operators

**Niklas Reinhardt**

NIKLAS.REINHARDT@IWR.UNI-HEIDELBERG.DE

*Interdisziplinäres Zentrum für wissenschaftliches Rechnen  
Universität Heidelberg  
Im Neuenheimer Feld 205, 69120 Heidelberg, Germany*

**Sven Wang**

SVEN.WANG@EPFL.CH

*Institute of Mathematics  
École Polytechnique Fédérale de Lausanne (EPFL)  
Station 8, CH-1015 Lausanne, Switzerland*

**Jakob Zech**

JAKOB.ZECH@UNI-HEIDELBERG.DE

*Interdisziplinäres Zentrum für wissenschaftliches Rechnen  
Universität Heidelberg  
Im Neuenheimer Feld 205, 69120 Heidelberg, Germany*

**Editor:** Benjamin Guedj

## Abstract

We present statistical convergence results for the learning of (possibly) non-linear mappings in infinite-dimensional spaces. Specifically, given a map  $G_0 : \mathcal{X} \rightarrow \mathcal{Y}$  between two separable Hilbert spaces, we analyze the problem of recovering  $G_0$  from  $n \in \mathbb{N}$  noisy input-output pairs  $(x_i, y_i)_{i=1}^n$  with  $y_i = G_0(x_i) + \varepsilon_i$ ; here the  $x_i \in \mathcal{X}$  represent randomly drawn “design” points, and the  $\varepsilon_i$  are assumed to be either i.i.d. white noise processes or subgaussian random variables in  $\mathcal{Y}$ . We provide general convergence results for least-squares-type empirical risk minimizers over compact regression classes  $\mathbf{G} \subseteq L^\infty(\mathcal{X}, \mathcal{Y})$ , in terms of their approximation properties and metric entropy bounds, which are derived using empirical process techniques. This generalizes classical results from finite-dimensional nonparametric regression to an infinite-dimensional setting. As a concrete application, we study an encoder-decoder based neural operator architecture termed FrameNet. Assuming  $G_0$  to be holomorphic, we prove algebraic (in the sample size  $n$ ) convergence rates in this setting, thereby overcoming the curse of dimensionality. To illustrate the wide applicability, as a prototypical example we discuss the learning of the non-linear solution operator to a parametric elliptic partial differential equation.

**Keywords:** nonparametric estimation, neural networks, operator learning, minimax convergence rates, empirical risk minimization, partial differential equations

## 1. Introduction

Learning non-linear relationships of high- and infinite-dimensional data is a fundamental problem in modern statistics and machine learning. In recent years, “Operator Learning” has emerged as a powerful tool for analyzing and approximating mappings  $G_0$  between *infinite-dimensional* spaces (Li et al., 2020; Hesthaven and Ubbiali, 2018; Bhattacharya et al., 2021; Lu et al., 2021; Raonic et al., 2023; Anandkumar et al., 2019; Owhadi and Yoo, 2019; Nelsen and Stuart, 2024; Kovachki et al., 2024b). The primary motivation for

considering truly infinite-dimensional data stems from applications in the natural sciences, where inputs and outputs of operators are elements in function spaces. For instance,  $G_0$  could be the operator relating an initial condition  $x$  of a dynamical system to the state  $G_0(x)$  of the system after a certain time, or a coefficient-to-solution map of a parametric partial differential equation (PDE).

For finite-dimensional inputs and outputs, nonparametric regression is the standard framework for inferring general, non-linear relationships. There, one aims to reconstruct some “ground truth”  $G_0 : \mathbb{R}^d \rightarrow \mathbb{R}^m$ ,  $d, m \in \mathbb{N}$ , from noisy data  $(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}^m$ ,  $i = 1, \dots, n$ , generated via  $y_i = G_0(x_i) + \varepsilon_i$ , where  $x_i$  are called the “design points” and  $\varepsilon_i$  are typically independent and identically distributed (i.i.d.) noise variables. In the framework of empirical risk minimization (ERM), one chooses a suitable function class  $\mathbf{G}$  of mappings from  $\mathbb{R}^d$  to  $\mathbb{R}^m$  and some loss function  $L : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$  measuring the discrepancy between predictions  $G(x_i)$  and the data  $y_i$ . Statistical estimation is achieved by minimizing

$$\hat{G}_n \in \arg \min_{G \in \mathbf{G}} J_n(G), \quad J_n(G) := \frac{1}{n} \sum_{i=1}^n L(G(x_i), y_i), \quad (1.1)$$

assuming that minimizers exist. In the finite-dimensional setting, statistically optimal convergence rates for such estimators were established for least-squares, maximum likelihood, and more generally “minimum contrast” estimators (van de Geer, 2000; Barron et al., 1999; Birgé and Massart, 1993); see also Schmidt-Hieber (2020a) where such results are shown for ERMs over neural network classes. However, as is well-known, both—approximation rates ( DeVore et al., 1989; DeVore and Lorentz, 1993) as well as statistical convergence rates (van de Geer, 2000; Giné and Nickl, 2016) over classical smoothness classes—deteriorate exponentially in terms of the dimension  $d$ . This renders computations practically infeasible for large  $d$ . This phenomenon is referred to as the *curse of dimensionality*, see also Section 4.1 ahead.

The framework for operator learning considered in this paper can be viewed as a direct extension of (1.1) to infinite dimensions. Given Hilbert spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , and a mapping  $G_0 : \mathcal{X} \rightarrow \mathcal{Y}$ , the goal is to reconstruct  $G_0$  from “training data”  $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$  with

$$y_i = G_0(x_i) + \varepsilon_i,$$

where the regression class  $\mathbf{G}$  is a suitable set of measurable mappings between  $\mathcal{X}$  and  $\mathcal{Y}$  and  $\varepsilon_i$  are centered noise variables, see Section 2 for details. This “supervised learning” setting underlies popular methods such as the PCA-Net (Hesthaven and Ubbiali, 2018; Bhattacharya et al., 2021).

We also mention the framework of “physics-informed learning” which is common in operator learning (relevant e.g. for the DeepONet Lu et al. 2021), but which is not considered in the present manuscript. Here, information on the ground truth  $G_0$  is not known in the form of input-output pairs, but instead is implicitly described via

$$\mathcal{N}(x, G_0(x)) = 0,$$

where  $\mathcal{N} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{Z}$  for a third vector space  $\mathcal{Z}$ . Typically,  $\mathcal{N}(x, \cdot)$  encodes a family of differential operators parametrized by  $x$  which represent the underlying physical model. In

this case, the loss to minimize is a residual of the form  $\sum_{i=1}^n \|\mathcal{N}(x_i, G(x_i))\|_{\mathcal{Z}}^2$ , thus leading to an “unsupervised learning problem”. The two cases, supervised and unsupervised learning, can also be combined. Various different (neural network based) architectures (i.e. regression classes  $\mathbf{G}$ ) have been proposed in recent years for the purpose of supervised or unsupervised operator learning.

### 1.1 Outline and Contributions

We provide statistical convergence results for operator learning which do not suffer from the curse of dimensionality, and which can be applied to prototypical problems in the PDE literature. We first develop our theory in an abstract setting for ERM over classes of mappings between separable Hilbert spaces, and later apply our theory to concrete examples. In doing so, we build upon and synthesize influential proof techniques from nonparametric statistics, in particular M-estimation (van de Geer, 2001; Nickl et al., 2020), approximation theory for parametric PDEs (Cohen et al., 2011; Cohen and DeVore, 2015; Herrmann et al., 2024), and empirical process theory (Talagrand, 2005; Dirksen, 2015).

To illustrate the scope of our contributions, we start by stating a convergence result for the elliptic “Darcy flow” problem on the  $d$ -dimensional torus. This is a standard example in PDE driven forward and inverse problems (Babuška et al., 2010; Chkifa et al., 2015b; Cliffe et al., 2011; Schwab and Stuart, 2012; Stuart, 2010; Nickl, 2023; Nickl et al., 2020). We aim for an informal exposition here, with full details given in Section 4: Denote by  $\mathbb{T}^d$  the  $d$ -dimensional torus, fix a smooth source function  $f : \mathbb{T}^d \rightarrow (0, \infty)$  and let  $a_{\min} > 0$ . For a sufficiently smooth and uniformly positive conductivity  $a : \mathbb{T}^d \rightarrow \mathbb{R}$ , denote by  $G_0(a)$  the unique solution of the elliptic PDE

$$-\nabla \cdot (a \nabla u) = f \quad \text{on } \mathbb{T}^d \quad \text{and} \quad \int_{\mathbb{T}^d} u(x) dx = 0. \quad (1.2)$$

Now let  $\gamma$  be some probability distribution on  $L^2(\mathbb{T}^d)$  such that

$$\text{supp}(\gamma) \subseteq \{a \in H^{\mathfrak{s}}(\mathbb{T}^d) : \inf_{x \in \mathbb{T}^d} a(x) \geq a_{\min}, \|a\|_{H^{\mathfrak{s}}(\mathbb{T}^d)} \leq R\}, \quad (1.3)$$

for some  $R > 0$ ,  $\mathfrak{s} > 2d + 1$ . Suppose we observe noisy input-output pairs  $(a_i, y_i)_{i=1}^n$  given by  $y_i = G_0(a_i) + \varepsilon_i$ , where the  $\varepsilon_i$  are independent  $L^2(\mathbb{T}^d)$ -Gaussian white noise processes (Section 2). The operator  $G_0$  can then be learned from this data as stated in the next theorem (for details see Section 4.2); it regards empirical risk minimizers  $\hat{G}_n$  over the so-called FrameNet  $\mathbf{G}_{\text{FN}}$ , which corresponds to a neural network based class of measurable mappings  $\mathcal{X} \rightarrow \mathcal{Y}$  (Section 3). Formally,  $\hat{G}_n$  is defined as a minimizer of the least squares objective

$$\hat{G}_n \in \arg \min_{G \in \mathbf{G}_{\text{FN}}} \frac{1}{n} \sum_{i=1}^n \|y_i - G(a_i)\|_{L^2(\mathbb{T}^d)}^2, \quad (1.4)$$

although a suitable modification is required to make this mathematically rigorous (Section 2.1).

**Theorem 1 (Informal)** *Consider the operator  $G_0$  from the Darcy problem on the  $d$ -dimensional torus  $\mathbb{T}^d$  ( $d \geq 2$ ), and suppose that  $\gamma$  satisfies (1.3) for some  $\mathfrak{s} > 3d/2 + 1$  and  $a_{\min} > 0$ . Fix  $\tau > 0$  (arbitrarily small).*

Then there exists a constant  $C$  such that for each  $n \in \mathbb{N}$  there exists a FrameNet class  $\mathbf{G}_{\text{FN}}(n)$  and any empirical risk minimizer  $\hat{G}_n$  in (1.4) satisfies<sup>1</sup>

$$\mathbb{E}_{G_0} \left[ \int \|\hat{G}_n(a) - G_0(a)\|_{L^2(\mathbb{T}^d)}^2 d\gamma(a) \right] \leq C n^{-\frac{2\mathfrak{s}+2-3d}{2\mathfrak{s}+2-d} + \tau}. \quad (1.5)$$

The most significant feature of the above statement is that, although  $G_0$  maps between infinite-dimensional spaces, the convergence rate in (1.5) is algebraic in  $n$ . Thus, it does not suffer from a curse of dimensionality with respect to the operator input  $a$ . The strong dependence on the spatial dimension  $d$  remains, but this reflects a curse of dimensionality that is generally unavoidable due to finite spatial smoothness, e.g., DeVore et al. (1989). Moreover, for infinitely smooth input data ( $\mathfrak{s} \rightarrow \infty$ ), the convergence rate gets arbitrarily close to  $n^{-1}$ . The classes  $\mathbf{G}_{\text{FN}}$ , whose existence is postulated by the theorem, can be precisely characterized in terms of the sparsity, depth, width and other network class parameters, which are chosen in terms of the statistical sample size  $n$ . We also note that the regularity assumption  $\mathfrak{s} > 3d/2 + 1$  was made here for convenience and can be weakened to  $\mathfrak{s} > 3d/2$ , see Theorem 27 and Remark 28 below.

To achieve Theorem 1 and several other related results, we build our theory in multiple steps. In Section 2, a general regression framework for mappings between Hilbert spaces is considered. Our first main result, Theorem 5, gives a non-asymptotic concentration upper bound on the empirical risk between  $\hat{G}_n$  and  $G_0$ , with respect to the design points  $a_i$ . The upper bound is quantified in terms of the metric entropy of the “regression class”  $\mathbf{G}$  and the best approximation of  $G_0$  from  $\mathbf{G}$ . Theorem 6 strengthens this statement to  $L^2(\gamma)$ -loss, for the case of random design points  $a_i \sim \gamma$ . These results provide an operator learning analogue to classical convergence rates in nonparametric regression. The proofs rely on probabilistic generic chaining techniques (Talagrand, 2014; Dirksen, 2015) and “slicing” arguments as introduced in van de Geer (2001), which we generalize to the current setting. We also note our proofs contrast existing nonparametric statistical analyses of neural networks (Schmidt-Hieber, 2020a) for real-valued regression, where generic chaining techniques were not required for obtaining optimal rates (up to log-factors).

In the second part of this work, we apply our statistical results to the specific deep operator network class  $\mathbf{G}_{\text{FN}}$  termed “FrameNet” and introduced in Herrmann et al. (2024). Together with their underlying decoder-encoder structure and feedforward neural network structure, FrameNet classes are defined in Section 3. These classes are known to satisfy good approximation properties for holomorphic operators, a property which is fulfilled for the Darcy problem (1.2) and more broadly a wide range of PDE based problems (Cohen et al., 2010, 2011; Cohen and DeVore, 2015; Jerez-Hanckes et al., 2017; Harbrecht et al., 2016; Henríquez and Schwab, 2021; Hiptmair et al., 2018; Spence and Wunsch, 2023; Cohen et al., 2018). In Section 3, we identify such operator holomorphy as a key regularity property which allows to derive “dimension-free” statistical convergence rates. By extending approximation theoretic results from Herrmann et al. (2024), see Theorem 18 below, as well as establishing metric entropy bounds for  $\mathbf{G}_{\text{FN}}$  based on Schmidt-Hieber (2020a), we obtain algebraic convergence results for ERMs over FrameNet classes, for reconstructing holomorphic operators  $G_0$ . Specifically, Theorem 23 bounds the  $L^2$ -risk  $\mathbb{E}[\|\hat{G}_n - G_0\|_{L^2(\gamma)}^2] \lesssim n^{-\kappa/(\kappa+1)}$ ,

---

1. Here and in the following  $\mathbb{E}_{G_0}$  denotes the expectation w.r.t. the random data  $(x_i, y_i)_i$  generated by the ground truth  $G_0$ . Similarly, we write  $\mathbb{P}_{G_0}$  for corresponding probabilities.

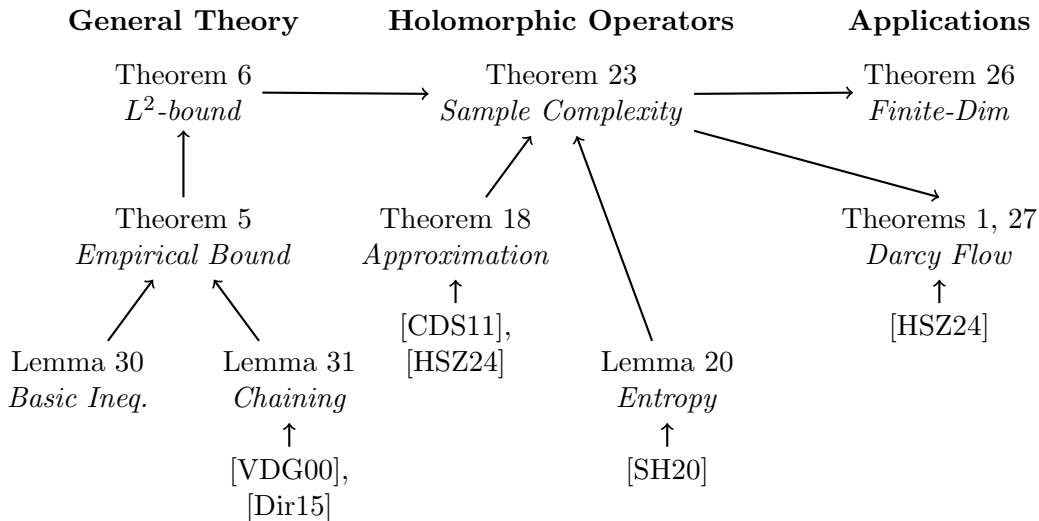


Figure 1: Dependency diagram of the main theorems, lemmas and references in this work.

where  $\kappa > 0$  denotes the *approximation* rate established in Herrmann et al. (2024). We treat the case of ReLU and RePU (Li, 2020) activation functions for both sparse and fully-connected architectures.

In Section 4, we illustrate the usefulness of our general theory in two concrete settings. First, we show how our theory recovers well-known minimax-optimal convergence rates for real-valued regression (i.e.,  $\mathcal{Y} = \mathbb{R}$ ) on  $d$ -dimensional domains. This proves that our abstract results from Section 2 cannot be improved *in general*, although matching lower bounds are yet unknown in the infinite-dimensional setting. Thereafter, Section 4.2 demonstrates how our theory can be used to yield the first algebraic convergence rates for a non-linear operator arising from PDEs—see in particular Theorem 27 and Remark 28, which underlie Theorem 1. Figure 1 summarizes the dependencies between the main theorems, lemmas and key references underlying the results in this work.

## 1.2 Existing Results

The approximation of mappings between infinite-dimensional spaces has been studied extensively in the context of Uncertainty Quantification, where  $G_0$  corresponds to the solution operator of a parameter dependent PDE. Various methodologies have been proposed and analyzed for this task, including for example compressed sensing (Doostan and Owhadi, 2011; Rauhut and Schwab, 2017), sparse-grid interpolation (Chkifa et al., 2013; Nobile et al., 2008), least-squares (Cohen et al., 2017; Chkifa et al., 2015a), and reduced basis methods (Quarteroni et al., 2016; Hesthaven et al., 2016). Recently, neural network approaches have become increasingly popular for this task as they provide a highly expressive and fast to evaluate parametrization of high-dimensional functions. These attributes make them particularly useful for learning surrogates in scientific applications (Hesthaven and Ubbiali, 2018; Lu et al., 2021; Li et al., 2020; Bhattacharya et al., 2021; Anandkumar et al.,

2019; O’Leary-Roseberry et al., 2022; O’Leary-Roseberry et al., 2024; Becker et al., 2024; Ciccì et al., 2022; Dal Santo et al., 2020; Kröpfl et al., 2022).

*Approximation Theory for Neural Operators.* First theoretical results on operator learning focused on the approximation error, establishing the existence of neural network architectures capable of approximating  $G_0$  up to a certain accuracy, with the error decreasing algebraic in terms of the number of learnable network parameters. For example, the works Schwab and Zech (2019); Kutyniok et al. (2022); Schwab and Zech (2023) showed that neural networks have sufficient expressivity to efficiently approximate certain (holomorphic) mappings  $G_0$ . Such results are based on the observation that the smoothness of  $G_0$  implies the image of this operator to have moderate  $n$ -widths, i.e. to be well approximated in moderate-dimensional linear subspaces (Dung et al., 2023; Cohen et al., 2010, 2011; Cohen and DeVore, 2015; Hoang and Schwab, 2014; Bachmayr et al., 2017). Specifically for DeepONets (Lu et al., 2021), such a result was obtained in Lanthaler et al. (2022), and for the presently considered architecture in Herrmann et al. (2024). Moreover, Schwab et al. (2025) provided a statement of this type for Lipschitz continuous operators by exploiting so-called superexpressivity of certain classes of neural networks.

*Statistical Theory for Neural Operators.* The analysis of sample complexity has received less attention so far. In Lanthaler (2023), the authors analyzed in particular the error of PCA encoders and decoders used for PCA-Net, but did not analyze the statistical error for the full operator. The work de Hoop et al. (2023) provides such a result for the estimation of linear and diagonalizable mappings from noisy data; for lower bounds see for example Chagny et al. (2025). For other work on “functional regression”, see, e.g., Greven and Scheipl (2017); Morris and Carroll (2006). An analysis for nonparametric regression of nonlinear mappings from noisy data in infinite dimensions was provided in Liu et al. (2024). There, the authors considered Lipschitz continuous mappings  $G_0$ , and proved consistency in the large data limit. Additionally they give convergence results, which, however, are subject to the curse of dimensionality. This is due to their very general assumption on the smoothness of  $G_0$ : It was shown very early (Mhaskar and Hahm, 1997), that the nonlinear  $n$ -width of Lipschitz operators in  $L^2$  decays only logarithmically, i.e. the number of (exact) data points needed for the reconstruction of the functional is exponential in the desired accuracy. Recently, Kovachki et al. (2024a) generalized these results and showed a generic curse of dimensionality for the reconstruction of Lipschitz operators and  $C^k$ -operators from exact data. Moreover, the authors show that under the existence of some intrinsic low-dimensionality allowing for fast approximation, also the dependence on the data complexity improves.

Concerning the case of noisy and holomorphic operators, we also refer to the recent works Adcock et al. (2025, 2024) who consider a setup similar to ours, and, unlike us, treat the more general case of Banach space valued functions. The authors derive upper bounds and concentration inequalities for the  $L^2$ -error, and lower bounds for the approximation error, in terms of a neural network based architecture. Key differences to our work include in particular that the results of Adcock et al. (2025, 2024) do not address convergence in the noise-polluted large data limit  $n \rightarrow \infty$ , nor do they directly yield convergence rates for concrete PDE models; the latter typically requires to exploit PDE regularity theory to show holomorphy in spaces of higher spatial regularity, leading to a multilevel decomposition of the operator. Moreover, the generalization results in these works either focus on linearly

parametrized models (Adcock et al., 2025, 2024), or are restricted to specific subsets of minimizers (Adcock et al., 2024, Theorem 3.2). Our analysis aims to fill this gap by providing generalization bounds for arbitrary minimizers and nonlinear parameter dependence, as is standard for neural networks.

*M-Estimation in Nonparametric Regression.* Convergence theory of M-estimators and (penalized) empirical risk minimizers was investigated around the 2000s in foundational works (van de Geer, 2000; van de Geer, 2001; Birgé and Massart, 1993; Barron et al., 1999). These works build on concentration inequalities for empirical processes using “chaining” techniques which date back to seminal contributions, see Talagrand (2014); Giné and Nickl (2016) and references therein. These techniques are known to produce minimax-optimal rates  $n^{-2s/(2s+d)}$  for ERMs over  $s$ -smooth Sobolev, Hölder and more generally, Besov smoothness classes of real-valued functions on bounded  $d$ -dimensional Euclidean domains. The analysis of neural-network based ERMs was initiated by the work Schmidt-Hieber (2020a), which considered regression over (compositional) Hölder classes on finite-dimensional domains, and was followed by several other works such as Suzuki (2019). We also mention Nickl et al. (2020); Nickl and Wang (2024); Agapiou and Wang (2024) which analyze ERMs in non-linear elliptic PDE-based inverse problems such as the “Darcy” flow problem studied here. The present Hilbert space setting falls outside the scope of such classical theory for scalar-valued functions. However, the derivation of our concentration inequalities for ERMs does build upon the same probabilistic empirical process machinery laid out above (Talagrand, 2014; Dirksen, 2015).

### 1.3 Notation

We write  $\mathbb{N} = \{1, 2, \dots\}$  and  $\mathbb{N}_0 = \{0, 1, 2, \dots\}$ . We write  $a_n \lesssim b_n, a_n \gtrsim b_n$  for real sequences  $(a_n)_{n \in \mathbb{N}}, (b_n)_{n \in \mathbb{N}}$  if  $a_n$  is respectively upper or lower bounded by a positive multiplicative constant which does not depend on  $n$  (but may well depend on other ambient parameters which we make explicit whenever confusion may arise). By  $a_n \simeq b_n$ , we mean that both  $a_n \lesssim b_n$  and  $a_n \gtrsim b_n$ .

For a pseudometric space  $(T, d)$  and any  $\delta > 0$ , let  $N(T, d, \delta)$  be the  $\delta$ -covering number of  $T$ , i.e. the minimal number of open  $\delta$ -balls in  $d$  needed to cover  $T$ . We denote the metric entropy of  $T$  by

$$H(T, d, \delta) = \log N(T, d, \delta).$$

Given a Borel probability measure  $\gamma$  on  $\mathcal{X}$  and a subset  $D \subseteq \mathcal{X}$ , we define the norms

$$\begin{aligned} \|G\|_{L^2(\mathcal{X}, \gamma; \mathbb{Y})}^2 &:= \int_{\mathcal{X}} \|G(x)\|_{\mathbb{Y}}^2 d\gamma(x), \\ \|G\|_{\infty, D} &:= \sup_{x \in D} \|G(x)\|_{\mathbb{Y}} \end{aligned}$$

and also write  $\|\cdot\|_{L^2(\gamma)}$  and  $\|\cdot\|_{\infty}$  if the underlying spaces are clear from context. The space of real-valued, square summable sequences indexed over  $\mathbb{N}$  is denoted by  $\ell^2(\mathbb{N})$ . The complexification of a real Hilbert space  $H$  is denoted by  $H_{\mathbb{C}}$  (Kirwan, 1997; Muñoz et al., 1999).

## 2. Regression in Hilbert Spaces

In the following we formalize a regression framework for mappings between Hilbert spaces.

### 2.1 Problem Formulation

Throughout, let  $\mathcal{X}$  and  $\mathcal{Y}$  denote two separable (real) Hilbert spaces with respective inner products  $\langle \cdot, \cdot \rangle_{\mathcal{X}}$ ,  $\langle \cdot, \cdot \rangle_{\mathcal{Y}}$  and suppose

$$G_0 : \mathcal{X} \rightarrow \mathcal{Y}$$

is some non-linear (Borel measurable) operator which we aim to reconstruct. The observed data are assumed to be noisy “input-output pairs”  $(x_i, y_i)_{i=1}^n \in (\mathcal{X} \times \mathcal{Y})^n$  given by

$$x_i \stackrel{\text{iid}}{\sim} \gamma, \quad i = 1, \dots, n, \quad (2.1a)$$

and

$$y_i = G_0(x_i) + \sigma \varepsilon_i, \quad i = 1, \dots, n, \quad (2.1b)$$

where  $\sigma > 0$  denotes a scalar “noise level”,  $\varepsilon_i$  are independent random noise variables and  $\gamma$  is a probability distribution on  $\mathcal{X}$ . The  $x_i \in \mathcal{X}$  are also referred to as the “design points”, and we write  $\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{X}^n$ . We will both derive results which are *conditional* on the design  $\mathbf{x}$ , as well as results for *random design*. To avoid confusions, we will use the notations  $P_{G_0}^{\mathbf{x}}$ ,  $\mathbb{E}_{G_0}^{\mathbf{x}}$  to denote probabilities and expectations under the distribution (2.1) with fixed design  $\mathbf{x}$ , and we use  $\mathbb{P}_{G_0}$ ,  $\mathbb{E}_{G_0}$  to denote probabilities and expectations with random design  $x_i \sim \gamma$ .

**Remark 2** *In practice, we will often deal with scenarios in which  $G_0$  is only defined on some measurable subset  $\mathcal{V} \subset \mathcal{X}$ , see e.g. the solution operator for the Darcy flow in Section 4.2.1. In this case, our results from Section 2.2 can be applied to any measurable extension of  $G_0$  onto  $\mathcal{X}$ . To apply the sample complexity results from Section 3.4 the extension additionally needs to be holomorphic in an open set containing  $\mathcal{V}$ , see Assumption 2 and Section 3 for the precise setup and statement.*

*White Noise Model.* We shall treat two different assumptions on the noise, the first being that the  $(\varepsilon_i)_{i=1}^n$  in (2.1) are independent copies of a  $\mathcal{Y}$ -white noise process. Recall that for any given separable Hilbert space  $\mathcal{Y}$ , the  $\mathcal{Y}$ -Gaussian white noise process is defined as the mean-zero Gaussian process  $\mathbb{W}_{\mathcal{Y}} = (\mathbb{W}_{\mathcal{Y}}(y) : y \in \mathcal{Y})$  indexed by  $\mathcal{Y}$  with “iso-normal” covariance structure

$$\mathbb{W}_{\mathcal{Y}}(y) \sim \mathcal{N}(0, \|y\|_{\mathcal{Y}}^2), \quad \text{Cov}(\mathbb{W}_{\mathcal{Y}}(y), \mathbb{W}_{\mathcal{Y}}(y')) = \langle y, y' \rangle_{\mathcal{Y}}, \quad \text{for all } y, y' \in \mathcal{Y}.$$

It is well-known that  $\mathbb{W}_{\mathcal{Y}}$  does not take values in  $\mathcal{Y}$  unless  $\dim(\mathcal{Y}) < \infty$ , but is interpreted as a stochastic process indexed by  $\mathcal{Y}$ , see Giné and Nickl (2016, p.19) for details. Nevertheless, we slightly abuse notation and use the common notation  $\langle \mathbb{W}_{\mathcal{Y}}, y \rangle_{\mathcal{Y}} := \mathbb{W}_{\mathcal{Y}}(y)$ .

Under this assumption, conditionally on  $x_i$  we interpret each observation  $y_i$  in (2.1) as a realisation of a Gaussian process  $(y_i(f) : f \in \mathcal{Y})$  with

$$\mathbb{E}[y_i(f)] = \langle G_0(x_i), f \rangle_{\mathcal{Y}}, \quad \text{Cov}(y_i(f), y_i(f')) = \langle f, f' \rangle_{\mathcal{Y}},$$

and we shall again use the notation  $\langle y_i, f \rangle_{\mathcal{Y}}$  to denote  $y_i(f)$  (see also Tsybakov 2009; Giné and Nickl 2016; Nickl et al. 2020 where this common viewpoint is explained in detail).

**Example 1** Let  $\mathcal{O} \subseteq \mathbb{R}^d$  be a bounded, smooth domain. Then, for  $\mathcal{Y} = L^2(\mathcal{O})$ , one can show that draws of an  $L^2(\mathcal{O})$ -white noise process a.s. take values in negative Sobolev spaces  $H^{-\kappa}$  for  $\kappa > d/2$  (Nickl, 2020; Castillo and Nickl, 2013).

*Sub-Gaussian Noise Model.* The second setting we consider is that of sub-Gaussian noise. We say that a random vector  $X$  taking values in  $\mathcal{Y}$  is sub-Gaussian with parameter  $\eta > 0$  if  $\mathbb{E}[X] = 0$  and

$$\mathbb{P}(\|X\|_{\mathcal{Y}} \geq t) \leq 2 \exp\left(-\frac{t^2}{2\eta^2}\right), \quad \text{for all } t \geq 0.$$

In the sub-Gaussian noise model, we assume that  $(\varepsilon_i)_{i=1}^n$  in (2.1) are independent sub-Gaussian variables in  $\mathcal{Y}$  with parameter  $\eta = 1$ .

### 2.1.1 EMPIRICAL RISK MINIMIZATION

Let  $\mathbf{G}$  be a class of (measurable) operators  $\mathbf{G} \ni G : \mathcal{X} \rightarrow \mathcal{Y}$ . We would like to study classical empirical risk minimizers of least-squares type over  $\mathbf{G}$ . Specifically, given regression data  $(x_i, y_i)_{i=1}^n$ , consider the empirical risk

$$\tilde{I}_n(G) := \frac{1}{n} \sum_{i=1}^n \|y_i - G(x_i)\|_{\mathcal{Y}}^2, \quad \tilde{I}_n : \mathbf{G} \rightarrow [0, \infty]. \quad (2.2)$$

However, this functional takes finite values almost surely only in the sub-Gaussian noise model. In the white noise model, since  $y_i \notin \mathcal{Y}$ , it holds  $\tilde{I}_n(G) = \infty$  almost surely—we thus consider a modified definition of least-squares type estimators which is common in the literature on regression with white noise (Nickl et al., 2020; Giné and Nickl, 2016). Instead of (2.2), we consider

$$I_n(G) = \frac{1}{n} \sum_{i=1}^n -2\langle G(x_i), y_i \rangle_{\mathcal{Y}} + \|G(x_i)\|_{\mathcal{Y}}^2, \quad I_n : \mathbf{G} \rightarrow \mathbb{R}, \quad (2.3)$$

which takes finite values a.s. also in the white noise model. Note that the latter objective function can be obtained from (2.2) by formally subtracting the term  $n^{-1} \sum_{i=1}^n \|y_i\|_{\mathcal{Y}}^2$  which exhibits no dependency on  $G$ . Therefore in the sub-Gaussian noise model the minimization of  $\tilde{I}_n$  and  $I_n$  are equivalent which is why we consider (2.3) in the following. We will denote minimizers of  $I_n(G)$  by  $\hat{G}_n$ .

Our assumptions on the class  $\mathbf{G}$  in the ensuing theorems will ensure that a *measurable choice* of minimizers  $\hat{G}_n$  of  $I_n$  exists, see Theorem 5 (i). However, the ERM  $\hat{G}_n$  will in general not be unique, since we do not impose convexity on  $\mathbf{G}$ . The reason is that our main application, the NN-based FrameNet class  $\mathbf{G}_{\text{FN}}$ , is non-convex.

**Remark 3 (Connection to maximum likelihood)** *In the white noise model, it follows from the Cameron-Martin theorem (see for example Giné and Nickl 2016, Theorem 2.6.13) that  $-nI_n(G)/(2\sigma^2)$  constitutes the negative log-likelihood of the (dominated) statistical model arising from (2.1) with white noise. In this case  $\hat{G}_n$  can also be interpreted as a (nonparametric) maximum likelihood estimator over the class  $\mathbf{G}$ .*

**Remark 4** Consider nonparametric regression of an unknown function  $f : \mathcal{O} \rightarrow \mathbb{R}$  for some bounded, smooth domain  $\mathcal{O}$ . Here, it is well-known that the observation white noise error model, where data is given by  $Y = f + \sigma\mathbb{W}$  (with  $\mathbb{W}$  a  $L^2(\mathcal{O})$ -white noise process) is asymptotically equivalent in a Le Cam-sense to an observation model with  $m$  “equally spaced” (random or deterministic) observation points throughout  $\mathcal{O}$ ,

$$Y_i = f(z_i) + \eta_i, \quad i = 1, \dots, m,$$

with i.i.d.  $N(0, 1)$  errors, where the equivalence holds for  $\sigma \asymp 1/\sqrt{m}$ , see Reiß (2008). Therefore, our observation model (2.1) may be viewed as a simplified proxy.

## 2.2 General Convergence Results

Let  $\mathbf{G}$  be a class of operators mapping from  $\mathcal{X}$  to  $\mathcal{Y}$ . For any fixed  $\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{X}^n$  and (Borel) measurable map  $G : \mathcal{X} \rightarrow \mathcal{Y}$ , we denote the empirical seminorm induced by  $\mathbf{x}$  with

$$\|G\|_n^2 = \frac{1}{n} \sum_{i=1}^n \|G(x_i)\|_{\mathcal{Y}}^2. \quad (2.4)$$

For any element  $G^* \in \mathbf{G}$  and  $\delta > 0$ , define the localized classes

$$\mathbf{G}_n^*(\delta) = \{G \in \mathbf{G} : \|G - G^*\|_n \leq \delta\},$$

and denote its metric entropy integral by

$$J(\delta) = J(\mathbf{G}_n^*(\delta), \|\cdot\|_n) := \int_0^\delta H^{\frac{1}{2}}(\mathbf{G}_n^*(\delta), \|\cdot\|_n, \rho) d\rho. \quad (2.5)$$

The following result provides a general convergence theorem for empirical risk minimizers with high probability, which relates the empirical risk of ERMs over some operator class  $\mathbf{G}$  to the metric entropy of  $\mathbf{G}$ . It can be viewed as a generalisation of classical convergence results for sieved M-estimators (van de Geer, 2000) to Hilbert space valued functions. The proof can be found in Appendix B.1.

**Theorem 5** For some measurable  $G_0 : \mathcal{X} \rightarrow \mathcal{Y}$ , let the data  $(x_i, y_i)_{i=1}^n$  arise from (2.1) either with white noise or with sub-Gaussian noise. Let  $\mathbf{G}$  be a class of measurable maps from  $\mathcal{X} \rightarrow \mathcal{Y}$ , let  $G^* \in \mathbf{G}$ , and let  $\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{X}^n$  be such that the following holds.

- (a) There exists a constant  $C > 0$  s.t.  $\mathbf{G}$  is compact with respect to some norm  $\|\cdot\|$  satisfying  $\|\cdot\|_n \leq C\|\cdot\|$ .
- (b) There exists  $\Psi_n : (0, \infty) \rightarrow [0, \infty)$  s.t.  $\Psi_n(\delta) \geq J(\mathbf{G}_n^*(\delta), \|\cdot\|_n)$  for all  $\delta > 0$  and

$$\delta \mapsto \frac{\Psi_n(\delta)}{\delta^2} \quad \text{is non-increasing for } \delta \in (0, \infty).$$

Then the following holds.

- (i) Minimizers  $\hat{G}_n$  of the empirical risk (2.3) exist, and there is a measurable selection (with respect to the data  $(x_i, y_i)_{i=1}^n$ ) of such a minimizer.

(ii) Fix any measurable selection  $\hat{G}_n$  from part (i). Then there exists a universal constant  $C_{\text{Ch}} > 0$  (see Lemma 32) such that for any  $G^* \in \mathbf{G}$  as above, any positive sequence  $(\delta_n)_{n \in \mathbb{N}}$  satisfying

$$\sqrt{n}\delta_n^2 \geq 32C_{\text{Ch}}\sigma\Psi_n(\delta_n), \quad (2.6)$$

and for any

$$R_n \geq \max \left\{ \delta_n, \frac{4C_{\text{Ch}}\sigma}{\sqrt{n}}, \sqrt{2}\|G^* - G_0\|_n \right\},$$

it holds that

$$\mathbb{P}_{G_0}^{\mathbf{x}}(\|\hat{G}_n - G_0\|_n \geq R_n) \leq 2 \exp\left(-\frac{nR_n^2}{16C_{\text{Ch}}^2\sigma^2}\right). \quad (2.7)$$

The lower bound for  $R_n$  in (ii), which determines the convergence rate of  $\|\hat{G}_n - G_0\|_n$ , is typically optimized by balancing the “stochastic term”  $\delta_n$  and the empirical approximation error  $\|G^* - G_0\|_n$ , see, e.g., Theorem 23 below. Note that Theorem 5 gives a convergence rate with respect to the *empirical seminorm*  $\|\cdot\|_n$ . Therefore, when the design points  $\mathbf{x}$  are random, the norm itself is also random, and the assumptions (a) and (b) in Theorem 5 have to be understood conditional on possible realizations of  $\mathbf{x}$ . Theorem 6 below will give a corresponding concentration inequality on the  $\|\cdot\|_{L^2(\gamma)}$ -error under the assumption of i.i.d. random design  $x_i \sim \gamma$ . In this setting, a sufficient condition for (b) to be satisfied almost surely is to take  $\Psi_n$  as an upper bound for the entropy integral with uniform entropy  $H(\mathbf{G}, \|\cdot\|_{\infty, \text{supp}(\gamma)}, \delta)$ .

The compactness of  $\mathbf{G}$  in part (a) is needed for the existence of the ERM  $\hat{G}_n$  in (2.3). The existence result for  $\hat{G}_n$  (see for example Nickl 2007, Proposition 5) requires compact *metric* spaces, which is why we assume compactness w.r.t. a *norm*  $\|\cdot\|$  stronger than the empirical *seminorm*  $\|\cdot\|_n$ .

In particular, compactness implies separability of  $\mathbf{G}$  with respect to  $\|\cdot\|_n$ , which in turn is needed to guarantee the measurability of certain suprema of empirical processes ranging over  $\mathbf{G}$ . See the proof of Theorem 5 below, in particular (B.4) and (B.6). The technical growth restriction in (b) on the function  $\Psi_n(\delta)$  (i.e., our upper bound for the entropy integral) is required for the “peeling device” in (B.4). In particular, this assumption is satisfied in case that  $H(\mathbf{G}, \|\cdot\|_{\infty}, \delta) \lesssim \delta^{-\alpha}$  for any  $0 < \alpha < 2$ , see Corollary 8 below for details.

Theorem 5 provides a concentration inequality for the empirical seminorm  $\|\hat{G}_n - G_0\|_n$ . Under the assumption of randomly chosen design points  $x_i \sim \gamma$ ,  $x_i \in \mathcal{X}$ , this statement can be extended to a convergence result for the  $L^2(\gamma)$ -norm. To this end, we also need slightly stronger technical assumptions on the class  $\mathbf{G}$  with respect to the  $\|\cdot\|_{\infty, \text{supp}(\gamma)}$ -norm.

**Assumption 1** For some probability measure  $\gamma$  on  $\mathcal{X}$ , assume  $\mathbf{x} = (x_1, \dots, x_n)$  arise from i.i.d. draws  $x_i \sim \gamma$ . Let  $\mathbf{G}$  be a class of measurable maps  $\mathcal{X} \rightarrow \mathcal{Y}$ ,  $G^* \in \mathbf{G}$  and  $\|G_0\|_{\infty, \text{supp}(\gamma)} < \infty$ . Suppose

(a)  $\mathbf{G}$  is compact with respect to  $\|\cdot\|_{\infty, \text{supp}(\gamma)}$ ,

(b) *There exists a (deterministic) upper bound  $\Psi_n : (0, \infty) \rightarrow [0, \infty)$  such that for a.e.  $\mathbf{x} \sim \gamma^n$ , it holds  $\Psi_n(\delta) \geq J(\mathbf{G}_n^*(\delta), \|\cdot\|_n)$  for all  $\delta > 0$  and*

$$\delta \mapsto \frac{\Psi_n(\delta)}{\delta^2} \quad \text{is non-increasing for } \delta \in (0, \infty). \quad (2.8)$$

Note that Assumption 1 is strictly stronger than assumptions (a) and (b) in Theorem 5, since  $\|\cdot\|_{\infty, \text{supp}(\gamma)}$  is stronger than  $\|\cdot\|_n$  and  $\Psi_n$  in (2.8) does not depend on  $\mathbf{x}$ . In particular, Assumption 1 implies the existence of a measurable ERM  $\hat{G}_n$  by Theorem 5 (i).

The uniform  $\|\cdot\|_{\infty, \text{supp}(\gamma)}$  assumptions are used to control the concentration of the empirical seminorm  $\|\cdot\|_n$  around the ‘‘population norm’’  $\|\cdot\|_{L^2(\gamma)}$ , see Lemma 33 and also Lemma 34 below. Let us write  $\mathbf{F} = \{G - G_0 : G \in \mathbf{G}\}$ . As an immediate consequence of Assumption 1, there exists some  $M_{\mathbf{F}} < \infty$  such that

$$\sup_{F \in \mathbf{F}} \|F\|_{\infty, \text{supp}(\gamma)} = \sup_{G \in \mathbf{G}} \|G - G_0\|_{\infty, \text{supp}(\gamma)} \leq M_{\mathbf{F}}. \quad (2.9)$$

We can now state our main concentration inequality for the convergence of  $\|\hat{G}_n - G_0\|_{L^2(\gamma)}$ , which in particular provides a bound for the mean squared error  $\mathbb{E}_{G_0}[\|\hat{G}_n - G_0\|_{L^2(\gamma)}^2]$  as well. The proof of Theorem 6 can be found in Appendix B.2.

**Theorem 6 ( $L^2(\gamma)$ -Concentration under Random Design)** *Consider the nonparametric regression model (2.1) either with white noise or with sub-Gaussian noise and any measurable empirical risk minimizer  $\hat{G}_n$  from (2.3). Suppose that  $G_0$ ,  $\mathbf{G}$ ,  $G^*$ ,  $\gamma$  and  $\Psi_n(\cdot)$  are such that Assumption 1 holds. Then there exists some universal constant  $C > 0$  such that for any positive sequences  $(\delta_n)_{n \in \mathbb{N}}$  and  $(\tilde{\delta}_n)_{n \in \mathbb{N}}$  with*

$$\sqrt{n}\delta_n^2 \geq C\sigma\Psi_n(\delta_n) \quad \text{and} \quad n\tilde{\delta}_n^2 \geq CM_{\mathbf{F}}^2 H(\mathbf{G}, \|\cdot\|_{\infty, \text{supp}(\gamma)}, \tilde{\delta}_n), \quad (2.10)$$

all  $G^* \in \mathbf{G}$  as above and all

$$R_n \geq C \max \left\{ \delta_n, \tilde{\delta}_n, \|G^* - G_0\|_{\infty, \text{supp}(\gamma)}, \frac{\sigma + M_{\mathbf{F}}}{\sqrt{n}} \right\} \quad (2.11)$$

we have,

$$\mathbb{P}_{G_0} \left( \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq R_n \right) \leq 2 \exp \left( -\frac{nR_n^2}{C^2(\sigma^2 + M_{\mathbf{F}}^2)} \right). \quad (2.12)$$

To keep the presentation simple, we have left the numerical constants in the preceding theorem implicit. However, they can be made explicit, see the proof for details. The following bound on the mean squared error is obtained upon integration of the concentration inequalities from Theorem 5 and Lemma 33. Note that directly integrating the  $L^2(\gamma)$ -concentration, cf. (2.12), gives an approximation term in the uniform norm  $\|\cdot\|_{\infty, \text{supp}(\gamma)}$  (following 2.11), which has weaker convergence properties in general, see Theorem 18. For the proof of Corollary 7, see Appendix B.3.

**Corollary 7 ( $L^2(\gamma)$ -Mean Squared Error)** *Consider the setting of Theorem 6 and assume in addition that Assumption 1 is fulfilled for all  $G^* \in \mathbf{G}$  (with the same  $\Psi_n$ ). Then, for some universal constant  $C > 0$  and all  $n \in \mathbb{N}$ ,*

$$\mathbb{E}_{G_0} \left[ \|\hat{G}_n - G_0\|_{L^2(\gamma)}^2 \right] \leq C \left( \delta_n^2 + \tilde{\delta}_n^2 + \frac{\sigma^2 + M_{\mathbf{F}}^2}{n} \right) + 8 \inf_{G^* \in \mathbf{G}} \|G^* - G_0\|_{L^2(\gamma)}^2. \quad (2.13)$$

To demonstrate the typical use-cases of our abstract results, we summarize in the following corollary the rates which can be obtained under algebraic approximation properties of  $\mathbf{G}$  and two concrete scalings of the metric entropy of  $\mathbf{G}$ . These correspond to the typical entropy bounds satisfied by (i) some fixed  $n$ -independent, infinite-dimensional regression class, and (ii) an  $N$ -dimensional approximation class, where  $N$  is chosen in terms of  $n$ . In the following corollary,  $\lesssim$  refers to an inequality involving a constant independent of  $n$ ,  $N$  and  $\delta$ . For a proof of Corollary 8, see Appendix B.4.

**Corollary 8** *Consider the setting of Corollary 7. Let  $\mathbf{G} = \mathbf{G}(N)$ ,  $N \in \mathbb{N}$ , be a sequence of regression classes<sup>2</sup> such that  $\inf_{G^* \in \mathbf{G}(N)} \|G^* - G_0\|_{L^2(\gamma)}^2 \lesssim N^{-\beta}$  for some  $\beta > 0$  and all  $N \in \mathbb{N}$ . Denote the entropy by  $H(N, \delta) := H(\mathbf{G}(N), \|\cdot\|_{\infty, \text{supp}(\gamma)}, \delta)$ .*

(i) *If  $H(N, \delta) \lesssim \delta^{-\alpha}$  for some  $0 < \alpha < 2$ , then for all  $n \in \mathbb{N}$  and all  $N = N(n) = \lceil n^{\frac{2}{(2+\alpha)\beta}} \rceil$ , it holds*

$$\mathbb{E}_{G_0} \left[ \|\hat{G}_n - G_0\|_{L^2(\gamma)}^2 \right] \lesssim n^{-\frac{2}{2+\alpha}}.$$

(ii) *If  $H(N, \delta) \lesssim N \log(\delta^{-1})$ , then for all  $n \in \mathbb{N}$ ,  $n \geq 2$  and all  $N(n) = \lceil n^{\frac{1}{\beta+1}} \rceil$ , it holds*

$$\mathbb{E}_{G_0} \left[ \|\hat{G}_n - G_0\|_{L^2(\gamma)}^2 \right] \lesssim \log(n) n^{-\frac{\beta}{\beta+1}}.$$

**Remark 9 (Effective smoothness)** *In classical nonparametric regression over  $s$ -smooth function classes on  $[0, 1]^d$ , the entropy assumption  $H(N, \delta) \lesssim \delta^{-\alpha}$  from part (i) is fulfilled for  $\alpha = d/s$ , which yields the minimax-optimal rate  $n^{-\frac{2s}{2s+d}}$ , see Section 4.1 for details. Since the rate only depends on  $\alpha$  (or equivalently  $\alpha^{-1}$ ), we can think of  $\alpha^{-1} = s/d$  as the “effective smoothness” of the statistical model at hand.*

The sub-Gaussian noise model poses a more restrictive regularity assumption on  $\varepsilon_i$  than the assumption of white noise. It is possible to get  $L^2(\gamma)$ -convergence for sub-Gaussian noise in the case the entropy integral  $J(\delta)$  in (2.5) is not finite, such that Assumption 1 (b) does not hold. The details are shown in the next theorem, its proof is deferred to Appendix B.5.

**Theorem 10** *Consider the nonparametric regression model (2.1) with sub-Gaussian noise and the empirical risk from (2.3). Let Assumption 1 (a) hold, that means suppose that*

---

2. We may think of  $N$  as the number of parameters of  $\mathbf{G}$ , e.g. the size of the FrameNet class  $\mathbf{G}_{\text{FN}}$  in Section 3 below.

$\|G_0\|_{\infty, \text{supp}(\gamma)} < \infty$  and  $\mathbf{G}$  is compact with respect to  $\|\cdot\|_{\infty, \text{supp}(\gamma)}$ . Then there exists some universal constant  $C_1 > 0$ , such that for any positive sequences  $(\delta_n)_{n \in \mathbb{N}}$  and  $(\tilde{\delta}_n)_{n \in \mathbb{N}}$  with

$$n\delta_n^4 \geq C_1^2 \sigma^2 M_{\mathbf{F}}^2 H \left( \frac{\delta_n^2}{8\sigma^2 + \delta_n^2} \right) \quad \text{and} \quad n\tilde{\delta}_n^2 \geq 6M_{\mathbf{F}}^2 H(\mathbf{G}, \|\cdot\|_{\infty, \text{supp}(\gamma)}, \tilde{\delta}_n), \quad (2.14)$$

all  $G^* \in \mathbf{G}$  and all

$$R_n \geq \max \left\{ \delta_n, \tilde{\delta}_n, \|G^* - G_0\|_{\infty, \text{supp}(\gamma)}, \frac{M_{\mathbf{F}}}{\sqrt{n}} \right\},$$

we have

$$\mathbb{P}_{G_0} \left( \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq R_n \right) \leq 4 \exp \left( -\frac{nR_n^4}{C_1^2 \sigma^2 (1 + M_{\mathbf{F}}^2)} \right) + 2 \exp \left( -\frac{nR_n^2}{C_1^2 M_{\mathbf{F}}^2} \right). \quad (2.15)$$

Furthermore, for all  $n \in \mathbb{N}$  there exists a universal constant  $C_2 > 0$  such that

$$\mathbb{E}_{G_0} \left[ \|\hat{G}_n - G_0\|_{L^2(\gamma)}^2 \right] \leq C_2 \left( \delta_n^2 + \tilde{\delta}_n^2 + \frac{(1 + \sigma)(1 + M_{\mathbf{F}}^2)}{\sqrt{n}} + \inf_{G^* \in \mathbf{G}} \|G^* - G_0\|_{L^2(\gamma)}^2 \right). \quad (2.16)$$

**Remark 11** Consider  $\alpha \geq 2$  in Corollary 8 (i). Then  $J(\delta)$  in (2.5) is not necessarily finite, and hence Assumption 1 (b) need not be satisfied. Therefore Theorem 6 cannot be applied. In the sub-Gaussian noise case, we may still use the  $L^2(\gamma)$ -bound in (2.16) however. Similar to the proof of Corollary 8, it can then be shown that  $\delta_n^2 \lesssim n^{-\frac{1}{2+\alpha}}$  and  $\tilde{\delta}_n^2 \lesssim n^{-\frac{2}{2+\alpha}}$  satisfy (2.14). This yields

$$\mathbb{E}_{G_0} \left[ \|\hat{G}_n - G_0\|_{L^2(\gamma)}^2 \right] \lesssim n^{-\frac{1}{2+\alpha}},$$

i.e. half the convergence rate of the ‘‘chaining regime’’ considered in Corollary 8.

### 3. Learning Holomorphic Operators with FrameNet

In this section we first recall the NN-based operator class FrameNet from Herrmann et al. (2024, Section 2), see Sections 3.1 and 3.2. This will provide the regression class  $\mathbf{G}_{\text{FN}}$  over which to estimate  $G_0$ . Similar to for example PCA-Net (Hesthaven and Ubbiali, 2018), FrameNet consists of mappings

$$G = \mathcal{D}_y \circ g \circ \mathcal{E}_x, \quad (3.1)$$

for a linear encoder  $\mathcal{E}_x : \mathcal{X} \rightarrow \ell^2(\mathbb{N})$ , a linear decoder  $\mathcal{D}_y : \ell^2(\mathbb{N}) \rightarrow \mathcal{Y}$ , and a coefficient map  $\ell^2(\mathbb{N}) \rightarrow \ell^2(\mathbb{N})$ . The encoder maps an  $x \in \mathcal{X}$  to its coefficients in some representation system of  $\mathcal{X}$ . Conversely, the decoder builds a  $y \in \mathcal{Y}$  out of a coefficient series in  $\ell^2(\mathbb{N})$ . These representation systems consist of frames that are fixed a priori. The coefficient map  $g$  is represented by a feedforward neural network, that will be trained by ERM.

Subsequently, in Sections 3.3-3.4, we apply our analysis to the learning of *holomorphic* operators  $G_0 : \mathcal{X} \rightarrow \mathcal{Y}$ . For such mappings, the FrameNet architecture was shown in Herrmann et al. (2024) to be capable of overcoming the curse of dimensionality in terms of the *approximation error*. We generalize this property to the case of bounded network parameters in Subsection 3.3 and furthermore show metric entropy bounds. This allows us to prove that FrameNet can overcome the curse of dimensionality in the learning of holomorphic operators, both in terms of the *approximation capability* and in terms of *sample complexity*.

### 3.1 Representation Systems

We briefly recall basic definitions and properties of frames. For more details see for instance Christensen (2016).

#### 3.1.1 FRAMES

**Definition 12** A family  $\Psi = \{\psi_j : j \in \mathbb{N}\} \subset \mathcal{X}$  is called a frame of  $\mathcal{X}$ , if the analysis operator

$$T_\Psi : \mathcal{X} \rightarrow \ell^2(\mathbb{N}), \quad v \mapsto (\langle v, \psi_j \rangle_{\mathcal{X}})_{j \in \mathbb{N}}$$

is bounded and boundedly invertible between  $\mathcal{X}$  and  $\text{range}(T_\Psi) \subset \ell^2(\mathbb{N})$ .

Every orthonormal basis of  $\mathcal{X}$  is trivially a frame. Since Definition 12 merely requires bounded invertibility on the range of  $T_\Psi$ ,  $T_\Psi$  need not be surjective, and in particular  $\Psi$  need not consist of linearly independent vectors. The *frame bounds* of  $\Psi$  are defined as

$$\Lambda_\Psi := \|T_\Psi\|_{\mathcal{X} \rightarrow \ell^2} = \sup_{0 \neq v \in \mathcal{X}} \frac{\|T_\Psi v\|_{\ell^2}}{\|v\|_{\mathcal{X}}}, \quad \lambda_\Psi := \inf_{0 \neq v \in \mathcal{X}} \frac{\|T_\Psi v\|_{\ell^2}}{\|v\|_{\mathcal{X}}}, \quad (3.2)$$

the *synthesis operator*  $T'_\Psi$  as

$$T'_\Psi : \ell^2(\mathbb{N}) \rightarrow \mathcal{X}, \quad (v_i)_{i \in \mathbb{N}} \mapsto \mathbf{v}^T \Psi := \sum_{i \in \mathbb{N}} v_i \psi_i,$$

and finally the *frame operator* as  $Q_\Psi := T'_\Psi T_\Psi : \mathcal{X} \rightarrow \mathcal{X}$ . The frame operator  $Q_\Psi$  is boundedly invertible, self-adjoint and positive, see Christensen (2016, Lemma 5.1.5). The family  $\tilde{\Psi} := Q_\Psi^{-1} \Psi$  is a frame of  $\mathcal{X}$ , called the (*canonical*) *dual frame* of  $\mathcal{X}$ . The analysis operator of the dual frame is  $\tilde{T}_\Psi := T_\Psi (T'_\Psi T_\Psi)^{-1}$  and its frame bounds are  $\lambda_{\tilde{\Psi}}^{-1}$  and  $\Lambda_{\tilde{\Psi}}^{-1}$ .

**Definition 13** A family  $\Psi = \{\psi_j : j \in \mathbb{N}\} \subset \mathcal{X}$  is called a *Riesz basis* of  $\mathcal{X}$  if there exists a bounded, bijective operator  $A : \mathcal{X} \rightarrow \mathcal{X}$  and an orthonormal basis  $(e_j)_{j \in \mathbb{N}}$  with  $\psi_j = A e_j$  for all  $j \in \mathbb{N}$ .

A Riesz basis is a frame  $\Psi$  which is also a basis. Equivalently, a Riesz basis is a frame with  $\ker(T'_\Psi) = 0$  and therefore  $\text{range}(T'_\Psi) = \ell^2(\mathbb{N})$ . Moreover, the dual frame  $\tilde{\Psi}$  of a Riesz basis is also a Riesz basis (Christensen, 2016, Section 5).

#### 3.1.2 ENCODER AND DECODER

Throughout the rest of this paper we fix frames and their duals on  $\mathcal{X}, \mathcal{Y}$  and denote them by

$$\Psi_{\mathcal{X}} = (\psi_j)_{j \in \mathbb{N}}, \quad \tilde{\Psi}_{\mathcal{X}} = (\tilde{\psi}_j)_{j \in \mathbb{N}}, \quad \Psi_{\mathcal{Y}} = (\eta_j)_{j \in \mathbb{N}}, \quad \tilde{\Psi}_{\mathcal{Y}} = (\tilde{\eta}_j)_{j \in \mathbb{N}}.$$

The corresponding analysis operators are  $T_{\Psi_{\mathcal{X}}}, T_{\tilde{\Psi}_{\mathcal{X}}}, T_{\Psi_{\mathcal{Y}}}$  and  $T_{\tilde{\Psi}_{\mathcal{Y}}}$ . We then introduce encoder and decoder maps via

$$\mathcal{E}_{\mathcal{X}} := T_{\Psi_{\mathcal{X}}} = \begin{cases} \mathcal{X} \rightarrow \ell^2(\mathbb{N}), \\ x \mapsto (\langle x, \psi_j \rangle_{\mathcal{X}})_{j \in \mathbb{N}}, \end{cases} \quad \mathcal{D}_{\mathcal{Y}} := T'_{\Psi_{\mathcal{Y}}} = \begin{cases} \ell^2(\mathbb{N}) \rightarrow \mathcal{Y}, \\ (y_j)_{j \in \mathbb{N}} \mapsto \sum_{j \in \mathbb{N}} y_j \eta_j. \end{cases} \quad (3.3)$$

In case  $\Psi_{\mathcal{X}}$  and  $\Psi_{\mathcal{Y}}$  are Riesz bases, the mappings in (3.3) are boundedly invertible.

### 3.1.3 SMOOTHNESS SCALES

The encoder mapping  $\mathcal{E}_X : \mathcal{X} \rightarrow \ell^2(\mathbb{N})$  maps an element to its coefficients in the frame representation. For computational purposes, this coefficient sequence must be truncated, as only finitely many coefficients can be considered. Consequently, it is essential to control the error resulting from discarding higher-order frame coefficients. To formalize this we next introduce scales of subspaces of  $\mathcal{X}$ ,  $\mathcal{Y}$ , whose elements exhibit a certain coefficient decay.

**Definition 14** *Let  $\theta = (\theta_j)_{j \in \mathbb{N}}$  be a strictly positive, monotonically decreasing sequence such that  $\theta^{1+\varepsilon} \in \ell^1(\mathbb{N})$  for all  $\varepsilon > 0$ .*

*For all  $r, t \geq 0$  we introduce the subspaces  $\mathcal{X}_r \subset \mathcal{X}$  and  $\mathcal{Y}_t \subset \mathcal{Y}$  via*

$$\mathcal{X}_r := \{x \in \mathcal{X} : \|x\|_{\mathcal{X}_r} < \infty\} \quad \text{and} \quad \mathcal{Y}_t := \{y \in \mathcal{Y} : \|y\|_{\mathcal{Y}_t} < \infty\}$$

where

$$\|x\|_{\mathcal{X}_r}^2 := \sum_{j \in \mathbb{N}} \langle x, \tilde{\psi}_j \rangle_{\tilde{\mathcal{X}}}^2 \theta_j^{-2r} \quad \text{and} \quad \|y\|_{\mathcal{Y}_t}^2 := \sum_{j \in \mathbb{N}} \langle y, \tilde{\eta}_j \rangle_{\tilde{\mathcal{Y}}}^2 \theta_j^{-2t}.$$

For every  $r \geq 0$ ,  $\mathcal{X}_r$  is a Hilbert space (Herrmann et al., 2024, Lemma 1).

## 3.2 FrameNet

In this subsection we recall the FrameNet architecture from Herrmann et al. (2024, Section 2). We start by formally introducing feedforward neural networks (NNs) following Opschoor et al. (2022a); Herrmann et al. (2024).

### 3.2.1 FEEDFORWARD NEURAL NETWORKS

**Definition 15** *A function  $f : \mathbb{R}^{p_0} \rightarrow \mathbb{R}^{p_{L+1}}$  is called a neural network, if there exist  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ , integers  $p_1, \dots, p_{L+1} \in \mathbb{N}$ ,  $L \in \mathbb{N}$  and real numbers  $w_{i,j}^l, b_j^l \in \mathbb{R}$  such that for all  $x = (x_i)_{i=1}^{p_0} \in \mathbb{R}^{p_0}$*

$$z_j^1 = \sigma \left( \sum_{i=1}^{p_0} w_{i,j}^1 x_i + b_j^1 \right), \quad j = 1, \dots, p_1, \quad (3.4a)$$

$$z_j^{l+1} = \sigma \left( \sum_{i=1}^{p_l} w_{i,j}^{l+1} z_i^l + b_j^{l+1} \right), \quad l = 1, \dots, L-1, \quad j = 1, \dots, p_{l+1},$$

$$f(x) = (z_j^{L+1})_{j=1}^{p_{L+1}} = \left( \sum_{i=1}^{p_L} w_{i,j}^{L+1} z_i^L + b_j^{L+1} \right)_{j=1}^{p_{L+1}}. \quad (3.4b)$$

We call  $\sigma$  the activation function,  $L$  the depth,  $p := \max_{l=0, \dots, L+1} p_l$  the width,  $w_{i,j}^l \in \mathbb{R}$  the weights, and  $b_j^l \in \mathbb{R}$  the biases of the NN.

While different NNs can realize the same function, for simplicity we refer to a function  $f : \mathbb{R}^{p_0} \rightarrow \mathbb{R}^{p_{L+1}}$  as an NN of type (3.4), if it allows for (at least) one such representation. Additionally, our analysis will require some further terminology: For a NN  $f : \mathbb{R}^{p_0} \rightarrow \mathbb{R}^{p_{L+1}}$  as in (3.4) its

- *size* is the number of nonzero parameters

$$\text{size}(f) := |\{(i, j, l) : w_{i,j}^l \neq 0\}| + |\{(j, l) : b_j^l \neq 0\}|,$$

- *maximum of parameters* is

$$\text{mpar}(f) := \max \left\{ \max_{i,j,l} |w_{i,j}^l|, \max_{j,l} |b_j^l| \right\},$$

- *maximum range* on  $\Omega \subseteq \mathbb{R}^{p_0}$  is

$$\text{mran}_\Omega(f) := \sup_{x \in \Omega} \|f(x)\|_2 = \sup_{x \in \Omega} \left( \sum_{j=1}^{p_{L+1}} (z_j^{L+1}(x))^2 \right)^{\frac{1}{2}}.$$

Throughout, for  $q \in \mathbb{N}$ ,  $q \geq 1$ , we consider the activation function

$$\sigma_q(x) := \max\{0, x\}^q, \quad x \in \mathbb{R}.$$

For  $q = 1$ ,  $\sigma_1$  is the *rectified linear unit* (ReLU), for  $q \geq 2$ ,  $\sigma_q$  is called *rectified power unit* (RePU).

**Remark 16** *Definition 15 introduces a NN as a function  $f : \mathbb{R}^{p_0} \rightarrow \mathbb{R}^{p_{L+1}}$ . Throughout, we also understand the realization of a NN as a map  $f : \ell^2(\mathbb{N}) \rightarrow \ell^2(\mathbb{N})$  via extension by zeros. This is equivalent to suitably padding the weight matrices  $(w_{i,j}^l)_{i,j}$  and bias vectors  $(b_j^l)_j$  in Definition 15 for  $l \in \{0, L+1\}$  with infinitely many zeros, see Herrmann et al. (2024, Remark 13).*

### 3.2.2 THE FRAME NET CLASS

By definition of frames, the encoding operator  $\mathcal{E}_X : \mathcal{X} \rightarrow \ell^2(\mathbb{N})$  in (3.3) is injective, and the decoding operator  $\mathcal{D}_Y : \ell^2(\mathbb{N}) \rightarrow \mathcal{Y}$  is surjective. Thus for every mapping  $G : \mathcal{X} \rightarrow \mathcal{Y}$ , there exists a coefficient map  $g : \ell^2(\mathbb{N}) \rightarrow \ell^2(\mathbb{N})$  such that

$$G = \mathcal{D}_Y \circ g \circ \mathcal{E}_X.$$

This motivates the introduction of the following function class.

Given  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ , positive integers  $L, p, s \in \mathbb{N}$ , and reals  $M, B \in \mathbb{R}$ , let

$$\begin{aligned} \mathbf{g}_{\text{FN}}(\sigma, L, p, s, M, B) &:= \{g : \mathbb{R}^{p_0} \rightarrow \mathbb{R}^{p_{L+1}} \text{ is NN with activation function } \sigma \text{ s.t.} \\ &\quad \text{depth}(g) \leq L, \text{ width}(g) \leq p, \text{ size}(g) \leq s, \\ &\quad \text{mpar}(g) \leq M, \text{ mran}_{[-1,1]^{p_0}}(g) \leq B, p_0, p_{L+1} \leq p\}. \end{aligned}$$

Instead of directly considering (3.1), for our analysis, contrary to Herrmann et al. (2024), with  $\theta$  from Definition 14, fixed  $R > 0$ , and  $U := [-1, 1]^{\mathbb{N}}$ , it will be convenient to introduce the linear scaling

$$S_r := \begin{cases} \times_{j \in \mathbb{N}} [-R\theta_j^r, R\theta_j^r] \rightarrow U \\ (x_j)_{j \in \mathbb{N}} \mapsto \left(\frac{x_j}{R\theta_j^r}\right)_{j \in \mathbb{N}}. \end{cases} \quad (3.5)$$

The FrameNet class then consists of all operators

$$\mathbf{G}_{\text{FN}}(\sigma, L, p, s, M, B) := \{G = \mathcal{D}_Y \circ g \circ S_r \circ \mathcal{E}_X : g \in \mathbf{g}_{\text{FN}}(\sigma, L, p, s, M, B)\}. \quad (3.6)$$

### 3.3 Approximation of Holomorphic Operators

Our statistical theory established in Section 2 shows that sample complexity depends on

- (i) the approximation quality of  $\mathbf{G}$  w.r.t  $G_0$ , cf. the term  $\inf_{G^* \in \mathbf{G}} \|G^* - G_0\|_{L^2(\gamma)}$  in (2.13),
- (ii) the metric entropy of  $\mathbf{G}$ , cf. the terms  $\delta_n^2$  and  $\tilde{\delta}_n^2$  in (2.13).

In this subsection we give results for both the approximation quality and the metric entropy of the FrameNet class  $\mathbf{G}_{\text{FN}}$ .

#### 3.3.1 SETTING

Let us start by making the assumptions on  $G_0$  and the sampling measure  $\gamma$  on  $\mathcal{X}$  more precise. Denote in the following  $U := [-1, 1]^{\mathbb{N}}$ , let  $r > \frac{1}{2}$ ,  $R > 0$ , and let the frames  $(\psi_j)_{j \in \mathbb{N}}$ ,  $(\eta_j)_{j \in \mathbb{N}}$  be as in Section 3, and  $(\theta_j)_{j \in \mathbb{N}}$  as in Definition 14. Then

$$\sigma_R^r := \begin{cases} U \rightarrow \mathcal{X}, \\ \mathbf{y} \mapsto R \sum_{j \in \mathbb{N}} \theta_j^r y_j \psi_j \end{cases} \quad (3.7)$$

yields a well-defined map, since by construction  $(y_j \theta_j^r)_{j \in \mathbb{N}} \in \ell^2(\mathbb{N})$  for every  $\mathbf{y} \in U$ , and  $(\psi_j)_{j \in \mathbb{N}}$  is a frame. Next, we introduce the ‘‘cubes’’

$$C_R^r(\mathcal{X}) = \left\{ a \in \mathcal{X} : \sup_{j \in \mathbb{N}} \theta_j^{-r} |\langle a, \tilde{\psi}_j \rangle_{\mathcal{X}}| \leq R \right\}. \quad (3.8)$$

Observe that with  $\sigma_R^r(U) := \{\sigma_R^r(\mathbf{y}) : \mathbf{y} \in U\}$ , clearly

$$C_R^r(\mathcal{X}) \subseteq \sigma_R^r(U),$$

but equality holds in general only if the  $(\psi_j)_{j \in \mathbb{N}}$  form a Riesz basis (Herrmann et al., 2024, Remark 10).

**Example 2** Let  $\mathcal{X} = \mathbb{R}$ ,  $r = 1$ ,  $R = 1$ ,  $\theta_1 = 3/2$  and  $\theta_2 = 1/2$ . Consider the frame  $\Psi = \{1, 1\}$  (which is not a basis) with dual analogue  $\tilde{\Psi} = \{1/2, 1/2\}$ . Then with  $U = [-1, 1]^2$

$$C_1^1(\mathbb{R}) = [-1, 1] \subsetneq [-2, 2] = \{\theta_1 y_1 \psi_1 + \theta_2 y_2 \psi_2 : \mathbf{y} \in U\}.$$

The holomorphy assumption on  $G_0$  can now be formulated as follows (Herrmann et al., 2024, Assumption 1). Recall that for a real Hilbert space  $\mathcal{X}$ , we denote by  $\mathcal{X}_{\mathbb{C}}$  its complexification.

**Assumption 2** For some  $r > 1$ ,  $R > 0$ ,  $t > 0$ ,  $C_{G_0} < \infty$  there exists an open set  $O_{\mathbb{C}} \subset \mathcal{X}_{\mathbb{C}}$  containing  $\sigma_R^r(U)$ , such that

- (a)  $\sup_{a \in O_{\mathbb{C}}} \|G_0(a)\|_{(\mathcal{Y}_{\mathbb{C}})_t} \leq C_{G_0}$ , and  $G_0 : O_{\mathbb{C}} \rightarrow \mathcal{Y}_{\mathbb{C}}$  is holomorphic,
- (b)  $\gamma$  is a probability measure on  $\mathcal{X}$  with  $\text{supp}(\gamma) \subseteq C_R^r(\mathcal{X})$ .

The assumption requires  $G_0$  to be holomorphic on a superset of  $\sigma_R^r(U)$ . The probability measure  $\gamma$  is allowed to have support on  $C_R^r(\mathcal{X})$  which is potentially smaller than  $\sigma_R^r(U)$ . In particular,  $G_0$  is then holomorphic on a superset of the support of  $\gamma$ .

**Example 3** Denote by  $\lambda$  the Lebesgue measure. Then  $\pi := \otimes_{j \in \mathbb{N}} \frac{\lambda}{2}$  is the uniform probability measure on  $U = [-1, 1]^{\mathbb{N}}$ , and

$$\gamma := (\sigma_R^r)_\# \pi \quad (3.9)$$

defines a probability measure on  $\mathcal{X}$  with support  $\sigma_R^r(U)$ . If  $\Psi_{\mathcal{X}}$  is a Riesz basis, then  $\text{supp}(\gamma) = C_R^r(\mathcal{X})$ , so that  $\text{supp}(\gamma) \subseteq C_R^r(\mathcal{X})$  as required in Assumption 2.

We emphasize that Assumption 2 does not require  $\Psi_{\mathcal{X}}$  to be a Riesz basis, and all results below remain valid for general frames. However, the Riesz basis property yields sharper rates when combined with the specific choice of  $\gamma$  from Example 3.

### 3.3.2 APPROXIMATION THEORY

Theorem 1 in Herrmann et al. (2024) establishes a convergence rate for approximating  $G_0$  within  $\mathbf{G}_{\text{FN}}$  (with unbounded weights), in terms of the number of trainable network parameters. Our main results on sample complexity depend on bounds on the metric entropy, which require bounded network weights. To address this, we now extend Herrmann et al. (2024, Theorem 1) to FrameNet architectures with bounded weights, as introduced in Section 3.

Similar to Schwab and Zech (2019); Opschoor et al. (2022b); Herrmann et al. (2024), the analysis builds on polynomial chaos expansions (Xiu and Karniadakis, 2002). Denote by  $(L_j)_{j \in \mathbb{N}}$  the univariate Legendre polynomials normalized such that  $\frac{1}{2} \int_{-1}^1 L_j(x)^2 dx = 1$  for all  $j \in \mathbb{N}$ . Then (see Abramowitz and Stegun 1948, Chapter 22)

$$\sup_{x \in [-1, 1]} |L_j(x)| \leq \sqrt{2j+1} \quad \forall j \in \mathbb{N}_0. \quad (3.10)$$

Next, let  $\mathcal{F}$  be the set of infinite-dimensional multiindices with finite support, i.e.

$$\mathcal{F} := \left\{ (\nu_j)_{j \in \mathbb{N}_0} \in \mathbb{N}_0^{\mathbb{N}} : |\text{supp } \nu| < \infty \right\}, \quad (3.11)$$

where  $\text{supp } \nu := \{j : \nu_j \neq 0\}$ . For finite sets of multiindices  $\Lambda \subseteq \mathcal{F}$ , their effective dimension and maximal order is defined as

$$d(\Lambda) := \sup\{|\text{supp } \nu| : \nu \in \Lambda\} \quad \text{and} \quad m(\Lambda) := \sup\{|\nu| : \nu \in \Lambda\},$$

where  $|\nu| := \sum_{j \in \mathbb{N}} \nu_j$ . Moreover, with  $U := [-1, 1]^{\mathbb{N}}$  for all  $\mathbf{y} \in U$  and  $\nu \in \mathcal{F}$ , we let  $L_\nu(\mathbf{y}) := \prod_{j \in \mathbb{N}} L_{\nu_j}(y_j)$  be the corresponding multivariate Legendre polynomial. This infinite product is well-defined, since for all but finitely many  $j$  holds  $\nu_j = 0$  so that  $L_{\nu_j} \equiv 1$ . The next Proposition gives an approximation result for multivariate Legendre polynomials. It is an extension of Opschoor et al. (2022a, Proposition 2.13) to the case of bounded network parameters. The proof is provided in Appendix D.1.

**Proposition 17 ( $\sigma_1$ -NN approximation of  $L_\nu$ )** Let  $\delta \in (0, 1/2)$  and  $\Lambda \subset \mathcal{F}$  be finite. Then there exists a  $\sigma_1$ -NN  $f_{\Lambda, \delta}$ , such that its outputs  $\{\tilde{L}_{\nu, \delta}\}_{\nu \in \Lambda}$  satisfy

$$\forall \nu \in \Lambda : \sup_{\mathbf{y} \in U} |L_\nu(\mathbf{y}) - \tilde{L}_{\nu, \delta}(\mathbf{y})| \leq \delta. \quad (3.12)$$

Furthermore, there exists a constant  $C > 0$  independent of  $\Lambda$ ,  $d(\Lambda)$ ,  $m(\Lambda)$  and  $\delta$  such that

$$\text{depth}(f_{\Lambda,\delta}) \leq C \left[ \log(d(\Lambda))d(\Lambda) \log(m(\Lambda))m(\Lambda) + m(\Lambda)^2 + \log(\delta^{-1})(\log(d(\Lambda)) + m(\Lambda)) \right]$$

$$\text{width}(f_{\Lambda,\delta}) \leq C|\Lambda|d(\Lambda)$$

$$\text{size}(f_{\Lambda,\delta}) \leq C \left[ |\Lambda|d(\Lambda)^2 \log(m(\Lambda)) + m(\Lambda)^3 d(\Lambda)^2 + \log(\delta^{-1})(m(\Lambda)^2 + |\Lambda|)d(\Lambda) \right]$$

$$\text{mpar}(f_{\Lambda,\delta}) \leq 1.$$

A similar result holds for RePU neural networks, see Proposition 51. Given  $q \in \mathbb{N}$ , and  $N \in \mathbb{N}$ , we introduce the two FrameNet classes

$$\begin{aligned} \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N) &:= \mathbf{G}_{\text{FN}}(\sigma_q, \text{width}_N, \text{depth}_N, \text{size}_N, M, B), \\ \mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_q, N) &:= \mathbf{G}_{\text{FN}}(\sigma_q, \text{width}_N, \text{depth}_N, \infty, M, B) \end{aligned} \quad (3.14a)$$

where for certain constants  $C_L, C_p, C_s, M, B \geq 1$  (to be determined later)

$$\text{depth}_N = \max\{1, \lceil C_L \log(N) \rceil\}, \quad \text{width}_N = \lceil C_p N \rceil, \quad \text{size}_N = \lceil C_s N \rceil, \quad (3.14b)$$

and  $\lceil x \rceil$  denotes the smallest integer larger or equal to  $x \in \mathbb{R}$ .

We emphasize that  $\mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$  corresponds to a sparsely connected architecture, whereas  $\mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_q, N)$  represents a fully connected architecture (because there is no constraint on its size). In particular, since every linear transformation in between activation functions has at most  $\text{width}_N + \text{width}_N^2$  parameters, we have

$$\text{size}(\mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_q, N)) \leq (\text{depth}_N + 1)(\text{width}_N + \text{width}_N^2) = O(\log(N)N^2) \quad \text{as } N \rightarrow \infty. \quad (3.15)$$

Thus  $\mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_q, N)$  can essentially be quadratically larger than  $\mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$ .

The proof of Theorem 18 below, given in Appendix D.3, crucially uses Proposition 17 to construct FrameNets that approximate the multivariate Legendre polynomials corresponding to the 'most relevant' Legendre coefficients of  $G_0$ . Hereby Assumption 2 guarantees sufficient decay of the Legendre coefficients by using that algebraic decay in the input space, see Assumption 2 (b), is inherited by the Legendre coefficients of the output (see Theorem 54 and 55) if the respective operator is holomorphic, see Assumption 2 (a).

Theorem 18 is an extension of Herrmann et al. (2024, Theorems 1 and 2) to the case of bounded network parameters.

**Theorem 18 (Sparse network approximation)** *Let  $G_0, \gamma$  satisfy Assumption 2 with  $r > 1, t > 0$ . Let  $q \geq 1$  be an integer and fix  $\tau > 0$  (arbitrarily small).*

(i) *There exists  $C > 0$  s.t. for all  $N \in \mathbb{N}$*

$$\inf_{G \in \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)} \left[ \|G - G_0\|_{\infty, \text{supp}(\gamma)}^2 \right] \leq CN^{-2 \min\{r-1, t\} + \tau}.$$

(ii) *Let  $\Psi_\chi$  be a Riesz basis and let  $\gamma$  be as in (3.9). Then there exists  $C > 0$  s.t. for all  $N \in \mathbb{N}$*

$$\inf_{G \in \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)} \left[ \|G - G_0\|_{L^2(\gamma)}^2 \right] \leq CN^{-2 \min\{r-\frac{1}{2}, t\} + \tau}.$$

**Remark 19** *Theorem 18 provides an approximation rate for  $\mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$  in terms of  $N$ . Since  $\mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_q, N) \supseteq \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$ , trivially the statement remains true for  $\mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_q, N)$ . However, as the size of  $\mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_q, N)$  can be quadratically larger by (3.15), the convergence rate in terms of network size is essentially halved for the fully connected architecture.*

### 3.3.3 ENTROPY BOUNDS FOR FRAME NET

In the following we bound the metric entropy (see 1.3) of FrameNets for ReLU activation functions. Recall that  $\Lambda_{\Psi_y}$  is the frame constant in (3.2). The proof of Lemma 20 is given in Appendix D.4.

**Lemma 20** (see **Schmidt-Hieber 2020a, Lemma 5**) *Let  $L, p, s, M, B \geq 1$ ,  $\sigma_R^r$  be as in (3.7) and  $U = [-1, 1]^N$ . Then  $\mathbf{G}_{\text{FN}} = \mathbf{G}_{\text{FN}}(\sigma_1, L, p, s, M, B)$  is compact with respect to  $\|\cdot\|_{\infty, \sigma_R^r(U)}$  and  $\|\cdot\|_n$ . Furthermore,  $\mathbf{G}_{\text{FN}}$  satisfies for all  $\delta > 0$*

$$H(\mathbf{G}_{\text{FN}}, \|\cdot\|_{\infty, \sigma_R^r(U)}, \delta) \leq (s+1) \log \left( 2^{L+6} \Lambda_{\Psi_y} L^2 M^{L+1} p^{L+4} \max\{1, \delta^{-1}\} \right). \quad (3.16)$$

*In particular there exist constants  $C_H^{\text{SP}}, C_H^{\text{FC}} > 0$  such that for the sparse and fully-connected FrameNet classes from (3.14) it holds*

$$\begin{aligned} H(\mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_1, N), \|\cdot\|_{\infty, \sigma_R^r(U)}, \delta) &\leq C_H^{\text{SP}} N (1 + \log(N)^2 + \log(\max\{1, \delta^{-1}\})) \\ H(\mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_1, N), \|\cdot\|_{\infty, \sigma_R^r(U)}, \delta) &\leq C_H^{\text{FC}} N^2 (1 + \log(N)^3 + \log(\max\{1, \delta^{-1}\})) \end{aligned}$$

for all  $\delta > 0$ .

**Remark 21** *The entropy bound is independent of the constant  $B$  bounding the maximum range of the network, see Subsection 3.2. However,  $B < \infty$  will be necessary to apply Theorem 6.*

For RePU activation, as mentioned before, the metric entropy bounds exhibit a worse dependency on the network parameters, due to the lack of global Lipschitz continuity of  $\sigma_q$  if  $q \geq 2$ . The proof of Lemma 22 is given in Appendix D.5.

**Lemma 22** *Let  $q \in \mathbb{N}$ ,  $q \geq 2$ , and  $L, p, s, M, B \geq 1$ ,  $\sigma_R^r$  be as in (3.7) and  $U = [-1, 1]^N$ . Then  $\mathbf{G}_{\text{FN}} = \mathbf{G}_{\text{FN}}(\sigma_q, L, p, s, M, B)$  is compact with respect to  $\|\cdot\|_{\infty, \sigma_R^r(U)}$  and  $\|\cdot\|_n$ . Furthermore,  $\mathbf{G}_{\text{FN}}$  satisfies for all  $\delta > 0$*

$$H(\mathbf{G}_{\text{FN}}, \|\cdot\|_{\infty, \sigma_R^r(U)}, \delta) \leq (s+1) \log \left( \Lambda_{\Psi_y} L q^{L+q} (2pM)^{4q^{2L+2}} \max\{1, \delta^{-1}\} \right). \quad (3.17)$$

*Consider the constant  $C_L$  from (3.14). Then there exists  $C_H^{\text{SP}}, C_H^{\text{FC}} > 0$  such that*

$$\begin{aligned} H(\mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N), \|\cdot\|_{\infty, \sigma_R^r(U)}, \delta) &\leq C_H^{\text{SP}} N^{1+2C_L \log(q)} (1 + \log(N)^2 + \log(\max\{1, \delta^{-1}\})) \\ H(\mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_q, N), \|\cdot\|_{\infty, \sigma_R^r(U)}, \delta) &\leq C_H^{\text{FC}} N^{2+2C_L \log(q)} (1 + \log(N)^2 + \log(\max\{1, \delta^{-1}\})) \end{aligned}$$

for all  $\delta > 0$ .

### 3.4 Statistical Theory

Using our statistical results from Theorem 6 and Corollaries 7-8 as well as the approximation result from Theorem 18, we can bound  $\mathbb{E}_{G_0}[\|\hat{G}_n - G_0\|_{L^2(\gamma)}^2]$  solely in terms of the statistical sample size  $n$ . This is formalized in the next two theorems, for ReLU and RePU activation.

Our result distinguishes between the sparse and fully connected architectures  $\mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$  and  $\mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_q, N)$  in (3.14). In practice, fully connected architectures are often preferred, due to their simpler implementation, and because training sparse NN architectures can run into problems like “bad” local minima, see for example Evci et al. (2019). Our theoretical upper bounds are sharper for sparse architectures; this is because the additional free parameters in the fully connected architecture increase the entropy of this class, but do not yield better approximation properties in our proofs.

**Theorem 23 (ERM for white noise and ReLU)** *Let  $G_0 : \mathcal{X} \rightarrow \mathcal{Y}$  and  $\gamma$  satisfy Assumption 2 for some  $r > 1$ ,  $t > 0$ . Fix  $\tau > 0$  (arbitrarily small) and set (cp. Example 3)*

$$\kappa := \begin{cases} 2 \min\{r - \frac{1}{2}, t\} & \text{if } \Psi_{\mathcal{X}} \text{ is a Riesz basis and } \gamma = (\sigma_R^r)_{\#}\pi, \\ 2 \min\{r - 1, t\} & \text{otherwise.} \end{cases}$$

For every  $n \in \mathbb{N}$ , let  $(x_i, y_i)_{i=1}^n$  be data generated by (2.1) either in the white noise model or in the sub-Gaussian noise model. Then there exist constants  $C_L, C_p, C_s, M, B \geq 1$  in (3.14) and  $C > 0$  (all independent of  $n$ ) such that

(i) **Sparse FrameNet:** with  $N = N(n) = \lceil n^{\frac{1}{\kappa+1}} \rceil$  and  $\mathbf{G} = \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_1, N)$ , there exists a measurable choice of an ERM  $\hat{G}_n$  of (2.3). Any such  $\hat{G}_n$  satisfies

$$\mathbb{E}_{G_0}[\|\hat{G}_n - G_0\|_{L^2(\gamma)}^2] \leq C n^{-\frac{\kappa}{\kappa+1} + \tau}, \quad (3.18)$$

(ii) **Fully connected FrameNet:** with  $N = N(n) = \lceil n^{\frac{1}{\kappa+2}} \rceil$  and  $\mathbf{G} = \mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_1, N)$ , there exists a measurable choice of an ERM  $\hat{G}_n$  of (2.3). Any such  $\hat{G}_n$  satisfies

$$\mathbb{E}_{G_0}[\|\hat{G}_n - G_0\|_{L^2(\gamma)}^2] \leq C n^{-\frac{\kappa}{\kappa+2} + \tau}. \quad (3.19)$$

**Proof** The proof follows directly from the entropy bounds in Lemma 20 and Corollary 8 (ii): First, Lemma 20 in particular verifies Assumption 1 for all  $G^* \in \mathbf{G}$ , which is required for Corollary 8. Applying the corollary with  $\beta = \kappa + \tau$  and absorbing the logarithmic term into  $\tau$  then gives (3.18). With  $\beta = \kappa/2 + \tau$  we obtain (3.19). Note that crucially  $M_{\mathbf{F}}$ , see (2.9), does not depend on  $N$ , because  $\|G_0\|_{\infty, \text{supp } \gamma} < \infty$  and  $\mathbf{G}_{\text{FN}}$  is universally bounded by the  $N$ -independent constant  $B$ , see (3.6).  $\blacksquare$

**Remark 24** *Consider the sparse FrameNet class from above. Similar to the proof of Corollary 8 (ii), one can show that the sequences*

$$(\delta_n)_{n \in \mathbb{N}}, \quad (\tilde{\delta}_n)_{n \in \mathbb{N}}, \quad \delta_n = \tilde{\delta}_n = N \log(N)^2 (1 + \log(n)) / n$$

satisfy the entropy conditions (2.10) for the sparse FrameNet class  $\mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_1, N)$ . By choosing  $N = \lceil n^{1/(\kappa+1)} \rceil$  to balance the approximation term in (2.11), Theorem 6 gives the following high probability bound for  $\|\hat{G}_n - G_0\|_{L^2(\gamma)}$ : For all  $\tau > 0$  there exists a constant  $C > 0$  such that for all  $n \in \mathbb{N}$  and all  $R_n \geq Cn^{-\kappa/(\kappa+1)+\tau}$  it holds

$$\mathbb{P}_{G_0} \left( \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq R_n \right) \leq 2 \exp \left( -\frac{nR_n^2}{C^2 (\sigma^2 + M_{\mathbf{F}}^2)} \right).$$

Similar high-probability bounds hold for the fully connected FrameNet class and the RePU activation function.

**Theorem 25 (ERM for white noise and RePU)** Consider the setting of Theorem 23 and let  $q \in \mathbb{N}$ ,  $q \geq 2$ . There exist constants  $C_L, C_p, C_s, M, B \geq 1$  in (3.14) and  $C > 0$  (all independent of  $N$ ) such that

- (i) **Sparse FrameNet:** with  $N = N(n) = \lceil n^{\frac{1}{\kappa+1+4C_L \log(q)}} \rceil$  and  $\mathbf{G} = \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$ , there exists a measurable choice of an ERM  $\hat{G}_n$  of (2.3). Any such  $\hat{G}_n$  satisfies

$$\mathbb{E}_{G_0} [\|\hat{G}_n - G_0\|_{L^2(\gamma)}^2] \leq Cn^{-\frac{\kappa}{\kappa+1+4C_L \log(q)} + \tau}. \quad (3.20)$$

- (ii) **Fully connected FrameNet:** with  $N = N(n) = \lceil n^{\frac{1}{\kappa+2+4C_L \log(q)}} \rceil$ ,  $\mathbf{G} = \mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_q, N)$ , there exists a measurable choice of ERM  $\hat{G}_n$  of (2.3). Any such  $\hat{G}_n$  satisfies

$$\mathbb{E}_{G_0} [\|\hat{G}_n - G_0\|_{L^2(\gamma)}^2] \leq Cn^{-\frac{\kappa}{\kappa+2+4C_L \log(q)} + \tau}. \quad (3.21)$$

**Proof** Similar to the proof of Theorem 23, the proof of Theorem 25 follows from the entropy bounds in Lemma 22 and Corollary 8 (ii) with  $\beta = \kappa/(1+2C_L \log(q)) + \tau$  for (3.20) and  $\beta = \kappa/(2+2C_L \log(q)) + \tau/2$  for (3.21). Again the logarithmic term is absorbed into  $\tau$ .  $\blacksquare$

A few remarks are in order. In the RePU case the activation function  $\sigma_q(x) = \max\{0, x\}^q$  has no global Lipschitz condition for  $q \geq 2$ . As a result, the entropy bounds obtained for the corresponding FrameNet class are larger than for ReLU. This leads to worse convergence rates. Moreover, for the RePU case, the convergence rate depends on the constant  $C_L$  in (3.14b). The proof of Theorem 18 shows that  $C_L$  depends on the decay properties of the Legendre coefficients  $(c_{\nu_{i,j}})_{i,j \in \mathbb{N}}$  of the function  $G_0 \circ \sigma_R^r$ , i.e.  $C_L$  depends on  $G_0$  (see (D.23) and (D.24) in Theorem 54). Explicit bounds on  $C_L$  are possible, see Zech (2018, Lemma 1.4.15).

## 4. Applications

We now present two applications of our results. First, in finite dimensional regression, our analysis recovers well-known minimax-optimal rates for standard smoothness classes. This indicates that our main statistical result in Theorem 6 is *in general optimal*. However, we do not claim optimality specifically for the approximation of holomorphic operators discussed in Section 3. Second, for an infinite dimensional problem we address the learning of a solution operator to a parameter dependent PDE.

### 4.1 Finite Dimensional Regression

Let  $d \in \mathbb{N}$  and  $\mathcal{X} = \mathbb{R}^d$ ,  $\mathcal{Y} = \mathbb{R}$ . Moreover, let  $D \subseteq \mathbb{R}^d$  be a bounded, open, smooth domain, and  $G_0 : D \rightarrow \mathbb{R}$  a ground-truth regression function. Suppose that

$$x_i \stackrel{\text{iid}}{\sim} \gamma \quad \text{and} \quad \varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1) \quad \forall i \in \mathbb{N}$$

are independent samples for some probability measure  $\gamma$  on  $D$ , and let

$$y_i = G_0(x_i) + \varepsilon_i \quad \forall i \in \mathbb{N}.$$

Given a regression class  $\mathbf{G}$  of measurable mappings from  $D \rightarrow \mathbb{R}$ , the least-squares problem is to determine

$$\hat{G}_n \in \arg \min_{G \in \mathbf{G}} \sum_{i=1}^n -2G(x_i)y_i + G(x_i)^2 = \arg \min_{G \in \mathbf{G}} \sum_{i=1}^n |G(x_i) - y_i|^2. \quad (4.2)$$

For  $\mathfrak{s} \geq 0$ , it is well-known that the minimax-optimal rate of recovering a ground-truth function  $G_0$  in  $\mathfrak{s}$ -smooth function classes, such as the Sobolev space  $H^{\mathfrak{s}}(D)$  or the Besov space  $B_{\infty, \infty}^{\mathfrak{s}}(D)$  (see Edmunds and Triebel 1996 for definitions), equals  $n^{-\frac{2\mathfrak{s}}{2\mathfrak{s}+d}}$  (Giné and Nickl, 2016; van de Geer, 2000).

Denote now by  $\mathbf{G}_R^{\mathfrak{s}}$  the ball of radius  $R > 0$  around the origin in either  $H^{\mathfrak{s}}(D)$  or  $B_{\infty, \infty}^{\mathfrak{s}}(D)$ . Then,

$$H(\mathbf{G}_R^{\mathfrak{s}}, \|\cdot\|_{\infty, \text{supp}(\gamma)}, \delta) \simeq \left(\frac{R}{\delta}\right)^{d/\mathfrak{s}} \quad \forall \delta \in (0, 1),$$

which holds for all  $\mathfrak{s} > 0$  if  $\mathbf{G}_R^{\mathfrak{s}}$  is the ball in  $B_{\infty, \infty}^{\mathfrak{s}}(D)$ , and for all  $\mathfrak{s} > d/2$  in case  $\mathbf{G}_R^{\mathfrak{s}}$  is the ball in  $H^{\mathfrak{s}}(D)$ , see Triebel (1995, Theorem 4.10.3). Corollary 8 (i) (with  $\alpha = d/\mathfrak{s} < 2$ ) then directly yields the following theorem. It recovers the minimax optimal rate for nonparametric least squares/maximum likelihood estimators.

**Theorem 26** *Let  $R > 0$  and  $\mathfrak{s} > d/2$ . Then, there exists  $C > 0$  such that for all  $G_0 \in \mathbf{G}_R^{\mathfrak{s}}$ , the estimator  $\hat{G}_n$  in (4.2) with  $\mathbf{G} = \mathbf{G}_R^{\mathfrak{s}}$  and data as in (4.1) satisfies*

$$\mathbb{E}_{G_0} [\|\hat{G}_n - G_0\|_{L^2(D)}^2] \leq C n^{-\frac{2}{2+\alpha}} = C n^{-\frac{2\mathfrak{s}}{2\mathfrak{s}+d}} \quad \forall n \in \mathbb{N}.$$

### 4.2 Parametric Darcy Flow

As a second application we apply Theorem 23 to the solution operator of the diffusion equation, extending the discussion of approximation errors in Herrmann et al. (2024, Section 7.1).

#### 4.2.1 SETUP

We recall the setup from Herrmann et al. (2024, Sections 7.1.1, 7.1.2).

Let  $d \in \mathbb{N}$ , and denote by  $\mathbb{T}^d \simeq [0, 1]^d$  the  $d$ -dimensional torus. In the following, all function spaces on  $\mathbb{T}^d$  are understood to be one-periodic in each variable. Fix  $\bar{a} \in L^\infty(\mathbb{T}^d)$  and  $f \in H^{-1}(\mathbb{T}^d)/\mathbb{R}$  such that for some constant  $a_{\min} > 0$

$$\text{ess inf}_{x \in \mathbb{T}^d} (\bar{a}(x) + a(x)) > a_{\min}. \quad (4.3)$$

We consider the ground truth  $G_0 : a \mapsto u$ , mapping  $a \in L^\infty(\mathbb{T}^d)$  to the solution  $u \in H^1(\mathbb{T}^d)$  of

$$-\nabla \cdot ((\bar{a} + a)\nabla u) = f \text{ on } \mathbb{T}^d \quad \text{and} \quad \int_{\mathbb{T}^d} u(x) \, dx = 0. \quad (4.4)$$

Then  $G_0 : \{a \in L^\infty(\mathbb{T}^d) : (4.3) \text{ holds}\} \rightarrow H^1(\mathbb{T}^d)$  is well-defined.

To represent  $a$  and  $u$ , we use Fourier expansions on  $\mathbb{T}^d$ . Denote for  $j \in \mathbb{N}_0$  and  $\mathbf{j} \in \mathbb{N}_0^d$ ,  $d \geq 2$ ,

$$\xi_0 := 1, \quad \xi_{2j}(x) := \sqrt{2} \cos(2\pi jx), \quad \xi_{2j-1}(x) := \sqrt{2} \sin(2\pi jx), \quad \xi_{\mathbf{j}}(x_1, \dots, x_d) := \prod_{k=1}^d \xi_{j_k}(x_k).$$

Then for  $r \geq 0$ ,  $\{\max\{1, |\mathbf{j}|\}^r \xi_{\mathbf{j}} : \mathbf{j} \in \mathbb{N}_0^d\}$  forms an ONB of  $H^r(\mathbb{T}^d)$  equipped with inner product

$$\langle u, v \rangle_{H^r(\mathbb{T}^d)} := \sum_{\mathbf{j} \in \mathbb{N}_0^d} \langle u, \xi_{\mathbf{j}} \rangle_{L^2} \langle v, \xi_{\mathbf{j}} \rangle_{L^2} \max\{1, |\mathbf{j}|\}^{2r}.$$

In the following, fix  $r_0, t_0 \geq 0$  and set

$$\begin{aligned} \mathcal{X} &:= H^{r_0}(\mathbb{T}^d), & \psi_{\mathbf{j}} &:= \max\{1, |\mathbf{j}|\}^{-r_0} \xi_{\mathbf{j}}, \\ \mathcal{Y} &:= H^{t_0}(\mathbb{T}^d), & \eta_{\mathbf{j}} &:= \max\{1, |\mathbf{j}|\}^{-t_0} \xi_{\mathbf{j}}, \end{aligned} \quad (4.5)$$

so that  $\Psi_{\mathcal{X}} := (\psi_{\mathbf{j}})_{\mathbf{j} \in \mathbb{N}_0^d}$ ,  $\Psi_{\mathcal{Y}} := (\eta_{\mathbf{j}})_{\mathbf{j} \in \mathbb{N}_0^d}$  form ONBs of  $\mathcal{X}$ ,  $\mathcal{Y}$  respectively. The encoder  $\mathcal{E}_{\mathcal{X}}$  and decoder  $\mathcal{D}_{\mathcal{Y}}$  are now as in (3.3). Direct calculation shows  $\mathcal{X}^r = H^{r_0+rd}$  and  $\mathcal{Y}^t = H^{t_0+td}$  for  $r, t \geq 0$ ; for more details see Herrmann et al. (2024, Section 7.1.2).

#### 4.2.2 SAMPLE COMPLEXITY

We now analyze the sample complexity for learning the PDE solution operator  $G_0$  in Section 4.2.1. For a proof of Theorem 27, see Appendix E.1.

**Theorem 27** *Let  $d \in \mathbb{N}$ ,  $d \geq 2$ ,  $\mathfrak{s} > \frac{3d}{2}$  and  $t_0 \in [0, 1]$ . Fix  $\tau_1 > 0$ ,  $\tau_2 \in (0, \min\{\mathfrak{s} - \frac{3d}{2}, \frac{\tau_1 d}{8}\})$  (both arbitrarily small), and set*

$$r_0 = \begin{cases} \frac{d}{2} + \tau_2 & \text{if } \mathfrak{s} \in (\frac{3d}{2}, 2d + 1 - t_0] \\ \frac{\mathfrak{s} + t_0 - 1}{2} & \text{if } \mathfrak{s} > 2d + 1 - t_0. \end{cases}$$

Moreover let  $f \in C^\infty(\mathbb{T}^d)$ , and let

- (a) **ground truth:**  $G_0 : a \mapsto u$  be given through (4.4),
- (b) **representation system:**  $\mathcal{E}_{\mathcal{X}}$ ,  $\mathcal{D}_{\mathcal{Y}}$  be as in (3.3) with the orthonormal basis in (4.5), and  $r_0, t_0$  from above,
- (c) **data:**  $\gamma$  be the measure defined in (3.7) and (3.9) with  $r = \frac{\mathfrak{s} - r_0}{d}$  such that  $\bar{a} + a$  satisfies (4.3) for all  $a \in \text{supp}(\gamma)$ , and let  $(x_i, y_i)_{i=1}^n$  be generated by (2.1) with the additive white noise model or the sub-Gaussian noise model,

(d) **regression class:**  $\mathbf{G} = \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_1, N)$  be the  $n$ -dependent sparse FrameNet architecture in (3.14) with  $N(n) = \lceil n^{\frac{1}{\kappa+1}} \rceil$ .

Then there exists a constant  $C > 0$  such that for all  $n \in \mathbb{N}$  there exists a measurable ERM  $\hat{G}_n \in \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_1, N(n))$  in (2.3), and any such  $\hat{G}_n$  satisfies

$$\mathbb{E}_{G_0}[\|\hat{G}_n - G_0\|_{L^2(H^{r_0}(\mathbb{T}^d), \gamma; H^{t_0}(\mathbb{T}^d))}^2] \leq Cn^{-\frac{\kappa}{\kappa+1} + \tau_1} \quad (4.6)$$

where

$$\kappa = \begin{cases} 2 \min\{\frac{\mathfrak{s}}{d} - 1, \frac{1-t_0}{d}\} & \text{if } \mathfrak{s} \in (\frac{3d}{2}, 2d + 1 - t_0], \\ \frac{\mathfrak{s}+1-t_0}{d} - 1 & \text{if } \mathfrak{s} > 2d + 1 - t_0. \end{cases} \quad (4.7)$$

**Remark 28** Consider the setting of Theorem 27, and let  $\text{supp}(\gamma) \subseteq C_R^r(\mathcal{X})$ . A slight modification of the proof of Theorem 27 similar to Herrmann et al. (2024), using the approximation bound in Theorem 18 (i) instead of (ii), then yields

$$\mathbb{E}_{G_0}[\|\hat{G}_n - G_0\|_{L^2(H^{r_0}(\mathbb{T}^d), \gamma; H^{t_0}(\mathbb{T}^d))}^2] \leq Cn^{-\frac{\kappa}{\kappa+1} + \tau_1} \quad (4.8)$$

where for some (small)  $\tau_2 > 0$

$$(r_0, \kappa) = \begin{cases} (\frac{d}{2} + \tau_2, \frac{2\mathfrak{s}}{d} - 3) & \text{if } \mathfrak{s} \in (\frac{3d}{2}, \frac{3d}{2} + 1 - t_0] \\ (\frac{\mathfrak{s}+t_0-\frac{d}{2}-1}{2}, \frac{\mathfrak{s}+1-t_0}{d} - \frac{3}{2}) & \text{if } \mathfrak{s} > \frac{3d}{2} + 1 - t_0. \end{cases}$$

Since

$$\begin{aligned} B_R(H^s(\mathbb{T}^d)) &= B_R(\mathcal{X}^r) = \{x \in \mathcal{X} : \|x\|_{\mathcal{X}^r} \leq R\} = \left\{ x \in \mathcal{X} : \sum_{j \in \mathbb{N}} \langle x, \tilde{\psi}_j \rangle_{\mathcal{X}}^2 \theta_j^{-2r} \leq R^2 \right\} \\ &\subseteq \left\{ a \in \mathcal{X} : \sup_{j \in \mathbb{N}} \theta_j^{-r} |\langle a, \tilde{\psi}_j \rangle_{\mathcal{X}}| \leq R \right\} = C_R^r(\mathcal{X}), \end{aligned}$$

in particular, (4.8) holds for any  $\gamma$  with  $\text{supp}(\gamma) \subseteq B_R(H^s(\mathbb{T}^d))$ . This shows Theorem 1.

Similar rates can also be obtained for this PDE model on a convex, polygonal domain  $D \subset \mathbb{T}^2$  with Dirichlet boundary conditions. The argument uses the Riesz basis constructed in Davydov and Stevenson (2005), but is otherwise similar to the torus, for details see Herrmann et al. (2024, Section 7.2). Moreover, using the RePU activation function, (4.6) holds with convergence rate  $\kappa/(\kappa + 1 + 4C_L \log(q))$ , where  $\kappa$  is from (4.7) and  $C_L$  from (3.14b). The proof is similar to Theorem 27 using Theorem 25 instead of Theorem 23. Finally, rates for the fully-connected class  $\mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_q, N)$  can be established using Theorems 23 (ii) and 25 (ii).

## 5. Conclusions

In this work, we established convergence theorems for empirical risk minimization to learn mappings  $G_0$  between infinite dimensional Hilbert spaces. Our setting assumes given data

in the form of  $n$  input-output pairs, with an additive noise model. We discuss both the case of Gaussian white noise and sub-Gaussian noise. Our main statistical result, Theorem 6, bounds the mean-squared  $L^2$ -error  $\mathbb{E}[\|\hat{G}_n - G_0\|_{L^2(\gamma)}^2]$  in terms of the approximation error, and algebraic rates in  $n$  depending only on the metric entropy of  $\mathbf{G}$ . This provides a general framework to study operator learning from the perspective of sample complexity.

In the second part of this work, we applied our statistical results to a specific operator learning architecture from Herrmann et al. (2024), termed FrameNet. As our main application, we showed that holomorphic operators  $G_0$  can be learned with NN-based surrogates without suffering from the curse of dimensionality, cf. Theorem 23. Such results have wide applicability, as the required holomorphy assumption is well-established in the literature, and has been verified for a variety of models including for example general elliptic PDEs (Cohen et al., 2010, 2011; Cohen and DeVore, 2015; Harbrecht et al., 2016), Maxwell’s equations (Jerez-Hanckes et al., 2017), the Calderon projector (Henríquez and Schwab, 2021), the Helmholtz equation (Hiptmair et al., 2018; Spence and Wunsch, 2023) and also nonlinear PDEs such as the Navier-Stokes equations (Cohen et al., 2018).

The main observation model considered in this work is white noise, which can be understood as an idealization of taking (equally spaced, noisy) point evaluations, see Remark 4. Our measurement model hence represents the process of gathering real-world data from physical systems rather than data generated from approximate solutions of PDEs, in which case the numerical error can be made arbitrarily small. Covering noise-free data with our statistical framework is left for future research.

Extending our results from Hilbert spaces to Banach spaces seems non-trivial. The white noise model is not well-defined in Banach spaces and would need to be replaced. Moreover, bounding the empirical risk in Theorem 5 uses the Hilbert space structure in an essential way: both the basic inequality (Lemma 30) and the chaining Lemma (Lemma 32) rely on the presence of an inner product and would need to be adapted in a Banach space setting.

We have also not discussed the training procedure for *finding* an empirical risk minimizer. The underlying optimization problem is in general non-convex and hence algorithms can typically only approximate the global minimizer up to some tolerance. However, we believe that large parts of our analysis extend straightforwardly to such approximate minimizers with the tolerance appearing as an additional additive term in the mean squared error bound. Full details are left for future work.

## Acknowledgments

SW and JZ gratefully acknowledge support from the German Research Foundation (DFG) within the Priority Programme SPP 2298, *Theoretical Foundations of Deep Learning* (project number 543965776).

## Appendix A. Auxiliary Probabilistic Lemmas

We recall the classical Bernstein inequality.

**Lemma 29 (Bernstein's inequality)** *Let  $X_1, \dots, X_n$  be independent, centered RVs with finite second moments  $\mathbb{E}[X_i^2] < \infty$  and uniform bound  $|X_i| \leq M$  for  $i = 1, \dots, n$ . Then it holds*

$$\mathbb{P} \left( \left| \sum_{i=1}^n X_i \right| \geq t \right) \leq 2 \exp \left( - \frac{t^2}{2 \sum_{i=1}^n \mathbb{E}[X_i^2] + \frac{2}{3} M t} \right), \quad t \geq 0.$$

For a proof of Bernstein's inequality, see Pollard (1984, page 193).

**Lemma 30 (Basic Inequality)** *Let  $\varepsilon_i, i = 1, \dots, n$  be i.i.d. white noise or sub-Gaussian noise. Then it holds for all  $G^* \in \mathbf{G}$*

$$\|\hat{G}_n - G_0\|_n^2 \leq \|G^* - G_0\|_n^2 + \frac{2\sigma}{n} \sum_{i=1}^n \langle \varepsilon_i, \hat{G}_n(x_i) - G^*(x_i) \rangle_{\mathfrak{Y}}.$$

**Proof** Let  $G^* \in \mathbf{G}$  be arbitrary. Using the definition of  $\hat{G}_n$  from (2.3), it holds

$$\begin{aligned} & \|\hat{G}_n - G_0\|_n^2 \\ &= \frac{1}{n} \sum_{i=1}^n \|\hat{G}_n(x_i)\|_{\mathfrak{Y}}^2 - 2 \langle G_0(x_i), \hat{G}_n(x_i) \rangle_{\mathfrak{Y}} + \|G_0(x_i)\|_{\mathfrak{Y}}^2 \\ &= \frac{1}{n} \sum_{i=1}^n \|\hat{G}_n(x_i)\|_{\mathfrak{Y}}^2 - 2 \langle G_0(x_i) + \sigma \varepsilon_i, \hat{G}_n(x_i) \rangle_{\mathfrak{Y}} + 2\sigma \langle \varepsilon_i, \hat{G}_n(x_i) \rangle_{\mathfrak{Y}} + \|G_0(x_i)\|_{\mathfrak{Y}}^2 \\ &\leq \frac{1}{n} \sum_{i=1}^n \|G^*(x_i)\|_{\mathfrak{Y}}^2 - 2 \langle G_0(x_i) + \sigma \varepsilon_i, G^*(x_i) \rangle_{\mathfrak{Y}} + 2\sigma \langle \varepsilon_i, \hat{G}_n(x_i) \rangle_{\mathfrak{Y}} + \|G_0(x_i)\|_{\mathfrak{Y}}^2 \\ &= \frac{1}{n} \sum_{i=1}^n \|G^*(x_i)\|_{\mathfrak{Y}}^2 - 2 \langle G_0(x_i), G^*(x_i) \rangle_{\mathfrak{Y}} + \|G_0(x_i)\|_{\mathfrak{Y}}^2 + 2\sigma \langle \varepsilon_i, \hat{G}_n(x_i) - G^*(x_i) \rangle_{\mathfrak{Y}} \\ &= \|G^* - G_0\|_n^2 + \frac{2\sigma}{n} \sum_{i=1}^n \langle \varepsilon_i, \hat{G}_n(x_i) - G^*(x_i) \rangle_{\mathfrak{Y}}, \end{aligned}$$

which shows the claim. ■

Next, we state a generic chaining result from Dirksen (2015, Theorem 3.2), originally derived for finite index sets, for countable index sets  $T$ . We restrict ourselves to the case of real-valued stochastic processes.

**Lemma 31** *Let  $T$  be a countable index set and  $d : T \times T \rightarrow [0, \infty)$  a pseudometric. Furthermore, let  $(X_t)_{t \in T}$  be an  $\mathbb{R}$ -valued stochastic process such that for some  $\alpha > 0$  and all  $s, t \in T$ ,*

$$\mathbb{P}(|X_t - X_s| \geq ud(t, s)) \leq 2 \exp(-u^\alpha), \quad u \geq 0. \tag{A.1}$$

Then, there exists a constant  $M > 0$  depending only on  $\alpha$ , such that for all  $t_0 \in T$  it holds that

$$\mathbb{P} \left( \sup_{t \in T} |X_t - X_{t_0}| \geq M \left( J_\alpha(T, d) + u \sup_{s, t \in T} d(s, t) \right) \right) \leq \exp \left( -\frac{u^\alpha}{\alpha} \right), \quad u \geq 1,$$

where  $J_\alpha$  denotes the metric entropy integral

$$J_\alpha(T, d) = \int_0^\infty (\log N(T, d, u))^\frac{1}{\alpha} du.$$

**Proof** Since  $T$  is countable, we can write  $T = \{t_j : j \in \mathbb{N}\}$ . Using this, we define  $T_n = \{t_j : j \leq n\}$  for  $n \in \mathbb{N}$ . Since  $T_n$  is finite (Dirksen, 2015, Theorem 3.2, Eq. 3.2 and its proof) gives  $\tilde{M} > 0$  s.t.

$$\left( \mathbb{E} \left[ \sup_{t \in T_n} |X_t - X_{t_0}|^p \right] \right)^\frac{1}{p} \leq \tilde{M} \left( J_\alpha(T_n, d) + \sup_{s, t \in T_n} d(s, t) p^\frac{1}{\alpha} \right) \quad (\text{A.2})$$

for all  $p \geq 1$ ,  $t_0 \in T$  and  $n \in \mathbb{N}$ . In (A.2), we used Dirksen (2015, Eq. (2.3)) to upper bound the  $\gamma_\alpha$ -functionals by the respective metric entropy integrals  $J_\alpha$ . The monotone convergence theorem shows (A.2) for  $T$  in the limit  $n \rightarrow \infty$ . Applying Dirksen (2015, Lemma A.1) then gives the claim with  $M = \exp(\alpha^{-1})\tilde{M}$ .  $\blacksquare$

We now use Lemma 31 to establish the following concentration bound, which is tailored towards the empirical processes appearing in our proofs, cf. Lemma 30. Note that this lemma can be viewed as a generalization of the key chaining Lemma 3.12 in Nickl and Wang (2024) to  $\mathcal{Y}$ -valued regression functions; the proof follows along the same lines.

**Lemma 32 (Chaining Lemma)** *Let  $\mathcal{X}, \mathcal{Y}$  be separable Hilbert spaces, and suppose  $\Theta$  is a (possibly uncountable) set parametrizing a class of maps*

$$\mathcal{H} = \{h_\theta : \mathcal{X} \rightarrow \mathcal{Y}, \theta \in \Theta\}.$$

Consider an empirical process of the form

$$Z_n(\theta) = \frac{1}{n} \sum_{i=1}^n \langle h_\theta(x_i), \varepsilon_i \rangle_{\mathcal{Y}},$$

where  $x_1, \dots, x_n \in \mathcal{X}$  are fixed elements and  $\varepsilon_1, \dots, \varepsilon_n$  are either **(i)** i.i.d. Gaussian white noise processes indexed by  $\mathcal{Y}$ , or **(ii)** i.i.d. sub-Gaussian random variables in  $\mathcal{Y}$  with parameter 1.

Recall the empirical seminorm  $\|\cdot\|_n$ . Suppose that

$$\sup_{\theta \in \Theta} \|h_\theta\|_n =: U < \infty, \quad (\text{A.3})$$

and define the metric entropy integral

$$J(\mathcal{H}, d_n) = \int_0^U \sqrt{\log N(\mathcal{H}, d_n, \tau)} d\tau, \quad d_n(\theta, \theta') = \|h_\theta - h_{\theta'}\|_n.$$

Let the space  $(\Theta, d_n)$  be separable. Then  $\sup_{\theta \in \Theta} Z_n(\theta)$  is measurable and there exists a universal constant  $C_{\text{Ch}} > 0$  such that for all  $\delta > 0$  with

$$\sqrt{n}\delta \geq C_{\text{Ch}}J(\mathcal{H}, d_n), \quad (\text{A.4})$$

it holds

$$\mathbb{P} \left( \sup_{\theta \in \Theta} |Z_n(\theta)| \geq \delta \right) \leq \exp \left( -\frac{8n\delta^2}{C_{\text{Ch}}^2 U^2} \right).$$

**Proof** In both cases, we will apply Dirksen (2015, Theorem 3.2), which we stated in Lemma 31.

**White noise case.** Let  $\theta, \theta' \in \Theta$  be arbitrary. Since  $\varepsilon_i, i = 1, \dots, n$  are independent white noise processes, we have  $Z_n(\theta) - Z_n(\theta') \sim \mathcal{N}(0, n^{-1}\|h_\theta - h_{\theta'}\|_n^2)$ , i.e. the increments of  $Z_n$  are normal (recall that  $\mathbf{x}$  is regarded as fixed here). Thus,

$$\mathbb{P} (|Z_n(\theta) - Z_n(\theta')| \geq t) \leq 2 \exp(-nt^2/(2\|h_\theta - h_{\theta'}\|_n^2)), \quad t \geq 0,$$

and

$$\begin{aligned} \mathbb{P} \left( |Z_n(\theta) - Z_n(\theta')| \geq \frac{\sqrt{2}td_n(\theta, \theta')}{\sqrt{n}} \right) &\leq 2 \exp \left( -\frac{d_n(\theta, \theta')^2 t^2}{\|h_\theta - h_{\theta'}\|_n^2} \right) \\ &= 2 \exp(-t^2), \quad t \geq 0, \end{aligned} \quad (\text{A.5})$$

which verifies the assumption (A.1) for  $\alpha = 2$  and  $\bar{d}_n := \sqrt{2}d_n/\sqrt{n}$ .

Eq. (A.5) shows that the process  $Z_n(\theta)$  is sub-Gaussian with respect to the pseudometric  $\bar{d}_n(\theta, \theta') = \sqrt{2}\|h_\theta - h_{\theta'}\|_n/\sqrt{n}$ . Therefore Giné and Nickl (2016, Theorem 2.3.7 (a)) yields that  $Z_n(\theta)$  is sample bounded and uniformly sample continuous. Since  $(\Theta, d_n)$  is separable, so is  $(\Theta, \bar{d}_n)$ . Thus it holds

$$\sup_{\theta \in \Theta} Z_n(\theta) = \sup_{\theta \in \Theta_0} Z_n(\theta) \quad \text{a.s.}, \quad (\text{A.6})$$

where  $\Theta_0 \subset \Theta$  denotes a countable, dense subset. The right hand side of (A.6) is measurable as a countable supremum. Therefore also the left hand side  $\sup_{\theta \in \Theta} Z_n(\theta)$  is measurable. Applying Lemma 31 (to the countable set  $\Theta_0$ ) and using (A.6) gives that for some universal constant  $M$  and all  $\theta_\dagger \in \Theta$ ,

$$\mathbb{P} \left( \sup_{\theta \in \Theta} |Z_n(\theta) - Z_n(\theta_\dagger)| \geq M \left( J(\mathcal{H}, \bar{d}_n) + \frac{tU}{\sqrt{n}} \right) \right) \leq \exp \left( -\frac{t^2}{2} \right), \quad t \geq 1.$$

Due to (A.3) it holds  $N(\mathcal{H}, d_n, \delta) = 1$  for all  $\delta \geq U$ , and thus

$$N(\mathcal{H}, \bar{d}_n, \tau) = N(\mathcal{H}, d_n, \sqrt{n}\tau/\sqrt{2}) = 1 \quad \forall \tau \geq \frac{\sqrt{2}U}{\sqrt{n}}.$$

Substituting  $\rho = \sqrt{n}\tau/\sqrt{2}$  we get

$$\begin{aligned} J(\mathcal{H}, \bar{d}_n) &= \int_0^{\sqrt{2}U/\sqrt{n}} \sqrt{\log N(\mathcal{H}, \bar{d}_n, \tau)} \, d\tau \\ &= \frac{\sqrt{2}}{\sqrt{n}} \int_0^U \sqrt{\log N\left(\mathcal{H}, \bar{d}_n, \frac{\sqrt{2}\rho}{\sqrt{n}}\right)} \, d\rho \\ &= \frac{\sqrt{2}}{\sqrt{n}} \int_0^U \sqrt{\log N(\mathcal{H}, d_n, \rho)} \, d\rho = \frac{\sqrt{2}J(\mathcal{H}, d_n)}{\sqrt{n}} \end{aligned}$$

and therefore

$$\mathbb{P}\left(\sup_{\theta \in \Theta} |Z_n(\theta) - Z_n(\theta_{\dagger})| \geq \frac{\sqrt{2}M}{\sqrt{n}} (J(\mathcal{H}, d_n) + tU)\right) \leq \exp\left(-\frac{t^2}{2}\right), \quad t \geq 1. \quad (\text{A.7})$$

Since  $Z_n(\theta_{\dagger}) \sim \mathcal{N}(0, n^{-1}\|h_{\theta_{\dagger}}\|_n^2)$ , it holds for all  $\theta_{\dagger} \in \Theta$

$$\mathbb{P}\left(|Z_n(\theta_{\dagger})| \geq \frac{MtU}{\sqrt{n}}\right) \leq \exp\left(\frac{-M^2U^2t^2}{2\|h_{\theta_{\dagger}}\|_n^2}\right) \leq \exp\left(-\frac{M^2t^2}{2}\right), \quad t \geq 0. \quad (\text{A.8})$$

Combining (A.7) and (A.8) yields for  $t \geq 1$

$$\begin{aligned} &\mathbb{P}\left(\sup_{\theta \in \Theta} |Z_n(\theta)| \geq \frac{3M}{\sqrt{n}} (J(\mathcal{H}, d_n) + tU)\right) \\ &\leq \mathbb{P}\left(\sup_{\theta \in \Theta} |Z_n(\theta) - Z_n(\theta_{\dagger})| \geq \frac{\sqrt{2}M}{\sqrt{n}} (J(\mathcal{H}, d_n) + tU)\right) \\ &\quad + \mathbb{P}\left(|Z_n(\theta_{\dagger})| \geq \frac{M}{\sqrt{n}} (J(\mathcal{H}, d_n) + tU)\right) \\ &\leq \exp\left(\frac{-t^2}{2}\right) + \mathbb{P}\left(|Z_n(\theta_{\dagger})| \geq \frac{MtU}{\sqrt{n}}\right) \\ &\leq \exp\left(\frac{-t^2}{2}\right) + \exp\left(\frac{-M^2t^2}{2}\right) \\ &= 2 \exp\left(\frac{-t^2}{2}\right), \end{aligned}$$

where we assumed without loss of generality that  $M \geq 1$ . Substitute  $\delta = 3M/\sqrt{n} (J(\mathcal{H}, d_n) + tU)$ , i.e.  $t = (\sqrt{n}\delta/3M - J(\mathcal{H}, d_n))/U$ . Because  $N(\mathcal{H}, d_n, \tau) \geq 2$  for  $\tau \leq U/2$ , we have that  $J(\mathcal{H}, d_n) \geq U/2\sqrt{\log(2)} > U/4$ . Therefore  $\sqrt{n}\delta \geq 15MJ(\mathcal{H}, d_n) := C_{\text{Ch}}J(\mathcal{H}, d_n)$  implies  $t \geq 1$  and thus

$$\mathbb{P}\left(\sup_{\theta \in \Theta} |Z_n(\theta)| \geq \delta\right) \leq 2 \exp\left(-\frac{(\sqrt{n}\delta/(3M) - J(\mathcal{H}, d_n))^2}{2U^2}\right) \leq 2 \exp\left(-\frac{8n\delta^2}{C_{\text{Ch}}^2U^2}\right) \quad (\text{A.9})$$

which gives the claim for the white noise case.

**Sub-Gaussian case.** For  $\theta, \theta' \in \Theta$  it holds

$$Z_n(\theta) - Z_n(\theta') = \frac{1}{n} \sum_{i=1}^n \langle \varepsilon_i, h_\theta(x_i) - h_{\theta'}(x_i) \rangle_{\mathbf{y}}.$$

Since the centered RVs  $n^{-1} \langle \varepsilon_i, h_\theta(x_i) - h_{\theta'}(x_i) \rangle_{\mathbf{y}}$  are i.i.d. sub-Gaussian with parameter  $n^{-1} \|h_\theta(x_i) - h_{\theta'}(x_i)\|_{\mathbf{y}}$ , the ‘generalized’ Hoeffding inequality for sub-Gaussian variables (see Vershynin 2018, Theorem 2.6.2) implies that for some universal constant  $c > 0$  the increment  $Z_n(\theta) - Z_n(\theta')$  is sub-Gaussian with parameter  $c \|h_\theta - h_{\theta'}\|_n / \sqrt{n}$ . Therefore

$$\mathbb{P} \left( \left| Z_n(\theta) - Z_n(\theta') \right| \geq \frac{\sqrt{2}tc}{\sqrt{n}} d_n(\theta, \theta') \right) \leq 2 \exp(-t^2), \quad t \geq 0.$$

From here on, the proof is similar to the white noise case and we obtain (A.9) by absorbing  $c$  into  $C_{\text{Ch}}$ .  $\blacksquare$

The following Lemma uses the ‘peeling device’ proof technique as developed in van de Geer (2000, Section 5) and is a generalization of van de Geer (2000, Lemma 5.13).

**Lemma 33** *Let  $x_i \stackrel{iid}{\sim} \gamma$  and consider the entropy  $H(\delta) = H(\mathbf{G}, \|\cdot\|_{\infty, \text{supp}(\gamma)}, \delta)$ . Let  $(\tilde{\delta}_n)_{n \in \mathbb{N}}$  be a positive sequence with*

$$n\tilde{\delta}_n^2 \geq 6M_{\mathbf{F}}^2 H(\tilde{\delta}_n).$$

*Then for  $R \geq \max\{8\tilde{\delta}_n, 18M_{\mathbf{F}}/\sqrt{n}\}$ , it holds that*

$$\mathbb{P} \left( \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq R, \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq 2\|\hat{G}_n - G_0\|_n \right) \leq 2 \exp \left( -\frac{nR^2}{320M_{\mathbf{F}}^2} \right).$$

**Proof** In the following we write  $\|\cdot\|_{\infty} = \|\cdot\|_{\infty, \text{supp}(\gamma)}$ . Let  $\mathbf{F} = \{G - G_0 : G \in \mathbf{G}\}$ . Then for  $R \geq 0$  it holds that

$$\begin{aligned} & \mathbb{P}_{G_0} \left( \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq R, \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq 2\|\hat{G}_n - G_0\|_n \right) \\ & \leq \mathbb{P}_{G_0} \left( \sup_{F \in \mathbf{F}, \|F\|_{L^2(\gamma)} \geq R} \|F\|_{L^2(\gamma)} \geq 2\|F\|_n \right). \end{aligned} \quad (\text{A.10})$$

For  $s \in \mathbb{N}$ , we define  $\mathbf{F}_s = \{F \in \mathbf{F} : sR \leq \|F\|_{L^2(\gamma)} \leq (s+1)R\}$ . As a union of disjoint events, it holds

$$\begin{aligned} \mathbb{P}_{G_0} \left( \sup_{F \in \mathbf{F}, \|F\|_{L^2(\gamma)} \geq R} \|F\|_{L^2(\gamma)} \geq 2\|F\|_n \right) &= \sum_{s=1}^{\infty} \mathbb{P}_{G_0} \left( \sup_{F \in \mathbf{F}_s} \|F\|_{L^2(\gamma)} \geq 2\|F\|_n \right) \\ &\leq \sum_{s=1}^{\infty} \mathbb{P}_{G_0} \left( \sup_{F \in \mathbf{F}_s} \left| \|F\|_{L^2(\gamma)} - \|F\|_n \right| \geq \frac{sR}{2} \right), \end{aligned}$$

where we used for  $F \in \mathbf{F}_s$  and  $s \in \mathbb{N}$

$$\|F\|_{L^2(\gamma)} \geq 2\|F\|_n \implies \|F\|_{L^2(\gamma)} - \|F\|_n \geq \frac{\|F\|_{L^2(\gamma)}}{2} \geq \frac{sR}{2}.$$

Now consider some covering  $\mathbf{F}_s^* = \{F_{s,j}\}_{j=1}^N$  with  $N = N(\mathbf{F}_s, \|\cdot\|_\infty, R/8)$ . Thus for arbitrary  $s \in \mathbb{N}$  and  $F \in \mathbf{F}_s$  there exists  $F_{s,j} \in \mathbf{F}_s^*$  s.t.  $\|F - F_{s,j}\|_\infty \leq R/8$ . We get

$$\begin{aligned} \left| \|F\|_{L^2(\gamma)} - \|F\|_n \right| &\leq \left| \|F\|_{L^2(\gamma)} - \|F_{s,j}\|_{L^2(\gamma)} \right| + \left| \|F_{s,j}\|_{L^2(\gamma)} - \|F_{s,j}\|_n \right| \\ &\quad + \left| \|F_{s,j}\|_n - \|F\|_n \right| \\ &\leq \frac{R}{4} + \left| \|F_{s,j}\|_{L^2(\gamma)} - \|F_{s,j}\|_n \right|. \end{aligned}$$

Therefore

$$\begin{aligned} &\sum_{s=1}^{\infty} \mathbb{P}_{G_0} \left( \sup_{F \in \mathbf{F}_s} \left| \|F\|_{L^2(\gamma)} - \|F\|_n \right| \geq \frac{sR}{2} \right) \\ &\leq \sum_{s=1}^{\infty} \mathbb{P}_{G_0} \left( \max_{j=1, \dots, N} \left| \|F_{s,j}\|_{L^2(\gamma)} - \|F_{s,j}\|_n \right| \geq \frac{sR}{4} \right) \\ &\leq \sum_{s=1}^{\infty} N(\mathbf{F}_s, \|\cdot\|_\infty, R/8) \max_{j=1, \dots, N} \mathbb{P}_{G_0} \left( \left| \|F_{s,j}\|_{L^2(\gamma)} - \|F_{s,j}\|_n \right| \geq \frac{sR}{4} \right), \quad (\text{A.11}) \end{aligned}$$

where we used  $sR/2 - R/4 \geq sR/4$  for  $s \in \mathbb{N}$ . Since  $F_{s,j} \in \mathbf{F}_s$ , we have  $\|F_{s,j}\|_{L^2(\gamma)} \geq sR$  and therefore

$$\begin{aligned} \left| \|F_{s,j}\|_{L^2(\gamma)} - \|F_{s,j}\|_n \right| &= \frac{\left| \|F_{s,j}\|_{L^2(\gamma)}^2 - \|F_{s,j}\|_n^2 \right|}{\|F_{s,j}\|_{L^2(\gamma)} + \|F_{s,j}\|_n} \\ &\leq \frac{1}{sR} \left| \|F_{s,j}\|_{L^2(\gamma)}^2 - \|F_{s,j}\|_n^2 \right|. \end{aligned}$$

Inserting into (A.11) gives

$$\begin{aligned} &\sum_{s=1}^{\infty} N(\mathbf{F}_s, \|\cdot\|_\infty, R/8) \max_{j=1, \dots, N} \mathbb{P}_{G_0} \left( \left| \|F_{s,j}\|_{L^2(\gamma)} - \|F_{s,j}\|_n \right| \geq \frac{sR}{4} \right) \\ &\leq \sum_{s=1}^{\infty} N(\mathbf{F}_s, \|\cdot\|_\infty, R/8) \max_{j=1, \dots, N} \mathbb{P}_{G_0} \left( \left| \|F_{s,j}\|_{L^2(\gamma)}^2 - \|F_{s,j}\|_n^2 \right| \geq \frac{s^2 R^2}{4} \right). \end{aligned}$$

Define the variables

$$Y_i = \frac{1}{n} \left( \|F_{s,j}\|_{L^2(\gamma)}^2 - \|F_{s,j}(x_i)\|_y^2 \right), \quad i = 1, \dots, n.$$

It holds for all  $i = 1, \dots, n$

$$\begin{aligned}
 \mathbb{E}[Y_i] &= 0, \\
 |Y_i| &\leq \frac{2M_{\mathbf{F}}^2}{n}, \\
 \mathbb{E}[Y_i^2] &= \frac{1}{n^2} \mathbb{E} \left[ \left( \|F_{s,j}\|_{L^2(\gamma)}^2 - \|F_{s,j}(x_i)\|_{\mathbb{Y}}^2 \right)^2 \right] \\
 &= \frac{1}{n^2} \mathbb{E} \left[ \left( \|F_{s,j}\|_{L^2(\gamma)}^4 - 2\|F_{s,j}\|_{L^2(\gamma)}^2 \|F_{s,j}(x_i)\|_{\mathbb{Y}}^2 + \|F_{s,j}(x_i)\|_{\mathbb{Y}}^4 \right) \right] \\
 &= \frac{1}{n^2} \mathbb{E} \left[ \|F_{s,j}(x_i)\|_{\mathbb{Y}}^4 \right] - \|F_{s,j}\|_{L^2(\gamma)}^4 \\
 &\leq \frac{M_{\mathbf{F}}^2}{n^2} \|F_{s,j}\|_{L^2(\gamma)}^2.
 \end{aligned}$$

Applying Bernstein's inequality (Lemma 29) for the variables  $Y_i$  yields

$$\begin{aligned}
 &\sum_{s=1}^{\infty} N(\mathbf{F}_s, \|\cdot\|_{\infty}, R/8) \max_{j=1, \dots, N} \mathbb{P}_{G_0} \left( \left| \|F_{s,j}\|_{L^2(\gamma)}^2 - \|F_{s,j}\|_n^2 \right| \geq \frac{s^2 R^2}{4} \right) \\
 &\leq \sum_{s=1}^{\infty} N(\mathbf{F}_s, \|\cdot\|_{\infty}, R/8) \max_{j=1, \dots, N} \exp \left( -\frac{ns^4 R^4}{32M_{\mathbf{F}}^2 (\|F_{s,j}\|_{L^2(\gamma)}^2 + \frac{s^2 R^2}{6})} \right) \\
 &\leq \sum_{s=1}^{\infty} \exp \left( H(\mathbf{F}_s, \|\cdot\|_{\infty}, R/8) - \frac{ns^2 R^2}{160M_{\mathbf{F}}^2} \right),
 \end{aligned}$$

where we used  $\|F_{s,j}\|_{L^2(\gamma)}^2 \leq (s+1)^2 R^2 \leq 4s^2 R^2$  for all  $j = 1, \dots, N$  and  $s \in \mathbb{N}$ , since  $F_{j,s} \in \mathbf{F}_s$ .

Since  $H(\delta)/\delta^2$  is non-increasing in  $\delta$  and  $R \geq 8\tilde{\delta}_n$ , we have

$$\frac{ns^2 R^2}{160M_{\mathbf{F}}^2} \geq \frac{n}{3M_{\mathbf{F}}^2} \left( \frac{R}{8} \right)^2 \geq 2H(\mathbf{G}, \|\cdot\|_{\infty}, R/8) \geq 2H(\mathbf{F}_s, \|\cdot\|_{\infty}, R/8).$$

Therefore we get

$$\begin{aligned}
 &\sum_{s=1}^{\infty} \exp \left( H(\mathbf{F}_s, \|\cdot\|_{\infty}, R/8) - \frac{ns^2 R^2}{160M_{\mathbf{F}}^2} \right) \\
 &\leq \sum_{s=1}^{\infty} \exp \left( -\frac{ns^2 R^2}{320M_{\mathbf{F}}^2} \right) \leq \sum_{s=1}^{\infty} \frac{1}{s^2} \exp \left( -\frac{nR^2}{320M_{\mathbf{F}}^2} \right) < 2 \exp \left( -\frac{nR^2}{320M_{\mathbf{F}}^2} \right). \quad (\text{A.12})
 \end{aligned}$$

Note that for  $x, y \geq 1$ , we have  $\exp(-xy) \leq \exp(-y)/x$ . Applying this in the case  $x = s^2 \geq 1$  and  $y = nR^2/(320M_{\mathbf{F}}^2) \geq 1$  for  $R \geq 18M_{\mathbf{F}}/\sqrt{n}$ , together with  $\sum_{s=1}^{\infty} 1/s^2 = \pi^2/6 < 2$ , gives (A.12). Combining (A.10)–(A.12) shows the claim.  $\blacksquare$

In the case of sub-Gaussian noise, standard bounds on the maximum of sub-Gaussian RVs suffice to bound the generalisation error.

**Lemma 34** Consider the sub-Gaussian noise model, i.e. suppose  $\|\varepsilon_i\|_{\mathbf{y}}$  are i.i.d. sub-Gaussian with parameter 1. Abbreviate  $H(\delta) = H(\mathbf{G}, \|\cdot\|_{\infty, \text{supp}(\gamma)}, \delta)$ . Then there exists a universal constant  $C > 0$  s.t. for all positive sequences  $(\delta_n)_{n \in \mathbb{N}}$  with

$$n\delta_n^4 \geq C^2 \sigma^2 M_{\mathbf{F}}^2 H\left(\frac{\delta_n^2}{8\sigma^2 + \delta_n^2}\right),$$

all  $G^* \in \mathbf{G}$  and  $R \geq \max\{\delta_n, \sqrt{2}\|G^* - G_0\|_n\}$ , it holds for every  $\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{X}^n$

$$\mathbb{P}_{G_0}^{\mathbf{x}}\left(\|\hat{G}_n - G_0\|_n \geq R\right) \leq 4 \exp\left(-\frac{nR^4}{C^2 \sigma^2 (1 + M_{\mathbf{F}}^2)}\right). \quad (\text{A.13})$$

**Proof** In the following we write  $\|\cdot\|_{\infty} = \|\cdot\|_{\infty, \text{supp}(\gamma)}$ . For  $R^2 \geq 2\|G^* - G_0\|_n^2$ , we use the basic inequality (Lemma 30) to obtain

$$\begin{aligned} \mathbb{P}_{G_0}^{\mathbf{x}}\left(\|\hat{G}_n - G_0\|_n^2 \geq R^2\right) &\leq \mathbb{P}_{G_0}^{\mathbf{x}}\left(\|\hat{G}_n - G_0\|_n^2 - \|G^* - G_0\|_n^2 \geq \frac{R^2}{2}\right) \\ &\leq \mathbb{P}_{G_0}^{\mathbf{x}}\left(\left|\frac{1}{n} \sum_{i=1}^n \langle \varepsilon_i, \hat{G}_n(x_i) - G^*(x_i) \rangle_{\mathbf{y}}\right| \geq \frac{R^2}{4\sigma}\right). \end{aligned} \quad (\text{A.14})$$

For  $R > 0$ , let  $\mathbf{G}^* = \{G - G^*, G \in \mathbf{G}\}$  and  $(G_j)_{j=1}^N$  with  $N = N(\mathbf{G}^*, \|\cdot\|_{\infty}, R^2/(8\sigma\mathbb{E}[\|\varepsilon_i\|_{\mathbf{y}}] + R^2))$  denote a minimal  $\|\cdot\|_{\infty}$ -cover of  $\mathbf{G}^*$ .

It holds for  $R > 0$

$$\begin{aligned} &\mathbb{P}_{G_0}^{\mathbf{x}}\left(\left|\frac{1}{n} \sum_{i=1}^n \langle \varepsilon_i, \hat{G}_n(x_i) - G^*(x_i) \rangle_{\mathbf{y}}\right| \geq \frac{R^2}{4\sigma}\right) \\ &\leq \mathbb{P}_{G_0}^{\mathbf{x}}\left(\sup_{G \in \mathbf{G}} \left|\frac{1}{n} \sum_{i=1}^n \langle \varepsilon_i, G(x_i) - G^*(x_i) \rangle_{\mathbf{y}}\right| \geq \frac{R^2}{4\sigma}\right) \\ &\leq \mathbb{P}_{G_0}^{\mathbf{x}}\left(\sup_{G \in \mathbf{G}} \left|\frac{1}{n} \sum_{i=1}^n \langle \varepsilon_i, G(x_i) - G^*(x_i) - G_{j^*}(x_i) \rangle_{\mathbf{y}}\right| + \max_{j=1, \dots, N} \left|\frac{1}{n} \sum_{i=1}^n \langle \varepsilon_i, G_j(x_i) \rangle_{\mathbf{y}}\right| \geq \frac{R^2}{4\sigma}\right) \\ &\leq \underbrace{\mathbb{P}_{G_0}^{\mathbf{x}}\left(\sup_{G \in \mathbf{G}} \left|\frac{1}{n} \sum_{i=1}^n \langle \varepsilon_i, G(x_i) - G^*(x_i) - G_{j^*}(x_i) \rangle_{\mathbf{y}}\right| \geq \frac{R^2}{8\sigma}\right)}_{(i)} \\ &\quad + \underbrace{\mathbb{P}_{G_0}^{\mathbf{x}}\left(\max_{j=1, \dots, N} \left|\frac{1}{n} \sum_{i=1}^n \langle \varepsilon_i, G_j(x_i) \rangle_{\mathbf{y}}\right| \geq \frac{R^2}{8\sigma}\right)}_{(ii)}. \end{aligned}$$

We estimate the terms (i) and (ii) separately. For (i), use  $\|G - G^* - G_{j^*}\|_\infty \leq R^2/(8\sigma\mathbb{E}[\|\varepsilon_i\|_{\mathcal{Y}}] + R^2)$  for all  $G \in \mathbf{G}$  and estimate

$$\begin{aligned} & \mathbb{P}_{G_0}^{\mathbf{x}} \left( \sup_{G \in \mathbf{G}} \left| \frac{1}{n} \sum_{i=1}^n \langle \varepsilon_i, G(x_i) - G^*(x_i) - G_{j^*}(x_i) \rangle_{\mathcal{Y}} \right| \geq \frac{R^2}{8\sigma} \right) \\ & \leq \mathbb{P}_{G_0}^{\mathbf{x}} \left( \left| \frac{1}{n} \sum_{i=1}^n \|\varepsilon_i\|_{\mathcal{Y}} \frac{R^2}{8\sigma\mathbb{E}[\|\varepsilon_i\|_{\mathcal{Y}}] + R^2} \right| \geq \frac{R^2}{8\sigma} \right) \\ & \leq \mathbb{P}_{G_0}^{\mathbf{x}} \left( \left| \frac{1}{n} \sum_{i=1}^n \|\varepsilon_i\|_{\mathcal{Y}} - \mathbb{E}[\|\varepsilon_i\|_{\mathcal{Y}}] \right| \geq \frac{R^2}{8\sigma} \right) \leq 2 \exp \left( -\frac{nR^4}{C^2\sigma^2} \right), \end{aligned}$$

where we used the Hoeffding inequality for sub-Gaussian random variables from Vershynin (2018, Theorem 2.6.2).

For (ii), it holds

$$\begin{aligned} & \mathbb{P}_{G_0}^{\mathbf{x}} \left( \max_{j=1, \dots, N} \left| \frac{1}{n} \sum_{i=1}^n \langle \varepsilon_i, G_j(x_i) \rangle_{\mathcal{Y}} \right| \geq \frac{R^2}{8\sigma} \right) \\ & \leq N \max_{j=1, \dots, N} \mathbb{P}_{G_0}^{\mathbf{x}} \left( \left| \frac{1}{n} \sum_{i=1}^n \langle \varepsilon_i, G_j(x_i) \rangle_{\mathcal{Y}} \right| \geq \frac{R^2}{8\sigma} \right) \\ & \leq 2N \exp \left( -\frac{2nR^4}{C^2\sigma^2 M_{\mathbf{F}}^2} \right) \leq 2 \exp \left( H \left( \frac{R^2}{8\sigma^2 + R^2} \right) - \frac{2nR^4}{C^2\sigma^2 M_{\mathbf{F}}^2} \right), \end{aligned}$$

where we used  $\mathbb{E}[\|\varepsilon_i\|_{\mathcal{Y}}] \leq \sigma$  (Cauchy–Schwarz) and Vershynin (2018, Theorem 2.6.2) for  $Y_i = n^{-1} \langle \varepsilon_i, G_j(x_i) \rangle_{\mathcal{Y}}$ ,  $i = 1, \dots, n$ .

Since  $H(\delta)/\delta^2$  is non-increasing, also  $H(\delta/(8\sigma^2 + \delta))/\delta^2$  is non-increasing and thus it holds for  $R \geq \delta_n$

$$\frac{nR^4}{C^2\sigma^2 M_{\mathbf{F}}^2} \geq H \left( \frac{R^2}{8\sigma^2 + R^2} \right).$$

This gives for  $R \geq \delta_n$

$$(ii) \leq 2 \exp \left( -\frac{nR^4}{C^2\sigma^2 M_{\mathbf{F}}^2} \right). \quad (\text{A.15})$$

Combining (A.14)–(A.15) gives the result. ■

## Appendix B. Proofs of Section 2

### B.1 Proof of Theorem 5

#### B.1.1 EXISTENCE AND MEASURABILITY OF $\hat{G}_n$

**White noise case.** For  $i = 1, \dots, n$ , denote the probability space of the RVs  $x_i$  and the white noise processes  $\varepsilon_i$  as  $(\Omega, \Sigma, \mathbb{P})$ . Furthermore, equip the Hilbert spaces  $\mathcal{X}$  and  $\mathcal{Y}$  with

Borel  $\sigma$ -algebras  $\mathcal{B}_X$  and  $\mathcal{B}_Y$  and the space  $\mathbf{G}$  with the Borel  $\sigma$ -algebra  $\mathcal{B}_{\mathbf{G}}$ . Then, there exists an orthonormal basis  $(\psi_j)_{j \in \mathbb{N}}$  of  $\mathcal{Y}$  and i.i.d. Gaussian variables  $Z_j \sim \mathcal{N}(0, 1)$  s.t.

$$\begin{aligned} \varepsilon_i &: \Omega \times \mathcal{Y} \rightarrow \mathbb{R}, \\ \varepsilon_i(\omega, y) &= \sum_{j=1}^{\infty} \langle y, \psi_j \rangle_{\mathcal{Y}} Z_j(\omega) \end{aligned}$$

for  $i = 1, \dots, n$ , see also Giné and Nickl (2016, Example 2.1.11)

Recall the noise level  $\sigma > 0$ . For  $i = 1, \dots, n$ , we define

$$\begin{aligned} u_i &: \Omega \times \mathbf{G} \rightarrow \mathbb{R}, \\ u_i(\omega, G) &= 2\sigma \varepsilon_i(\omega, G(x_i(\omega))) + 2 \langle G_0(x_i(\omega)), G(x_i(\omega)) \rangle_{\mathcal{Y}} - \|G(x_i(\omega))\|_{\mathcal{Y}}^2. \end{aligned} \quad (\text{B.1})$$

We aim to apply Nickl (2007, Proposition 5) to  $u := 1/n \sum_{i=1}^n u_i$  in order to get existence and measurability of  $\hat{G}_n$  in (2.3).

Per assumption, the metric space  $(\mathbf{G}, \|\cdot\|)$  is compact. We show that  $u_i$  from (B.1) is measurable in the first component and continuous in the second component for all  $i = 1, \dots, n$ . Then Nickl (2007, Proposition 5) shows the claim. Consider an arbitrary  $G \in \mathbf{G}$ . We show that  $u_i(\cdot, G)$  is  $(\Sigma, \mathcal{B}_{\mathbb{R}})$ -measurable, where  $\mathcal{B}_{\mathbb{R}}$  is the Borel  $\sigma$ -algebra in  $\mathbb{R}$ .

The RVs  $x_i$  are  $(\Sigma, \mathcal{B}_X)$ -measurable by definition. The maps  $G$  and  $G_0$  are assumed to be  $(\mathcal{B}_X, \mathcal{B}_Y)$ -measurable. Furthermore, because of their continuity, the scalar product  $\langle \cdot, \cdot \rangle_{\mathcal{Y}}$  is  $(\mathcal{B}_Y, \mathcal{B}_{\mathbb{R}})$ -measurable in both components and also the norm  $\|\cdot\|_{\mathcal{Y}}$  is  $(\mathcal{B}_Y, \mathcal{B}_{\mathbb{R}})$ -measurable. Therefore, since the composition of measurable functions is measurable, the latter two summands in (B.1) are  $(\Sigma, \mathcal{B}_{\mathbb{R}})$ -measurable.

Proceeding with the first summand, the RVs  $Z_j$  are  $(\Sigma, \mathcal{B}_{\mathbb{R}})$ -measurable by definition for all  $j \in \mathbb{N}$ . Therefore the products  $\langle \cdot, \psi_j \rangle_{\mathcal{Y}} Z_j$  are  $(\Sigma \otimes \mathcal{B}_Y, \mathcal{B}_{\mathbb{R}})$ -measurable for all  $j \in \mathbb{N}$ . Thus  $\varepsilon_i$  is, as the pointwise limit,  $(\Sigma \otimes \mathcal{B}_Y, \mathcal{B}_{\mathbb{R}})$ -measurable. Then, as the composition of measurable functions, the first summand in (B.1) is  $(\Sigma, \mathcal{B}_{\mathbb{R}})$ -measurable. Therefore  $u_i(\cdot, G)$ ,  $i = 1, \dots, n$ , and thus  $u(\cdot, G)$  is  $(\Sigma, \mathcal{B}_{\mathbb{R}})$ -measurable for all  $G \in \mathbf{G}$ .

We proceed and show that  $u(\omega, \cdot)$  is continuous w.r.t.  $\|\cdot\|$ . Therefore choose  $G, G' \in \mathbf{G}$  and  $\omega \in \Omega$ . Then it holds for  $i = 1, \dots, n$  and  $x_i = x_i(\omega)$

$$\begin{aligned} |u_i(\omega, G) - u_i(\omega, G')| &\leq 2\sigma |\varepsilon_i(\omega, G(x_i) - G'(x_i))| + 2 |\langle G_0(x_i), G(x_i) - G'(x_i) \rangle_{\mathcal{Y}}| \\ &\quad + |\|G(x_i)\|_{\mathcal{Y}} - \|G'(x_i)\|_{\mathcal{Y}}| (\|G(x_i)\|_{\mathcal{Y}} + \|G'(x_i)\|_{\mathcal{Y}}) \\ &\leq 2\sigma |\varepsilon_i(\omega, G(x_i) - G'(x_i))| \\ &\quad + \sqrt{n} (2\|G_0(x_i)\|_{\mathcal{Y}} + \|G(x_i)\|_{\mathcal{Y}} + \|G'(x_i)\|_{\mathcal{Y}}) \|G - G'\|_n \\ &\leq 2\sigma |\varepsilon_i(\omega, G(x_i) - G'(x_i))| \\ &\quad + C\sqrt{n} (2\|G_0(x_i)\|_{\mathcal{Y}} + \|G(x_i)\|_{\mathcal{Y}} + \|G'(x_i)\|_{\mathcal{Y}}) \|G - G'\|, \end{aligned} \quad (\text{B.2})$$

where we used that  $\|G(x_i) - G'(x_i)\|_{\mathcal{Y}} \leq \sqrt{n} \|G - G'\|_n$  for  $i = 1, \dots, n$  and  $\|\cdot\|_n \leq C \|\cdot\|$  at the last inequality. Furthermore, Giné and Nickl (2016, page 40 and Proposition 2.3.7(a)) yields that the white noise processes  $\varepsilon_i$  are a.s. sample  $d_{\varepsilon_i}$ -continuous w.r.t. their intrinsic

pseudometrics  $d_{\varepsilon_i} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ ,  $d_{\varepsilon_i}(y, y') = \mathbb{E}[\langle \varepsilon_i, y - y' \rangle_{\mathcal{Y}}^2]^{1/2} = \|y - y'\|_{\mathcal{Y}}$ . Since for all  $G, G' \in \mathbf{G}$  we have

$$\|G(x_i) - G'(x_i)\|_{\mathcal{Y}} \leq \sqrt{n} \|G - G'\|_n \leq C\sqrt{n} \|G - G'\|,$$

the white noise processes are also a.s. sample  $d$ -continuous, where  $d$  is the metric induced by  $\|\cdot\|$ .

Together with (B.2) this shows that there exists a null-set  $\Omega_0 \subset \Omega$  such that for all  $\omega \in \Omega \setminus \Omega_0$ ,  $u(\omega, \cdot)$  is continuous w.r.t.  $\|\cdot\|$ . Now we choose versions  $\tilde{x}_i$  and  $\tilde{\varepsilon}_i$  s.t.  $u(\omega, \cdot) = 0$  for all  $\omega \in \Omega_0$ . Then  $G \mapsto u(\omega, G) : (\mathbf{G}, \|\cdot\|) \rightarrow \mathbb{R}$  is continuous for all  $\omega \in \Omega$ . Applying Nickl (2007, Proposition 5) gives an  $(\Sigma, \mathcal{B}_{\mathbf{G}})$ -measurable MLSE  $\hat{G}_n$  in (2.3) with the desired minimization property.

**Sub-Gaussian case.** For  $i = 1, \dots, n$ , consider the functions  $u_i$  from (B.1). The measurability of  $u_i(\cdot, G)$ ,  $i = 1, \dots, n$  follows from the measurability of the scalar product  $\langle \cdot, \cdot \rangle_{\mathcal{Y}}$ , the norm  $\|\cdot\|_{\mathcal{Y}}$ ,  $G_0$  and all  $G \in \mathbf{G}$ . Note that in contrast to the white noise case,  $\varepsilon_i$  are RVs in  $\mathcal{Y}$  for all  $i = 1, \dots, n$  and therefore measurable without any further investigation. Also, since  $\varepsilon_i(\omega) \in \mathcal{Y}$  for all  $\omega$ , Cauchy–Schwarz immediately shows that  $u_i(\omega, \cdot)$  and therefore  $u(\omega, \cdot)$  is continuous w.r.t.  $\|\cdot\|$  for all  $\omega \in \Omega \setminus \Omega^0$ . Choosing versions  $\tilde{x}_i$  and  $\tilde{\varepsilon}_i$  and applying Nickl (2007, Proposition 5) as above gives the existence of an  $(\Sigma, \mathcal{B}_{\mathbf{G}})$ -measurable LSE  $\hat{G}_n$  in (2.3) and therefore finishes the proof.

### B.1.2 CONCENTRATION INEQUALITY FOR $\hat{G}_n$

The proof follows ideas developed in van de Geer (2000, Section 10.3) as well as Nickl and Wang (2024), where generic chaining bounds from Dirksen (2015) were used to bound the relevant empirical processes appearing below.

*Slicing argument.* Recall the definition (2.4) of the empirical norm. For  $R_n^2 \geq 2\|G^* - G_0\|_n^2$ , we have

$$\mathbb{P}_{G_0}^{\mathbf{x}}(\|\hat{G}_n - G_0\|_n^2 \geq R_n^2) \leq \mathbb{P}_{G_0}^{\mathbf{x}}(2(\|\hat{G}_n - G_0\|_n^2 - \|G^* - G_0\|_n^2) \geq R_n^2). \quad (\text{B.3})$$

As a union of disjoint events, it holds

$$\begin{aligned} & \mathbb{P}_{G_0}^{\mathbf{x}}(2(\|\hat{G}_n - G_0\|_n^2 - \|G^* - G_0\|_n^2) \geq R_n^2) \\ &= \sum_{s=0}^{\infty} \mathbb{P}_{G_0}^{\mathbf{x}}\left(2^{2s} R_n^2 \leq 2(\|\hat{G}_n - G_0\|_n^2 - \|G^* - G_0\|_n^2) < 2^{2s+2} R_n^2\right). \end{aligned}$$

Applying the basic inequality (Lemma 30) and defining the empirical process  $(X_G : G \in \mathbf{G})$  indexed by the operator class  $\mathbf{G}$  as

$$X_G = \frac{1}{n} \sum_{i=1}^n \langle \varepsilon_i, G(x_i) - G^*(x_i) \rangle_{\mathcal{Y}},$$

gives

$$\begin{aligned}
 & \sum_{s=0}^{\infty} \mathbb{P}_{G_0}^{\mathbf{x}} \left( 2^{2s} R_n^2 \leq 2(\|\hat{G}_n - G_0\|_n^2 - \|G^* - G_0\|_n^2) < 2^{2s+2} R_n^2 \right) \\
 & \leq \sum_{s=0}^{\infty} \mathbb{P}_{G_0}^{\mathbf{x}} \left( |X_{\hat{G}_n}| \geq \frac{2^{2s-2} R_n^2}{\sigma}, \|\hat{G}_n - G_0\|_n^2 - \|G^* - G_0\|_n^2 < 2^{2s+1} R_n^2 \right) \\
 & \leq \sum_{s=0}^{\infty} \mathbb{P}_{G_0}^{\mathbf{x}} \left( \sup_{G \in \mathbf{G}_n^*(2^{s+3/2} R_n)} |X_G| \geq \frac{2^{2s-2} R_n^2}{\sigma} \right). \tag{B.4}
 \end{aligned}$$

In (B.4), we additionally used that if  $2\|G^* - G_0\|_n^2 \leq R_n^2$  and

$$\|\hat{G}_n - G_0\|_n^2 - \|G^* - G_0\|_n^2 < 2^{2s+1} R_n^2$$

then  $\|\hat{G}_n - G_0\|_n^2 \leq 2^{2s+1} R_n^2 + R_n^2/2$  and thus

$$\|\hat{G}_n - G^*\|_n^2 \leq 2(\|\hat{G}_n - G_0\|_n^2 + \|G^* - G_0\|_n^2) \leq 2^{2s+2} R_n^2 + 2R_n^2 \leq 2^{2s+3} R_n^2.$$

*Concentration inequality for each slice.* We wish to apply Lemma 32 to bound the probabilities in (B.4). Let  $C_{\text{Ch}}$  be the generic constant from this lemma and let  $\delta_n$  satisfy (2.6). Then, due to  $\delta \mapsto \Psi_n(\delta)/\delta^2$  being non-increasing, (2.6) gives for all  $R_n \geq \delta_n$  and  $s \in \mathbb{N}_0$

$$\sqrt{n}(2^{2s+3} R_n^2) \geq 32C_{\text{Ch}}\sigma\Psi_n(2^{s+3/2} R_n)$$

so that with  $\Theta := \mathbf{G}_n^*(2^{s+3/2} R_n)$

$$\sqrt{n} \frac{2^{2s-2} R_n^2}{\sigma} \geq C_{\text{Ch}} \Psi_n(2^{s+3/2} R_n) \geq C_{\text{Ch}} J(\Theta, \|\cdot\|_n). \tag{B.5}$$

With  $h_G := G - G^*$  and  $\mathcal{H} := \{h_G : G \in \Theta\}$  we have  $J(\Theta, \|\cdot\|_n) = J(\mathcal{H}, \|\cdot\|_n)$  and thus (B.5) verifies assumption (A.4) of Lemma 32 for  $\delta = 2^{2s-2} R_n^2/\sigma$ .

Furthermore, since  $\Theta = \mathbf{G}_n^*(2^{s+3/2} R_n) \subset \mathbf{G}$  for  $s \in \mathbb{N}_0$ ,  $R_n > 0$  and  $(\mathbf{G}, \|\cdot\|_n)$  is compact, the space  $(\Theta, \|\cdot\|_n)$  is separable, which verifies the last assumption of Lemma 32. Applying this lemma with  $U = 2^{s+3/2} R_n$  shows that the (uncountable) suprema in (B.4) are measurable and that for all  $R_n \geq \delta_n$ ,

$$\begin{aligned}
 \sum_{s=0}^{\infty} \mathbb{P}_{G_0}^{\mathbf{x}} \left( \sup_{G \in \mathbf{G}_n^*(2^{s+3/2} R_n)} |X_G| \geq \frac{2^{2s-2} R_n^2}{\sigma} \right) & \leq \sum_{s=0}^{\infty} \exp\left(-\frac{8n2^{4s-4} R_n^4}{C_{\text{Ch}}^2 \sigma^2 2^{2s+3} R_n^2}\right) \\
 & \leq \sum_{s=0}^{\infty} \exp\left(-\frac{2^{2s-4} n R_n^2}{C_{\text{Ch}}^2 \sigma^2}\right) \\
 & \leq \sum_{s=0}^{\infty} 2^{-2s} \exp\left(-\frac{n R_n^2}{16C_{\text{Ch}}^2 \sigma^2}\right) \\
 & < 2 \exp\left(-\frac{n R_n^2}{16C_{\text{Ch}}^2 \sigma^2}\right). \tag{B.6}
 \end{aligned}$$

In (B.6) we additionally used  $\exp(-xy) \leq \exp(-x)/y$ , which holds for all  $x, y \geq 1$ , i.e. for  $R_n \geq 4C_{\text{Ch}}\sigma/\sqrt{n}$ . Combining (B.3)–(B.4) and (B.6) gives (2.7) and therefore shows the claim.

## B.2 Proof of Theorem 6

We use the concentration inequality from Theorem 5 for the empirical error and combine it with a key concentration result for the empirical norm around the  $L^2(\gamma)$ -norm, proved in Lemma 33. For any  $R > 0$ ,

$$\begin{aligned} \mathbb{P}_{G_0} \left( \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq R \right) &\leq \mathbb{P}_{G_0} \left( \|\hat{G}_n - G_0\|_n \geq \frac{R}{2} \right) \\ &\quad + \mathbb{P}_{G_0} \left( \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq R, \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq 2\|\hat{G}_n - G_0\|_n \right). \end{aligned} \quad (\text{B.7})$$

Now suppose that  $R_n$  satisfies

$$R_n \geq 2 \max \left\{ \delta_n, 4\tilde{\delta}_n, \sqrt{2}\|G^* - G_0\|_{\infty, \text{supp}(\gamma)}, \frac{4C_{ch}\sigma}{\sqrt{n}}, \frac{9M_{\mathbf{F}}}{\sqrt{n}} \right\}.$$

This is implied by (2.11) for an appropriate choice of  $C$ . In particular,  $R_n \geq 2\sqrt{2}\|G^* - G_0\|_n$  for all  $\mathbf{x} \in \mathcal{X}^n$ . Since Theorem 5 holds for  $\gamma^n$ -almost every  $\mathbf{x} \in \mathcal{X}^n$ , taking expectations over  $\mathbf{x} \sim \gamma^n$  gives

$$\mathbb{P}_{G_0} \left( \|\hat{G}_n - G_0\|_n \geq \frac{R_n}{2} \right) \leq 2 \exp \left( -\frac{nR_n^2}{64C_{Ch}^2\sigma^2} \right). \quad (\text{B.8})$$

Furthermore, Lemma 33 gives

$$\begin{aligned} &\mathbb{P}_{G_0} \left( \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq R_n, \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq 2\|\hat{G}_n - G_0\|_n \right) \\ &\leq 2 \exp \left( -\frac{nR_n^2}{320M_{\mathbf{F}}^2} \right), \end{aligned} \quad (\text{B.9})$$

since we have assumed  $R_n \geq \max\{8\tilde{\delta}_n, 18M_{\mathbf{F}}/\sqrt{n}\}$ . Combining (B.8) and (B.9) shows (2.12) for some  $C$ .

## B.3 Proof of Corollary 7

It remains to show (2.13). We use (B.7) to estimate

$$\begin{aligned} &\mathbb{E}_{G_0} [\|\hat{G}_n - G_0\|_{L^2(\gamma)}^2] \\ &= \int_0^\infty \mathbb{P}_{G_0} \left( \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq \sqrt{R_n} \right) dR_n \\ &\leq \int_0^\infty \int_{\mathcal{X}^n} \mathbb{P}_{G_0}^{\mathbf{x}} \left( \|\hat{G}_n - G_0\|_n \geq \frac{\sqrt{R_n}}{2} \right) d\gamma(\mathbf{x}) dR_n \\ &\quad + \int_0^\infty \mathbb{P}_{G_0} \left( \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq \sqrt{R_n}, \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq 2\|\hat{G}_n - G_0\|_n \right) dR_n. \end{aligned} \quad (\text{B.10})$$

For  $G^* \in \mathbf{G}$  set  $R_{n,1} = 2 \max\{\delta_n, \sqrt{2}\|G^* - G_0\|_n, \frac{4C_{\text{Ch}}\sigma}{\sqrt{n}}\}$ , where  $C_{\text{Ch}}$  is from Lemma 32. Then Theorem 5 gives

$$\begin{aligned}
 & \int_0^\infty \int_{\mathcal{X}^n} \mathbb{P}_{G_0}^{\mathbf{x}} \left( \|\hat{G}_n - G_0\|_n \geq \frac{\sqrt{R_n}}{2} \right) d\gamma(x) dR_n \\
 & \leq \int_{\mathcal{X}^n} R_{n,1}^2 d\gamma(\mathbf{x}) + 2 \int_{R_{n,1}^2}^\infty \exp\left(-\frac{nR_n}{64C_{\text{Ch}}^2\sigma^2}\right) dR_n \\
 & \leq 4\delta_n^2 + 8\|G^* - G_0\|_{L^2(\gamma)}^2 + \frac{16C_{\text{Ch}}^2\sigma^2}{n} + 2 \int_{R_{n,1}^2}^\infty \exp\left(-\frac{nR_n}{64C_{\text{Ch}}^2\sigma^2}\right) dR_n \\
 & \leq 4\delta_n^2 + 8\|G^* - G_0\|_{L^2(\gamma)}^2 + \frac{80C_{\text{Ch}}^2\sigma^2}{n}.
 \end{aligned}$$

Lemma 33 gives with  $R_{n,2} = \max\{8\tilde{\delta}_n, 18M_{\mathbf{F}}/\sqrt{n}\}$

$$\begin{aligned}
 & \int_0^\infty \mathbb{P}_{G_0} \left( \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq \sqrt{R_n}, \|\hat{G}_n - G_0\|_{L^2(\gamma)} \geq 2\|\hat{G}_n - G_0\|_n \right) dR_n \\
 & \leq R_{n,2}^2 + 2 \int_{R_{n,1}^2}^\infty \exp\left(-\frac{nR_n}{320M_{\mathbf{F}}^2}\right) dR_n \\
 & \leq 64\tilde{\delta}_n^2 + \frac{964M_{\mathbf{F}}^2}{n}. \tag{B.11}
 \end{aligned}$$

Combining (B.10)–(B.11) and taking an infimum over  $G^* \in \mathbf{G}$  shows (2.13) for some  $C$  and thus finishes the proof of Corollary 7.

#### B.4 Proof of Corollary 8

We construct sequences  $(\delta_n)_{n \in \mathbb{N}}$  and  $(\tilde{\delta}_n)_{n \in \mathbb{N}}$  satisfying (2.10) and balance the approximation term and the entropy terms in (2.12) by choosing  $N(n)$  appropriately.

**Proof of (i):** Choosing  $\tilde{\delta}_n^2 \simeq n^{-2/(2+\alpha)}$  immediately shows the second part of (2.10). For  $J(\delta)$  from (2.5) it holds

$$J(\delta) \leq \int_0^\delta \sqrt{H(\rho, N)} d\rho \lesssim \delta^{1-\alpha/2} =: \Psi_n(\delta).$$

Hence  $\delta_n^2 \simeq n^{-2/(2+\alpha)}$  satisfies the first part of (2.10). Therefore Corollary 7 shows

$$\mathbb{E}_{G_0} \left[ \|\hat{G}_n - G_0\|_{L^2(\gamma)}^2 \right] \lesssim N^{-\beta} + n^{-\frac{2}{2+\alpha}}.$$

Choosing  $N = \lceil n^{\frac{2}{(2+\alpha)\beta}} \rceil$  gives the claim.

**Proof of (ii):** Choosing  $\tilde{\delta}_n^2 \simeq N(1 + \log(n))/n$  shows in particular  $\log(n) \gtrsim \log(\tilde{\delta}_n^{-1})$ , which in turn shows that  $\tilde{\delta}_n$  satisfies the second part of (2.10). For  $J(\delta)$  from (2.5) it holds

$$J(\delta) \lesssim \int_0^\delta \sqrt{H(\rho, N)} d\rho \lesssim \sqrt{N}\delta (1 + \log(\delta^{-1})) =: \Psi_n(\delta).$$

Hence  $\delta_n^2 \simeq N(1 + \log(n))/n$  satisfies the first part of (2.10). Therefore Corollary 7 shows

$$\mathbb{E}_{G_0} \left[ \|\hat{G}_n - G_0\|_{L^2(\gamma)}^2 \right] \lesssim N^{-\beta} + \frac{N(1 + \log(n))}{n}.$$

Choosing  $N(n) = \lceil n^{\frac{1}{\beta+1}} \rceil$  for all  $n \geq 1$  gives the claim.

## B.5 Proof of Theorem 10

The concentration inequality (2.15) follows immediately from (B.7), (B.9) (which holds for sub-Gaussian noise) and (A.13) from Lemma 34. Application of Lemma 34 yields

$$\begin{aligned} & \int_0^\infty \int_{\mathcal{X}^n} \mathbb{P}_{G_0}^{\mathbf{x}} \left( \|\hat{G}_n - G_0\|_n \geq \frac{\sqrt{\delta}}{2} \right) d\gamma(\mathbf{x}) d\delta \\ & \leq 4\delta_n^2 + 8\|G^* - G_0\|_{L^2(\gamma)}^2 + 4 \int_0^\infty \exp\left(-\frac{n\delta^2}{16C_1^2\sigma^2(1+M_{\mathbf{F}}^2)}\right) d\delta \\ & \leq 4\delta_n^2 + 8\|G^* - G_0\|_{L^2(\gamma)}^2 + 16C_1\sigma(1+M_{\mathbf{F}})\sqrt{\frac{\pi}{n}}. \end{aligned}$$

Using (B.10) and (B.11), and minimizing over  $G^* \in \mathbf{G}$  gives the bound on the mean-squared error (2.16) for some  $C_2$ .

## Appendix C. Neural Network Theory

In this section we recap elementary operations and approximation theory for neural networks, based on Petersen and Voigtlaender (2018). For a NN  $f : \mathbb{R}^{p_0} \rightarrow \mathbb{R}^{p_{L+1}}$  (see Definition 15) we denote by  $\text{size}_{\text{in}}(f)$  and  $\text{size}_{\text{out}}(f)$  the number of nonzero weights and biases of the first and last layer respectively.

### C.1 Operations on Neural Networks

In this subsection, we recap elementary operations on NNs. We start with the parallelization of two NNs following Opschoor et al. (2022a, Section 2.2.1).

**Definition 35 (Parallelization)** *Let  $q \in \mathbb{N}$ ,  $q \geq 2$  and  $\sigma \in \{\sigma_1, \sigma_q\}$ . Let  $f$  and  $g$  be two  $\sigma$ -NNs realizing the functions  $f$  and  $g$  with the same depth  $L$ . Furthermore, define the input dimensions of  $f$  and  $g$  as  $n_f$  and  $n_g$  and the output dimensions as  $m_f$  and  $m_g$ .<sup>3</sup> Then there exists a  $\sigma$ -NN  $(f, g)$ , called the parallelization of  $f$  and  $g$ , which simultaneously realizes  $f$  and  $g$ , i.e.*

$$(f, g) : \mathbb{R}^{n_f} \times \mathbb{R}^{n_g} \rightarrow \mathbb{R}^{m_f} \times \mathbb{R}^{m_g} : (\mathbf{x}, \tilde{\mathbf{x}}) \mapsto (f(\mathbf{x}), g(\tilde{\mathbf{x}})).$$

*It holds*

$$\begin{aligned} \text{size}((f, g)) &= \text{size}(f) + \text{size}(g), \\ \text{depth}((f, g)) &= \text{depth}(f) = \text{depth}(g), \\ \text{width}((f, g)) &= \text{width}(f) + \text{width}(g), \\ \text{mpar}((f, g)) &= \max\{\text{mpar}(f), \text{mpar}(g)\}, \\ \text{mran}_\Omega((f, g))^2 &= \text{mran}_\Omega(f)^2 + \text{mran}_\Omega(g)^2. \end{aligned} \tag{C.1}$$

3. Using the syntax from (3.4a)–(3.4b), it holds  $n_f = p_0(f)$ ,  $n_g = p_0(g)$ ,  $m_f = p_{L+1}(f)$  and  $m_g = p_{L+1}(g)$ .

Let  $N \in \mathbb{N}$ ,  $N \geq 3$ . We extend Definition 35 to parallelize  $N$   $\sigma$ -neural networks  $f_i$ ,  $i = 1, \dots, N$  with equal depth and denote the resulting  $\sigma$ -NN as  $(\{f_i\}_{i=1}^N)$ . It holds that

$$\begin{aligned}
 \text{size} \left( \left( \{f_i\}_{i=1}^N \right) \right) &= \sum_{i=1}^N \text{size}(f_i), & (\text{C.2}) \\
 \text{size}_{\text{in}} \left( \left( \{f_i\}_{i=1}^N \right) \right) &= \sum_{i=1}^N \text{size}_{\text{in}}(f_i), \\
 \text{size}_{\text{out}} \left( \left( \{f_i\}_{i=1}^N \right) \right) &= \sum_{i=1}^N \text{size}_{\text{out}}(f_i), \\
 \text{depth} \left( \left( \{f_i\}_{i=1}^N \right) \right) &= \text{depth}(f_1), \\
 \text{width} \left( \left( \{f_i\}_{i=1}^N \right) \right) &= \sum_{i=1}^N \text{width}(f_i), & (\text{C.3}) \\
 \text{mpar} \left( \left( \{f_i\}_{i=1}^N \right) \right) &= \max_{i=1, \dots, N} \text{mpar}(f_i), \\
 \text{mran}_{\Omega} \left( \left( \{f_i\}_{i=1}^N \right) \right)^2 &= \sum_{i=1}^N \text{mran}_{\Omega}(f_i)^2.
 \end{aligned}$$

Next, we recall the concatenation of NNs, see Petersen and Voigtlaender (2018, Definition 2.2).

**Lemma 36 (Concatenation)** *Let  $q \in \mathbb{N}$ ,  $q \geq 2$  and  $\sigma \in \{\sigma_1, \sigma_q\}$ . Let  $f$  and  $g$  be two  $\sigma$ -NNs. Furthermore, let the output dimension  $m_g$  of  $g$  equal the input dimension  $n_f$  of  $f$ . Then there exists a  $\sigma$ -NN  $f \cdot g$  realizing the composition  $f \circ g : x \mapsto f(g(x))$  of the functions  $f$  and  $g$ . It holds*

$$\begin{aligned}
 \text{depth}(f \cdot g) &= \text{depth}(f) + \text{depth}(g), \\
 \text{width}(f \cdot g) &= \max\{\text{width}(f), \text{width}(g)\}, \\
 \text{mran}_{\Omega}(f \cdot g) &= \text{mran}_{g(\Omega)}(f).
 \end{aligned}$$

There is no simple control over the size and the weight bound of the concatenation  $f \cdot g$  in Definition 36. The reason is that the network  $f \cdot g$  multiplies network weights and biases of the NNs  $f$  and  $g$  at layer  $l = \text{depth}(g) + 1$  (for details see Petersen and Voigtlaender 2018, Definition 2.2). In the following we use *sparse concatenation* to get control over the size and the weights. We first introduce the realization of the identity map and separate the analysis for the  $\sigma_1$ - and  $\sigma_q$ -case. The following lemma is proven in Petersen and Voigtlaender (2018, Remark 2.4).

**Lemma 37 ( $\sigma_1$ -realization of identity map)** *Let  $d \in \mathbb{N}$  and  $L \in \mathbb{N}$ . Then there exists a  $\sigma_1$ -identity network  $\text{Id}_{\mathbb{R}^d}$  of depth  $L$ , which exactly realizes the identity map  $\text{Id}_{\mathbb{R}^d} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ ,  $\mathbf{x} \mapsto \mathbf{x}$ . It holds*

$$\text{size}(\text{Id}_{\mathbb{R}^d}) \leq 2d(L + 1), \tag{C.4}$$

$$\text{width}(\text{Id}_{\mathbb{R}^d}) \leq 2d, \tag{C.5}$$

$$\text{mpar}(\text{Id}_{\mathbb{R}^d}) \leq 1.$$

We proceed with the analogous result for the RePU activation function. The following lemma follows from the construction in Li (2020, Theorem 2.5 part 2) with a parallelization argument.

**Lemma 38 ( $\sigma_q$ -realization of identity map)** *Let  $q \in \mathbb{N}$  with  $q \geq 2$ . Further let  $d \in \mathbb{N}$  and  $L \in \mathbb{N}$  be arbitrary. Then there exists a  $\sigma_q$ -NN  $\text{Id}_{\mathbb{R}^d}$  of depth  $L$ , which exactly realizes the identity map  $\text{Id}_{\mathbb{R}^d}$ . It holds*

$$\begin{aligned} \text{size}(\text{Id}_{\mathbb{R}^d}) &\leq C_q d L, \\ \text{width}(\text{Id}_{\mathbb{R}^d}) &\leq C_q d, \\ \text{mpar}(\text{Id}_{\mathbb{R}^d}) &\leq C_q, \end{aligned} \tag{C.6}$$

where  $C_q$  is independent of  $d$  (but does depend on  $q$ ).

We next define the sparse concatenation of two NNs, see Petersen and Voigtlaender (2018, Definition 2.5).

**Definition 39 ( $\sigma_1$ -sparse concatenation)** *Let  $f$  and  $g$  be two  $\sigma_1$ -NNs. Furthermore, let the output dimension  $m_g$  of  $g$  equal the input dimension  $n_f$  of  $f$ . Then there exists a  $\sigma_1$ -NN  $f \circ g$  with*

$$f \circ g := f \cdot \text{Id}_{\mathbb{R}_{m_g}} \bullet g$$

realizing the composition  $f \circ g : x \mapsto f(g(x))$  of the functions  $f$  and  $g$ .<sup>4</sup> It holds

$$\text{size}(f \circ g) \leq \text{size}(g) + \text{size}_{\text{out}}(g) + \text{size}_{\text{in}}(f) + \text{size}(f) \leq 2\text{size}(f) + 2\text{size}(g), \tag{C.7}$$

$$\text{size}_{\text{in}}(f \circ g) \leq \begin{cases} \text{size}_{\text{in}}(g) & \text{depth}(g) \geq 1, \\ 2\text{size}_{\text{in}}(g) & \text{depth}(g) = 0, \end{cases}$$

$$\text{size}_{\text{out}}(f \circ g) \leq \begin{cases} \text{size}_{\text{out}}(f) & \text{depth}(f) \geq 1, \\ 2\text{size}_{\text{out}}(f) & \text{depth}(f) = 0, \end{cases} \tag{C.8}$$

$$\text{depth}(f \circ g) = \text{depth}(f) + \text{depth}(g) + 1, \tag{C.9}$$

$$\text{width}(f \circ g) \leq 2 \max\{\text{width}(f), \text{width}(g)\}, \tag{C.10}$$

$$\text{mpar}(f \circ g) \leq \max\{\text{mpar}(f), \text{mpar}(g)\}, \tag{C.11}$$

$$\text{mran}_{\Omega}(f \circ g) = B_{g(\Omega)}(f). \tag{C.12}$$

**Proof** The bounds in (C.9), (C.10) and (C.12) follow from Definition 37 with the NN calculus from Definition 36. The bounds on the sizes in (C.7)–(C.8) and the weight and bias bound (C.11) follow from the specific structure of the  $\sigma_1$ -identity network, see Petersen and Voigtlaender (2018, Remark 2.6).  $\blacksquare$

We proceed with the sparse concatenation of  $\sigma_q$ -NNs.

---

4. The symbol  $\circ$  does mean either the functional concatenation of  $f$  and  $g$  or the sparse concatenation of the NN  $f$  and  $g$  (which realizes the function  $f \circ g$ ).

**Definition 40 ( $\sigma_q$ -sparse concatenation)** Let  $q \in \mathbb{N}$  with  $q \geq 2$ . Let  $f$  and  $g$  be two  $\sigma_q$ -NNs. Furthermore, let the output dimension  $m_g$  of  $g$  equal the input dimension  $n_f$  of  $f$ . Then there exists a  $\sigma_q$ -NN  $f \circ g$  with

$$f \circ g := f \cdot \text{Id}_{\mathbb{R}^{m_g}} \cdot g$$

realizing the composition  $f \circ g : x \mapsto f(g(x))$  of the functions  $f$  and  $g$ . It holds

$$\begin{aligned} \text{size}(f \circ g) &\leq \text{size}(g) + (C_q - 1)\text{size}_{\text{out}}(g) + (C_q - 1)\text{size}_{\text{in}}(f) + \text{size}(f) \\ &\leq C_q \text{size}(f) + C_q \text{size}(g), \end{aligned} \tag{C.13}$$

$$\begin{aligned} \text{size}_{\text{in}}(f \circ g) &\leq \begin{cases} \text{size}_{\text{in}}(g) & \text{depth}(g) \geq 1, \\ C_q \text{size}_{\text{in}}(g) & \text{depth}(g) = 0, \end{cases} \\ \text{size}_{\text{out}}(f \circ g) &\leq \begin{cases} \text{size}_{\text{out}}(f) & \text{depth}(f) \geq 1, \\ C_q \text{size}_{\text{out}}(f) & \text{depth}(f) = 0, \end{cases} \end{aligned} \tag{C.14}$$

$$\text{depth}(f \circ g) = \text{depth}(f) + \text{depth}(g) + 1, \tag{C.15}$$

$$\text{width}(f \circ g) \leq C_q \max\{\text{width}(f), \text{width}(g)\}, \tag{C.16}$$

$$\text{mpar}(f \circ g) \leq C_q \max\{\text{mpar}(f), \text{mpar}(g)\}, \tag{C.17}$$

$$\text{mran}_{\Omega}(f \circ g) = B_{g(\Omega)}(f) \tag{C.18}$$

with a constant  $C_q > 1$  depending only on  $q$ .

**Proof** The bounds in (C.15), (C.16) and (C.18) follow from Definition 38 with the NN calculus from Definition 36. The bounds on the sizes in (C.13)–(C.14) and the weight and bias bound in (C.17) hold because of the specific structure of the  $\sigma_q$ -identity network  $\text{Id}_{\mathbb{R}^n}$ , see Li (2020, Eq. 2.57).  $\blacksquare$

Definitions 39 and 40 show that one can control the size, as well as the weights and biases of the concatenation of two NN by inserting one additional identity layer between the two networks. We end this subsection by introducing summation and scalar multiplication networks.

**Definition 41 (Summation networks)** Let  $q \in \mathbb{N}$ ,  $q \geq 2$ ,  $\sigma \in \{\sigma_1, \sigma_q\}$  and  $d, m \in \mathbb{N}$ . Then there exists a  $\sigma$ -NN  $\Sigma_m$  such that for  $x_1, \dots, x_m \in \mathbb{R}^d$

$$\Sigma_m(x_1, \dots, x_m) = \sum_{i=1}^m x_i,$$

with  $\text{depth}(\Sigma_m) = 0$ ,  $\text{width}(\Sigma_m) = md$ ,  $\text{size}(\Sigma_m) = md$  and  $\text{mpar}(\Sigma_m) = 1$ .

**Proof** Set  $\mathbf{w}^1 = (\mathbb{1}_d, \dots, \mathbb{1}_d)$  and  $\mathbf{b}^1 = 0$  with the  $d \times d$  identity matrices  $\mathbb{1}_d$ .  $\blacksquare$

**Definition 42 (Scalar multiplication networks)** Let  $q \in \mathbb{N}$ ,  $q \geq 2$  and  $\sigma \in \{\sigma_1, \sigma_q\}$ . Let  $\alpha \in \mathbb{R}$  and  $d \in \mathbb{N}$ . Then there exists a  $\sigma$ -NN  $SM_{\alpha}$  with

$$SM_{\alpha}(x) = \alpha x, \quad x \in \mathbb{R}^d.$$

Furthermore, there exists a constant  $C_q$  only depending on  $q$  such that

$$\begin{aligned} \text{depth}(SM_\alpha) &\leq C_q \max\{1, \log(|\alpha|)\}, \\ \text{width}(SM_\alpha) &\leq C_q d, \\ \text{size}(SM_\alpha) &\leq C_q d \max\{1, \log(|\alpha|)\}, \\ \text{mpar}(SM_\alpha) &\leq C_q. \end{aligned}$$

**Proof** For the proof of the  $\sigma_1$ -case, see Elbrachter et al. (2021, Lemma A.1). We prove the RePU case in the following.

Without loss of generality we can choose  $\alpha > 1$ . For  $\alpha < 0$  we set  $SM_\alpha = -SM_{-\alpha}$ . For  $0 < \alpha < 1$  we set  $SM_\alpha = \alpha \text{Id}_{\mathbb{R}^d}$ , where we directly multiply the weights and biases of the identity network (with depth  $L = 1$ ) with  $\alpha$ .

Therefore let  $\alpha > 1$ . Let  $K$  be the maximum integer smaller than  $\log_2(\alpha)$ , and set  $\tilde{\alpha} = 2^{-K+1}\alpha < 1$ . Furthermore, set  $A_1 = (\text{Id}_{\mathbb{R}^d}, \text{Id}_{\mathbb{R}^d})$  and  $A_2 = \Sigma_2$  with the one-layered identity network  $\text{Id}_{\mathbb{R}^d}$  and the summation network  $\Sigma_2$  from Definition 41. We notice that

$$A_2 \cdot A_1 x = 2x \quad \forall x \in \mathbb{R}^d.$$

Using the bounds for  $\Sigma_2$  from Definition 41 and  $\text{Id}_{\mathbb{R}^d}$  from Definition 38, we have  $\text{depth}(A_2 \cdot A_1) = 1$ ,  $\text{width}(A_2 \cdot A_1) \leq C_q d$ ,  $\text{size}(A_2 \cdot A_1) \leq C_q d$  and  $\text{mpar}(A_2 \cdot A_1) \leq C_q$ . Setting  $A_{2k+1} = A_1$ ,  $A_{2k+2} = A_2$  for  $k = 1, \dots, K$  and  $A_{2K+3} = \tilde{\alpha} \text{Id}_{\mathbb{R}^d}$ , we get

$$A_{2K+3} \circ A_{2K+2} \cdot A_{2K+1} \circ A_{2K} \cdot A_{2K-1} \circ \dots \circ A_2 \cdot A_1 x = \alpha x \quad \forall x \in \mathbb{R}^d.$$

Applying the  $\sigma_q$ -NN calculus for concatenation (Definition 36) and sparse concatenation (Definition 40) gives the desired bounds for  $SM_\alpha$ .  $\blacksquare$

## C.2 Neural Network Approximation Theory

In this subsection, we summarize approximation results for ReLU and RePU neural networks. In recent years, the expressivity and approximation properties of neural network architectures have been extensively studied in the literature (Mhaskar, 1996; Pinkus, 1999; Yarotsky, 2017; Poggio et al., 2017; Petersen and Voigtlaender, 2018; Elbrachter et al., 2021; Opschoor et al., 2022a; De Ryck et al., 2021). However, with few exceptions (Schmidt-Hieber, 2020a; De Ryck et al., 2021; Elbrachter et al., 2021), most of these works do not provide bounds on the size of the weights, which are crucial for controlling the entropy. Therefore, we revisit some of these arguments to provide complete proofs of our results.

We start with the well-known result that ReLU-NNs can approximate the multiplication map exponentially fast. The following proposition was shown in Elbrachter et al. (2021, Proposition III.3).

**Proposition 43 ( $\sigma_1$ -NN approximation of multiplication)** *Let  $D \in \mathbb{R}$ ,  $D \geq 1$  and  $\delta \in (0, 1/2)$ . Then there exists a  $\sigma_1$ -NN  $\tilde{\times}_{\delta, D} : [-D, D]^2 \rightarrow \mathbb{R}$  satisfying*

$$\sup_{x, y \in [-D, D]} |\tilde{\times}_{\delta, D}(x, y) - xy| \leq \delta.$$

Furthermore, there exists a constant  $C$ , independent of  $D$  and  $\delta$ , such that

$$\text{depth}(\tilde{\times}_{\delta,D}) \leq C(\log(D) + \log(\delta^{-1})), \quad (\text{C.19})$$

$$\text{width}(\tilde{\times}_{\delta,D}) \leq 5, \quad (\text{C.20})$$

$$\text{size}(\tilde{\times}_{\delta,D}) \leq C(\log(D) + \log(\delta^{-1})), \quad (\text{C.21})$$

$$\text{mpar}(\tilde{\times}_{\delta,D}) \leq 1. \quad (\text{C.22})$$

The next proposition shows that the multiplication map can be exactly realized by a  $\sigma_q$ -NN, which follows directly from Li (2020, Theorem 2.5 and Eq. 2.59).

**Proposition 44 ( $\sigma_q$ -NN approximation of multiplication)** *Let  $q \in \mathbb{N}$ ,  $q \geq 2$ . There exists a  $\sigma_q$ -NN  $\tilde{\times} : \mathbb{R}^2 \rightarrow \mathbb{R}$  with  $\text{depth}(\tilde{\times}) = 1$  exactly realizing the multiplication of two numbers, i.e.*

$$\tilde{\times}(x, y) = xy \quad \forall x, y \in \mathbb{R}.$$

We can extend the above results to the multiplication of  $N$  numbers, see Opschoor et al. (2022a, Proposition 2.6).

**Proposition 45 ( $\sigma_1$ -NN multiplication of  $N$  numbers)** *Let  $N \in \mathbb{N}$  with  $N \geq 2$ . Furthermore, let  $D \in \mathbb{R}$ ,  $D \geq 1$  and  $\delta \in (0, 1/2)$ . Then there exists a  $\sigma_1$ -NN  $\tilde{\prod}_{\delta,D} : [-D, D]^N \rightarrow \mathbb{R}$  such that*

$$\sup_{(y_i)_{i=1}^N \in [-D, D]^N} \left| \prod_{j=1}^N y_j - \tilde{\prod}_{\delta,D}(y_1, \dots, y_N) \right| \leq \delta. \quad (\text{C.23})$$

Furthermore, there exists a constant  $C$  independent of  $N$ ,  $\delta$  and  $D$  such that

$$\text{depth}\left(\tilde{\prod}_{\delta,D}\right) \leq C \log(N) (\log(N) + N \log(D) + \log(\delta^{-1})),$$

$$\text{width}\left(\tilde{\prod}_{\delta,D}\right) \leq 5N, \quad (\text{C.24})$$

$$\text{size}\left(\tilde{\prod}_{\delta,D}\right) \leq CN (\log(N) + N \log(D) + \log(\delta^{-1})), \quad (\text{C.25})$$

$$\text{mpar}\left(\tilde{\prod}_{\delta,D}\right) \leq 1.$$

**Proof** Analogous to Opschoor et al. (2022a, Proposition 2.6) we construct  $\tilde{\prod}_{\delta,D}$  as a binary tree of  $\tilde{\times}_{\delta,D}$ -networks from Proposition 43. We modify the proof of Opschoor et al. (2022a) to get a construction with bounded weights.

Define  $\tilde{N} := \min\{2^k : k \in \mathbb{N}, 2^k \geq N\}$ . We now consider the multiplication of  $\tilde{N}$  numbers with  $y_{N+1}, \dots, y_{\tilde{N}} := 1$ . This can be implemented by a zero-layer network with

$$w_{i,j}^1 = \begin{cases} 1 & i = j \leq N, \\ 0 & \text{otherwise,} \end{cases}$$

$$b_j^1 = \begin{cases} 0 & j \leq N, \\ 1 & N < j \leq \tilde{N}. \end{cases}$$

For  $l = 0, \dots, \log_2 \tilde{N} - 1$  we define the mapping  $R^l$

$$R^l : [-D_l, D_l]^{\tilde{N}} \mapsto [-D_{l+1}, D_{l+1}]^{\tilde{N}},$$

$$R^l(y_1^l, \dots, y_{2^{\log_2(\tilde{N})-l}}^l) := \left( \tilde{\times}_{\delta', D^l}(y_1^l, y_2^l), \dots, \tilde{\times}_{\delta', D^l}(y_{2^{\log_2(\tilde{N})-l-1}}^l, y_{2^{\log_2(\tilde{N})-l}}^l) \right) \quad (\text{C.26})$$

with  $\delta' := \delta/(\tilde{N}^2 D^{2\tilde{N}})$  and  $D_l := 2^l D^{2^l}$ . We now set

$$\widetilde{\prod}_{\delta, D} := R^{\log_2(\tilde{N})-1} \circ \dots \circ R^0. \quad (\text{C.27})$$

Eq. (C.27) shows that the map  $R^l$  describes the multiplications on level  $l$  of the binary tree  $\widetilde{\prod}_{\delta, D}$ . In order for (C.27) to be well-defined, we have to show that the outputs of the NN  $R^l$  are admissible inputs for the NN  $R^{l+1}$ .

We therefore denote with  $y_j^l, j = 1, \dots, 2^{\log_2(\tilde{N})-l}, l = 1, \dots, \log_2(\tilde{N}) - 1$ , the output of the network  $R^{l-1} \circ \dots \circ R^0$  applied to the input  $y_k^0 = y_k, k = 1, \dots, \tilde{N}$ . Then we have to show  $|y_j^l| \leq D_l$  for  $l = 0, \dots, \log_2(\tilde{N}) - 1$  and  $j = 1, \dots, 2^{\log_2(\tilde{N})-l}$ . We will show this claim by induction. For  $l = 0$  it holds  $|y_j^l| \leq D = D_0$ . Now assume  $|y_j^l| \leq D_l$  for arbitrary but fixed  $l \in \{0, \dots, \log_2(\tilde{N}) - 2\}$  and all  $j = 1, \dots, 2^{\log_2(\tilde{N})-l}$ . Then it holds

$$\begin{aligned} |y_j^{l+1}| &= |\tilde{\times}_{\delta', D^l}(y_{2j-1}^l, y_{2j}^l)| \\ &= |y_{2j-1}^l \cdot y_{2j}^l + \delta'| \leq D_l^2 + \underbrace{\delta'}_{\leq 1 \leq D_l^2} \leq 2D_l^2 = D_{l+1} \end{aligned}$$

for all  $j = 1, \dots, 2^{\log_2(\tilde{N})-(l+1)}$ , which shows the claim. We proceed by showing the error bound in (C.23). Therefore define

$$z_j^l := \prod_{k=1}^{2^l} y_{k+2^l(j-1)}$$

for  $l = 0, \dots, \log_2(\tilde{N})$  and  $j = 1, \dots, 2^{\log_2(\tilde{N})-l}$ . The quantities  $z_j^l$  describe the exact computations up to level  $l$  of the binary tree, i.e. the output of level  $l-1$ , if one uses standard multiplication instead of the multiplication networks  $\tilde{\times}$  in the first  $l-1$  levels. We now prove

$$|y_j^l - z_j^l| \leq 4^l D^{2^{l+1}} \delta', \quad j = 1, \dots, 2^{\log_2(\tilde{N})-l} \quad (\text{C.28})$$

by induction over  $l = 0, \dots, \log_2 \tilde{N}$ . Inserting  $l = \log_2(\tilde{N})$  then shows the error bound in (C.23) using the definition of  $\delta'$ .

We have  $y_j^0 = y_j = z_j^0$  for  $j = 1, \dots, \tilde{N}$ , therefore (C.28) holds for  $l = 0$ . Now assume (C.28) for arbitrary but fixed  $l \in \{0, \dots, \log_2(\tilde{N}) - 1\}$ . For  $j = 1, \dots, 2^{\log_2(\tilde{N}) - (l+1)}$  it holds

$$\begin{aligned}
 |y_j^{l+1} - z_j^{l+1}| &= \left| \tilde{\times}_{D_l, \delta'}(y_{2j-1}^l, y_{2j}^l) - z_{2j-1}^l z_{2j}^l \right| \\
 &= \left| y_{2j-1}^l y_{2j}^l + \delta' - z_{2j-1}^l z_{2j}^l \right| \\
 &= \left| (y_{2j-1}^l - z_{2j-1}^l + z_{2j-1}^l) \cdot (y_{2j}^l - z_{2j}^l + z_{2j}^l) + \delta' - z_{2j-1}^l z_{2j}^l \right| \\
 &= \left| (y_{2j-1}^l - z_{2j-1}^l) \cdot (y_{2j}^l - z_{2j}^l) + (y_{2j-1}^l - z_{2j-1}^l) z_{2j}^l \right. \\
 &\quad \left. + (y_{2j}^l - z_{2j}^l) z_{2j-1}^l + \delta' \right| \\
 &\leq \underbrace{\left| y_{2j-1}^l - z_{2j-1}^l \right|}_{\text{use(C.28)}} \cdot \underbrace{\left| y_{2j}^l - z_{2j}^l \right|}_{\leq 1} + \underbrace{\left| y_{2j-1}^l - z_{2j-1}^l \right|}_{\text{use(C.28)}} \cdot \underbrace{\left| z_{2j}^l \right|}_{\leq D^{2^l}} \\
 &\quad + \underbrace{\left| y_{2j}^l - z_{2j}^l \right|}_{\text{use(C.28)}} \cdot \underbrace{\left| z_{2j-1}^l \right|}_{\leq D^{2^l}} + |\delta'| \\
 &\leq \delta' \left( 4^l D^{2^{l+1}} \left( 1 + 2D^{2^l} \right) + 1 \right) \\
 &\leq 4 \cdot 4^l D^{2^{l+1}} D^{2^{l+1}} \delta' \leq 4^{l+1} D^{2^{l+2}} \delta',
 \end{aligned}$$

which shows (C.28) for  $l + 1$  and therefore the claim.

We proceed by calculating the depth of  $\tilde{\Pi}_{\delta, D}$ . Since  $\tilde{\Pi}_{\delta, D}$  concatenates the maps  $\tilde{\times}_{\delta', D_l}$ , we can repeatedly use (C.9) and get

$$\text{depth} \left( \tilde{\Pi}_{\delta, D} \right) \leq \sum_{j=0}^{\log_2(\tilde{N})-1} \text{depth}(\tilde{\times}_{\delta', D_j}) + \log_2(\tilde{N}) - 1.$$

We use the depth bound for  $\tilde{\times}$  from (C.19) and calculate

$$\begin{aligned}
 \text{depth} \left( \tilde{\Pi}_{\delta, D} \right) &\leq C \sum_{j=0}^{\log_2(\tilde{N})-1} \log(D_j \delta'^{-1}) + \log_2(\tilde{N}) \\
 &= C \sum_{j=0}^{\log_2(\tilde{N})-1} \log\left(2^j D^{2^j} \delta'^{-1}\right) + \log_2(\tilde{N}) \\
 &= C \log \left( \prod_{j=0}^{\log_2(\tilde{N})-1} 2^j D^{2^j} \delta'^{-1} \right) + \log_2(\tilde{N}) \\
 &\leq C \log \left( 2^{\frac{\log_2(\tilde{N}) \cdot (\log_2(\tilde{N})-1)}{2}} D^{2^{\log_2(\tilde{N})}} (\delta')^{-\log_2(\tilde{N})} \right) + \log_2(\tilde{N}) \\
 &\leq C \log \left( 2^{(\log_2(\tilde{N}))^2} D^{\tilde{N}} \tilde{N}^{2 \log_2(\tilde{N})} D^{2 \tilde{N} \log_2(\tilde{N})} \delta^{-\log_2(\tilde{N})} \right) + \log_2(\tilde{N}) \\
 &\leq C \log \left( 2^{(\log_2(\tilde{N}))^2} \tilde{N}^{2 \log_2(\tilde{N})} D^{3 \tilde{N} \log_2(\tilde{N})} \delta^{-\log_2(\tilde{N})} \right) + \log_2(\tilde{N}) \\
 &\leq C \log(N) \left( \log(N) + N \log(D) + \log(\delta^{-1}) \right). \tag{C.29}
 \end{aligned}$$

The constant  $C$  changes from line to line in (C.29).

For a bound on the width we use the fact that  $\widetilde{\prod}_{\delta, D}$  is a parallelization of at most  $\tilde{N}/2 \leq N$  networks  $\tilde{\times}_{\delta', D^l}$  in each layer  $l \in \{1, \dots, \log_2(\tilde{N})\}$ . With (C.3) and the width bound of  $\tilde{\times}$  in (C.20) it holds

$$\text{width} \left( \widetilde{\prod}_{\delta, D^l} \right) \leq N \text{width} (\tilde{\times}_{\delta', D^l}) \leq 5N.$$

For a bound on the size  $\text{size}(\widetilde{\prod}_{\delta, D})$  we observe that level  $l$ ,  $l = 0, \dots, \log_2(\tilde{N}) - 1$ , of the binary tree  $\widetilde{\prod}$  consists of  $2^{\log_2(\tilde{N})-l-1}$  product networks  $\tilde{\times}_{\delta', D^l}$ . We calculate

$$\begin{aligned} \text{size} \left( \widetilde{\prod}_{\delta, D} \right) &\leq \sum_{l=0}^{\log_2(\tilde{N})-1} 2^{\log_2(\tilde{N})-l-1} (s_{in}(\tilde{\times}_{\delta', D^l}) + s_{out}(\tilde{\times}_{\delta', D^l}) + \text{size}(\tilde{\times}_{\delta', D^l})) \\ &\leq \sum_{l=0}^{\log_2(\tilde{N})-1} 2^{\log_2(\tilde{N})-l-1} 3C \log(D_l \delta'^{-1}) \\ &\leq C \sum_{l=0}^{\log_2(\tilde{N})-1} 2^{\log_2(\tilde{N})-l-1} \log(2^l D^{2^l} \tilde{N}^2 D^{2\tilde{N}} \delta^{-1}) \\ &\leq C \sum_{l=0}^{\log_2(\tilde{N})-1} 2^{\log_2(\tilde{N})-l-1} (l + 2^l \log(D) + \log(\tilde{N}) + \tilde{N} \log(D) + \log(\delta^{-1})) \\ &\leq C (\tilde{N} \log_2(\tilde{N}) + \tilde{N} \log_2(\tilde{N}) \log(D) + \tilde{N} \log(\tilde{N}) + \tilde{N}^2 \log(D) + \tilde{N} \log(\delta^{-1})) \\ &\leq CN (\log(N) + N \log(D) + \log(\delta^{-1})). \end{aligned} \tag{C.30}$$

In (C.30) we used (C.7) to bound the size of a sparse concatenation and (C.21) for the size of the product network  $\tilde{\times}_{\delta, D}$ .

For the bound on the weights and biases, we get  $\text{mpar}(\widetilde{\prod}_{\delta, D}) \leq 1$  because of  $\text{mpar}(\tilde{\times}_{\delta, D}) \leq 1$ , see (C.22), and the NN calculus for sparse concatenation in (C.11) and parallelization in (C.1).  $\blacksquare$

We continue with the RePU-case.

**Proposition 46 ( $\sigma_q$ -NN of multiplication of  $N$  numbers)** *Let  $N, q \in \mathbb{N}$  with  $N, q \geq 2$ . Then there exists a  $\sigma_q$ -NN  $\widetilde{\prod}: \mathbb{R}^N \rightarrow \mathbb{R}$  such that*

$$\widetilde{\prod}(y_1, \dots, y_N) = \prod_{j=1}^N y_j.$$

Furthermore, there exists a constant  $C_q$  independent of  $N$  such that

$$\text{depth} \left( \widetilde{\Pi}_{\delta, D} \right) \leq C_q \log(N), \quad (\text{C.31})$$

$$\text{width} \left( \widetilde{\Pi}_{\delta, D} \right) \leq C_q N, \quad (\text{C.32})$$

$$\text{size} \left( \widetilde{\Pi}_{\delta, D} \right) \leq C_q N, \quad (\text{C.33})$$

$$\text{mpar} \left( \widetilde{\Pi}_{\delta, D} \right) \leq C_q. \quad (\text{C.34})$$

**Proof** The construction is similar to the ReLU case. We define  $\widetilde{\Pi}$  as a binary tree of product networks  $\tilde{\times}$ , see (C.26) and (C.27). The binary tree has a maximum of  $2N$  binary networks  $\tilde{\times}$ , a maximum height of  $\log_2(2N)$  and a maximum width of  $N$ . Therefore (C.31)–(C.34) follow with the NN calculus rules from Definition 40.  $\blacksquare$

We proceed and state the approximation results for univariate polynomials. We start with the ReLU case. The following proposition was shown in Elbrachter et al. (2021, Proposition III.5).

**Proposition 47 ( $\sigma_1$ -NN approximation of polynomials)** *Let  $m \in \mathbb{N}$  and  $a = (a_i)_{i=0}^m \in \mathbb{R}^{m+1}$ . Further let  $D \in \mathbb{R}$ ,  $D \geq 1$  and  $\delta \in (0, 1/2)$ . Define  $a_\infty = \max\{1, \|a\|_\infty\}$ . Then there exists a  $\sigma_1$ -NN  $\tilde{p}_{\delta, D} : [-D, D] \rightarrow \mathbb{R}$  satisfying*

$$\sup_{x \in [-D, D]} \left| \tilde{p}_{\delta, D}(x) - \sum_{i=0}^m a_i x^i \right| \leq \delta.$$

Furthermore, there exists a constant  $C$  independent of  $m$ ,  $a_i$ ,  $D$  and  $\delta$  such that

$$\begin{aligned} \text{depth}(\tilde{p}_{\delta, D}) &\leq Cm (m \log(D) + \log(\delta^{-1}) + \log(m) + \log(a_\infty)), \\ \text{width}(\tilde{p}_{\delta, D}) &\leq 9, \\ \text{size}(\tilde{p}_{\delta, D}) &\leq Cm (m \log(D) + \log(\delta^{-1}) + \log(m) + \log(a_\infty)), \\ \text{mpar}(\tilde{p}_{\delta, D}) &\leq 1. \end{aligned}$$

In the RePU-case we get the well-known result that polynomials can be exactly realized by  $\sigma_q$ -NNs, see Li et al. (2019).

**Proposition 48 ( $\sigma_q$ -NN realization of polynomials)** *Let  $m, q \in \mathbb{N}$ ,  $q \geq 2$  and  $a = (a_i)_{i=0}^m \in \mathbb{R}^{m+1}$ . Set  $a_\infty = \max\{1, \max_{i=0, \dots, m} a_i\}$ . Then there exists a  $\sigma_q$ -NN  $\tilde{p} : \mathbb{R} \rightarrow \mathbb{R}$  satisfying*

$$\tilde{p}(x) = \sum_{i=0}^m a_i x^i \quad \forall x \in \mathbb{R}.$$

Furthermore, there exists a constant  $C_q$  only depending on  $q$  such that

$$\begin{aligned} \text{depth}(\tilde{p}) &\leq C_q (\log(a_\infty) + m), \\ \text{width}(\tilde{p}) &\leq C_q, \\ \text{size}(\tilde{p}) &\leq C_q (\log(a_\infty) + m), \\ \text{mpar}(\tilde{p}) &\leq C_q. \end{aligned}$$

**Proof** We use Horner's method for polynomial evaluation and write

$$\sum_{i=0}^m a_i x^i = a_\infty \left( \frac{a_0}{a_\infty} + x \left( \frac{a_1}{a_\infty} + \cdots + x \left( \frac{a_{m-1}}{a_\infty} + x \frac{a_m}{a_\infty} \right) \cdots \right) \right). \quad (\text{C.35})$$

Following (C.35), we build  $\tilde{p}$  via

$$\tilde{p} = SM_{a_\infty} \circ \Sigma_2 \left( \frac{a_0}{a_\infty}, \tilde{\times} \left( \text{Id}_{\mathbb{R}}, \Sigma_2 \left( \frac{a_1}{a_\infty}, \dots, SM_{a_m a_\infty^{-1}}(\text{Id}_{\mathbb{R}}) \right) \cdots \right) \right).$$

The bounds for  $\tilde{p}$  follow from the respective bounds for  $\Sigma_2$  from Definition 41,  $\text{Id}_{\mathbb{R}}$  from Lemma 38 and  $SM_\alpha$  from Definition 42.  $\blacksquare$

We now use Propositions 47 and 48 to get an approximation result for univariate Legendre polynomials.

**Corollary 49 ( $\sigma_1$ -NN approximation of  $L_j$ )** *Let  $j \in \mathbb{N}_0$  and  $\delta \in (0, 1/2)$ . Then there exists a  $\sigma_1$ -NN  $\tilde{L}_{j,\delta} : [-1, 1] \rightarrow \mathbb{R}$  with*

$$\sup_{x \in [-1, 1]} |\tilde{L}_{j,\delta}(x) - L_j(x)| \leq \delta.$$

Furthermore, there exists a constant  $C$  such that it holds

$$\text{depth}(\tilde{L}_{j,\delta}) \leq Cj(j + \log(\delta^{-1})), \quad (\text{C.36})$$

$$\text{width}(\tilde{L}_{j,\delta}) \leq 9, \quad (\text{C.37})$$

$$\text{size}(\tilde{L}_{j,\delta}) \leq Cj(j + \log(\delta^{-1})), \quad (\text{C.38})$$

$$\text{mpar}(\tilde{L}_{j,\delta}) \leq 1.$$

**Proof** For  $j \in \mathbb{N}$ ,  $l \in \mathbb{N}_0$ ,  $l \leq j$ , denote the coefficients of  $L_j$  with  $c_l^j$ . In Opschoor et al. (2020, Eq. (4.17)) the bound  $\sum_{l=0}^j |c_l^j| \leq 4^j$  is derived. With  $c^j = (c_l^j)_{l=0}^j$  it holds  $\|c^j\|_\infty \leq \sum_{l=0}^j |c_l^j| \leq 4^j$ . The result now follows with Proposition 47.  $\blacksquare$

We continue with the  $\sigma_q$ -case.

**Corollary 50 ( $\sigma_q$ -NN approximation of  $L_j$ )** *Let  $j \in \mathbb{N}_0$ . Then there exists a  $\sigma_q$ -NN  $\tilde{L}_j : \mathbb{R} \rightarrow \mathbb{R}$  with*

$$\tilde{L}_j(x) = L_j(x) \quad \forall x \in \mathbb{R}.$$

*Furthermore, there exists a constant  $C_q$  only depending on  $q$  such that it holds*

$$\text{depth}(\tilde{L}_j) \leq C_q j, \tag{C.39}$$

$$\text{width}(\tilde{L}_j) \leq C_q, \tag{C.40}$$

$$\text{size}(\tilde{L}_j) \leq C_q j, \tag{C.41}$$

$$\text{mpar}(\tilde{L}_j) \leq C_q.$$

**Proof** The bounds follow similar to the  $\sigma_1$ -case using Proposition 48. ■

## Appendix D. Proofs of Section 3

### D.1 Proof of Proposition 17

We proceed analogously to the proof of Opschoor et al. (2022a, Proposition 2.13). We define  $f_{\Lambda, \delta}$  as a composition of two subnetworks,  $f_{\Lambda, \delta} := f_{\Lambda, \delta}^{(1)} \circ f_{\Lambda, \delta}^{(2)}$ . Corollary 49 ensures the existence of (arbitrarily good)  $\sigma_1$ -NN approximations of univariate Legendre polynomials. All relevant approximations are evaluated, in parallel, by subnetwork  $f_{\Lambda, \delta}^{(2)}$ , i.e.

$$f_{\Lambda, \delta}^{(2)}(\mathbf{y}) := \left( \left\{ \text{Id}_{\mathbb{R}} \circ \tilde{L}_{\nu_j, \delta'}(y_j) \right\}_{(j, \nu_j) \in T} \right), \tag{D.1}$$

where we used

$$\begin{aligned} T &:= \{(j, \nu_j) \in \mathbb{N}^2 : \boldsymbol{\nu} \in \Lambda, j \in \text{supp } \boldsymbol{\nu}\}, \\ \delta' &:= (2d(\Lambda))^{-1} (2m(\Lambda) + 2)^{-d(\Lambda)+1} \delta \end{aligned} \tag{D.2}$$

and  $\mathbf{y} = (y_j)_{(j, \nu_j) \in T}$ . In (D.1) the big round brackets denote a parallelization and we use the identity networks to synchronize the depth. The subnetwork  $f_{\Lambda, \delta}^{(1)}$  takes the output of  $f_{\Lambda, \delta}^{(2)}$  as input and computes, in parallel, tensorized Legendre polynomials using the multiplication networks  $\tilde{\Pi}_{\cdot, \cdot}$  introduced in Proposition 45. With  $M_{\boldsymbol{\nu}} := 2|\boldsymbol{\nu}|_1 + 2$  we define

$$\begin{aligned} f_{\Lambda, \delta}^{(1)}((z_k)_{k \leq |T|}) &= f_{\Lambda, \delta}^{(1)}\left(f_{\Lambda, \delta}^{(2)}(\mathbf{y})\right) \\ &:= \left( \left\{ \text{Id}_{\mathbb{R}} \circ \tilde{\Pi}_{\delta/2, M_{\boldsymbol{\nu}}} \left( \left\{ \text{Id}_{\mathbb{R}} \circ \tilde{L}_{\nu_j, \delta'}(y_j) \right\}_{j \in \text{supp } \boldsymbol{\nu}} \right) \right\}_{\boldsymbol{\nu} \in \Lambda} \right). \end{aligned} \tag{D.3}$$

The multiplication networks in (D.3) are well-defined, since

$$\sup_{y_j \in [-1, 1]} |\tilde{L}_{\nu_j, \delta'}(y_j)| \leq 2\nu_j + 2 \leq 2|\boldsymbol{\nu}|_1 + 2 = M_{\boldsymbol{\nu}}, \tag{D.4}$$

where we used (3.10) and  $\delta' < 1$ .

We will first show the error bound in (3.12). Let  $\nu \in \Lambda$  be arbitrary. We use the shorthand notation  $\|\cdot\| := \|\cdot\|_{L^\infty([-1,1]^{|T|})}$  and calculate

$$\begin{aligned}
 & \left\| L_\nu - \tilde{L}_{\nu, \delta} \right\| \\
 & \leq \left\| L_\nu - \prod_{j \in \text{supp } \nu} \tilde{L}_{\nu_j, \delta'} \right\| + \left\| \prod_{j \in \text{supp } \nu} \tilde{L}_{\nu_j, \delta'} - \widetilde{\prod}_{\delta/2, M_\nu} \left( \left\{ \tilde{L}_{\nu_j, \delta'} \right\}_{j \in \text{supp } \nu} \right) \right\| \\
 & \leq \sum_{k \in \text{supp } \nu} \left\| \prod_{\substack{j \in \text{supp } \nu: \\ j < k}} \tilde{L}_{\nu_j, \delta'} \right\| \cdot \left\| L_{\nu_k} - \tilde{L}_{\nu_k, \delta'} \right\| \cdot \left\| \prod_{\substack{j \in \text{supp } \nu: \\ j > k}} L_{\nu_j} \right\| + \frac{\delta}{2} \\
 & \leq d(\Lambda) M_\nu^{d(\Lambda)-1} \delta' + \frac{\delta}{2} \leq \left( \frac{M_\nu}{2m(\Lambda) + 2} \right)^{d(\Lambda)-1} \frac{\delta}{2} + \frac{\delta}{2} \leq \delta,
 \end{aligned}$$

where we used (D.4),  $M_\nu \leq 2m(\Lambda) + 2$  and the definition of  $\delta'$ .

We proceed and calculate the depth  $L$  of  $f_{\Lambda, \delta}$ . Since  $f_{\Lambda, \delta} = f_{\Lambda, \delta}^{(1)} \circ f_{\Lambda, \delta}^{(2)}$ , it holds  $\text{depth}(f_{\Lambda, \delta}) \leq \text{depth}(f_{\Lambda, \delta}^{(1)}) + \text{depth}(f_{\Lambda, \delta}^{(2)}) + 1$ , see (C.9). We start with a depth bound of  $f_{\Lambda, \delta}^{(2)}$ . Denoting by  $C$  a universal multiplicative constant that is allowed to change from line to line, it holds that

$$\begin{aligned}
 \text{depth} \left( f_{\Lambda, \delta}^{(2)} \right) &= 1 + \max_{\substack{\nu \in \Lambda \\ j \in \text{supp } \nu}} \text{depth} \left( \tilde{L}_{\nu_j, \delta'} \right) \\
 &\leq C \max_{\substack{\nu \in \Lambda \\ j \in \text{supp } \nu}} \nu_j (\nu_j + \log(\delta'^{-1})) \\
 &\leq Cm(\Lambda) (m(\Lambda) + \log(\delta'^{-1})) \\
 &\leq Cm(\Lambda) (\log(d(\Lambda)) + d(\Lambda) \log(m(\Lambda)) + m(\Lambda) + \log(\delta^{-1})) \\
 &\leq Cm(\Lambda) (\log(d(\Lambda)) + d(\Lambda) \log(m(\Lambda)) + m(\Lambda) + \log(\delta^{-1})). \quad (\text{D.5})
 \end{aligned}$$

In (D.5) we used the depth bound for univariate Legendre polynomials, (C.36), at the first inequality. Furthermore, we used  $\nu_j \leq m(\Lambda)$ . For the depth of  $f_{\Lambda, \delta}^{(1)}$  it holds

$$\begin{aligned}
 \text{depth} \left( f_{\Lambda, \delta}^{(1)} \right) &= 1 + \max_{\nu \in \Lambda} \text{depth} \left( \widetilde{\prod}_{\delta/2, M_\nu} \right) \\
 &\leq 1 + C \max_{\nu \in \Lambda} \log(|\text{supp } \nu|) (\log(|\text{supp } \nu|) + |\text{supp } \nu| \log(M_\nu) + \log(\delta^{-1})) \\
 &\leq 1 + C \log(d(\Lambda)) (\log(d(\Lambda)) + d(\Lambda) \log(m(\Lambda)) + \log(\delta^{-1})), \quad (\text{D.6})
 \end{aligned}$$

where we used  $|\text{supp } \nu| \leq d(\Lambda)$  for all  $\nu \in \Lambda$ ,  $M_\nu \leq 4m(\Lambda)$  and the depth bound for  $\sigma_1$ -multiplication networks from Proposition 45. Combining the two depth bounds (D.5) and

(D.6), we get

$$\begin{aligned}
 \text{depth}(f_{\Lambda,\delta}) &= 1 + \text{depth}\left(f_{\Lambda,\delta}^{(1)}\right) + \text{depth}\left(f_{\Lambda,\delta}^{(2)}\right) \\
 &\leq Cm(\Lambda) (\log(d(\Lambda)) + d(\Lambda) \log(m(\Lambda)) + m(\Lambda) + \log(\delta^{-1})) \\
 &\quad + C \log(d(\Lambda)) (\log(d(\Lambda)) + d(\Lambda) \log(m(\Lambda)) + \log(\delta^{-1})) \\
 &\leq C \left[ \log(d(\Lambda))d(\Lambda) \log(m(\Lambda))m(\Lambda) + m(\Lambda)^2 + \log(\delta^{-1})(\log(d(\Lambda)) + m(\Lambda)) \right].
 \end{aligned}$$

For the width  $\text{width}(f_{\Lambda,\delta})$  we use  $\text{width}(f_{\Lambda,\delta}) \leq 2 \max\{\text{width}(f_{\Lambda,\delta}^{(1)}), \text{width}(f_{\Lambda,\delta}^{(2)})\}$ , see (C.10). This leaves us to calculate  $\text{width}(f_{\Lambda,\delta}^{(2)})$  and  $\text{width}(f_{\Lambda,\delta}^{(1)})$ . It holds

$$\begin{aligned}
 \text{width}\left(f_{\Lambda,\delta}^{(2)}\right) &\leq \sum_{(j,\nu_j) \in T} \text{width}\left(\text{Id}_{\mathbb{R}} \circ \tilde{L}_{\nu_j,\delta'}\right) \\
 &\leq 2 \sum_{(j,\nu_j) \in T} \text{width}\left(\tilde{L}_{\nu_j,\delta'}\right) \leq 18|T|, \tag{D.7}
 \end{aligned}$$

where we used (C.5) for the width of the  $\sigma_1$ -identity network and (C.37) for the width of  $\tilde{L}_{\nu_j,\delta'}$ . For  $\text{width}(f_{\Lambda,\delta}^{(1)})$  it holds

$$\begin{aligned}
 \text{width}\left(f_{\Lambda,\delta}^{(1)}\right) &\leq \sum_{\nu \in \Lambda} \text{width}\left(\text{Id}_{\mathbb{R}} \circ \widetilde{\Pi}_{\delta/2, M_\nu}\right) \\
 &\leq 2 \sum_{\nu \in \Lambda} \text{width}\left(\widetilde{\Pi}_{\delta/2, M_\nu}\right) \\
 &\leq \sum_{\nu \in \Lambda} 10d(\Lambda) = 10|\Lambda|d(\Lambda), \tag{D.8}
 \end{aligned}$$

again using (C.5) and (C.24) for the width of the multiplication network  $\widetilde{\Pi}$ . Combining (D.7) and (D.8) gives

$$\text{width}(f_{\Lambda,\delta}) \leq 36|\Lambda|d(\Lambda),$$

where  $|T| \leq |\Lambda|d(\Lambda)$  was used.

To estimate  $\text{size}(f_{\Lambda,\delta})$ , we use (C.7) and find  $\text{size}(f_{\Lambda,\delta}) \leq 2\text{size}(f_{\Lambda,\delta}^{(1)}) + 2\text{size}(f_{\Lambda,\delta}^{(2)})$ . We calculate

$$\begin{aligned}
 \text{size}\left(f_{\Lambda,\delta}^{(2)}\right) &= \text{size}\left(\left\{\text{Id}_{\mathbb{R}} \circ \tilde{L}_{\nu_j,\delta'}(y_j)\right\}_{(j,\nu_j) \in T}\right) \\
 &= \sum_{(j,\nu_j) \in T} \text{size}\left(\text{Id}_{\mathbb{R}} \circ \tilde{L}_{\nu_j,\delta'}(y_j)\right) \\
 &\leq 2m(\Lambda)d(\Lambda) \max_{(j,\nu_j) \in T} \left(\text{size}(\text{Id}_{\mathbb{R}}) + \text{size}\left(\tilde{L}_{\nu_j,\delta'}(y_j)\right)\right) \\
 &\leq 10m(\Lambda)d(\Lambda) \max_{(j,\nu_j) \in T} \text{size}\left(\tilde{L}_{\nu_j,\delta'}(y_j)\right) \\
 &\leq Cd(\Lambda)m(\Lambda)^2 (m(\Lambda) + \log(\delta'^{-1})) \\
 &\leq Cd(\Lambda)m(\Lambda)^2 (\log(d(\Lambda)) + d(\Lambda) \log(m(\Lambda)) + m(\Lambda) + \log(\delta^{-1})). \tag{D.9}
 \end{aligned}$$

In (D.9) we used the NN calculus rules for the sizes of a sparse concatenation in (C.7) and a parallelization in (C.2). Furthermore, we used  $|T| \leq m(\Lambda)d(\Lambda)$  at the first equality and

$$\text{size}(\text{Id}_{\mathbb{R}}) \leq 4 \max_{(j, \nu_j) \in T} \text{depth} \left( \tilde{L}_{\nu_j, \delta'}(y_j) \right) \leq 4 \max_{(j, \nu_j) \in T} \text{size} \left( \tilde{L}_{\nu_j, \delta'}(y_j) \right), \quad (\text{D.10})$$

which follows from (C.4). At the third inequality in (D.9) we used the size bound for the univariate Legendre polynomials from (C.38).

For  $\text{size}(f_{\Lambda, \delta}^{(1)})$  it holds

$$\begin{aligned} \text{size} \left( f_{\Lambda, \delta}^{(1)} \right) &= \sum_{\nu \in \Lambda} \text{size} \left( \text{Id}_{\mathbb{R}} \circ \tilde{\Pi}_{\delta/2, M_{\nu}} \right) \\ &\leq 2 \sum_{\nu \in \Lambda} \left( \text{size}(\text{Id}_{\mathbb{R}}) + \text{size} \left( \tilde{\Pi}_{\delta/2, M_{\nu}} \right) \right) \\ &\leq 10|\Lambda| \max_{\nu \in \Lambda} \text{size} \left( \tilde{\Pi}_{\delta/2, M_{\nu}} \right) \\ &\leq C|\Lambda| \max_{\nu \in \Lambda} |\text{supp } \nu| \left( \log(|\text{supp } \nu|) + |\text{supp } \nu| \log(M_{\nu}) + \log(\delta^{-1}) \right) \\ &\leq C|\Lambda|d(\Lambda) \left( \log(d(\Lambda)) + d(\Lambda) \log(m(\Lambda)) + \log(\delta^{-1}) \right). \end{aligned} \quad (\text{D.11})$$

In (D.11), we used the size bound for  $\tilde{\Pi}$  from (C.25) and the argument from (D.10). Additionally we used  $M_{\nu} = 2|\nu|_1 + 2 \leq 4m(\Lambda)$ . Combining (D.9) and (D.11) shows the size bound for  $f_{\Lambda, \delta}$ .

The network  $f_{\Lambda, \delta}$  consists of sparse concatenations and parallelizations of the networks  $\tilde{\Pi}$  and  $\tilde{L}_j$ . Because we have  $\text{mpar}(\tilde{\Pi}) \leq 1$  and  $\text{mpar}(\tilde{L}_j) \leq 1$ , the NN calculus rules (C.11) and (C.1) yield  $\text{mpar}(f_{\Lambda, \delta}) \leq 1$ . This finishes the proof.

## D.2 RePU-Realization of Tensorized Legendre Polynomials

We show a result analogous to Proposition 17 for the RePU-realization of tensorized Legendre polynomials. The construction is similar to Opschoor et al. (2022a, Proposition 2.13).

**Proposition 51** ( $\sigma_q$ -NN approximation of  $L_{\nu}$ ) *Consider the setting of Proposition 17. Let  $q \in \mathbb{N}$ ,  $q \geq 2$ . Then there exists a  $\sigma_q$ -NN  $f_{\Lambda}$  such that the outputs  $\{\tilde{L}_{\nu}\}_{\nu \in \Lambda}$  of  $f_{\Lambda}$  satisfy*

$$\forall \nu \in \Lambda, \forall \mathbf{y} \in U : \quad \tilde{L}_{\nu}(\mathbf{y}) = L_{\nu}(\mathbf{y}).$$

Furthermore, there exists a constant  $C_q > 0$  depending only on  $q$  such that

$$\begin{aligned} \text{depth}(f_{\Lambda}) &\leq C_q \left( m(\Lambda) + \log(d(\Lambda)) \right), \\ \text{width}(f_{\Lambda}) &\leq C_q |\Lambda| d(\Lambda), \\ \text{size}(f_{\Lambda, \delta}) &\leq C_q d(\Lambda) \left( |\Lambda| + m(\Lambda)^2 \right), \\ \text{mpar}(f_{\Lambda, \delta}) &\leq C_q. \end{aligned}$$

**Proof** Similar to the proof of Proposition 17 we define  $f_\Lambda$  as a composition of two subnetworks  $f_\Lambda^{(1)}$  and  $f_\Lambda^{(2)}$ . It holds

$$f_\Lambda^{(2)}(\mathbf{y}) := \left( \left\{ \text{Id}_{\mathbb{R}} \circ \tilde{L}_{\nu_j}(y_j) \right\}_{(j, \nu_j) \in T} \right)$$

and

$$f_\Lambda^{(1)}((z_k)_{k \leq |T|}) = f_\Lambda^{(1)}\left(f_\Lambda^{(2)}(\mathbf{y})\right) := \left( \left\{ \text{Id}_{\mathbb{R}} \circ \widetilde{\prod} \left( \left\{ \text{Id}_{\mathbb{R}} \circ \tilde{L}_{\nu_j}(y_j) \right\}_{j \in \text{supp } \boldsymbol{\nu}} \right) \right\}_{\boldsymbol{\nu} \in \Lambda} \right)$$

with  $T$  from (D.2) and  $\mathbf{y} = (y_j)_{(j, \nu_j) \in T}$ . Furthermore, we use the  $\sigma_q$ -NNs  $\tilde{L}_j$  from Corollary 50 and  $\widetilde{\prod}$  from Proposition 46. The calculations are similar to the proof of Proposition 17. It holds

$$\text{depth}\left(f_\Lambda^{(2)}\right) = 1 + \max_{\substack{\boldsymbol{\nu} \in \Lambda \\ j \in \text{supp } \boldsymbol{\nu}}} \text{depth}\left(\tilde{L}_{\nu_j}\right) \leq C_q \max_{\substack{\boldsymbol{\nu} \in \Lambda \\ j \in \text{supp } \boldsymbol{\nu}}} \nu_j \leq C_q m(\Lambda). \quad (\text{D.12})$$

In (D.12) we used the depth bound for univariate Legendre polynomials, (C.39). Furthermore, we used  $\nu_j \leq m(\Lambda)$  for all  $\boldsymbol{\nu} \in \Lambda$  and  $j \in \text{supp } \boldsymbol{\nu}$ . For the depth of  $f_\Lambda^{(1)}$  it holds

$$\begin{aligned} \text{depth}\left(f_\Lambda^{(1)}\right) &= 1 + \max_{\boldsymbol{\nu} \in \Lambda} \text{depth}\left(\widetilde{\prod}\right) \leq 1 + C_q \max_{\boldsymbol{\nu} \in \Lambda} \log(|\text{supp } \boldsymbol{\nu}|) \\ &\leq 1 + C_q \log(d(\Lambda)), \end{aligned} \quad (\text{D.13})$$

where we used  $|\text{supp } \boldsymbol{\nu}| \leq d(\Lambda)$  for all  $\boldsymbol{\nu} \in \Lambda$  and the depth bound for  $\sigma_q$ -multiplication networks from Proposition 46. Combining the two depth bounds from (D.12) and (D.13), we get

$$\text{depth}(f_\Lambda) \leq C_q (m(\Lambda) + \log(d(\Lambda))).$$

For the width  $\text{width}(f_\Lambda)$  we calculate

$$\text{width}\left(f_\Lambda^{(2)}\right) = \sum_{(j, \nu_j) \in T} \text{width}\left(\text{Id}_{\mathbb{R}} \circ \tilde{L}_{\nu_j}\right) \leq C_q \sum_{(j, \nu_j) \in T} \text{width}\left(\tilde{L}_{\nu_j}\right) \leq C_q |T|, \quad (\text{D.14})$$

where we used (C.16) for the width of a  $\sigma_q$ -sparse concatenation and (C.40) for the width of  $\tilde{L}_{\nu_j}$ . For  $\text{width}(f_\Lambda^{(1)})$  it holds

$$\begin{aligned} \text{width}\left(f_\Lambda^{(1)}\right) &\leq \sum_{\boldsymbol{\nu} \in \Lambda} \text{width}\left(\text{Id}_{\mathbb{R}} \circ \widetilde{\prod}\right) \leq C_q \sum_{\boldsymbol{\nu} \in \Lambda} \text{width}\left(\widetilde{\prod}\right) \\ &\leq C_q \sum_{\boldsymbol{\nu} \in \Lambda} d(\Lambda) = C_q |\Lambda| d(\Lambda) \end{aligned} \quad (\text{D.15})$$

using (C.16) and (C.32) for the width of the multiplication network  $\widetilde{\prod}$ . Combining (D.14) and (D.15) gives

$$\text{width}(f_\Lambda) \leq C_q |\Lambda| d(\Lambda),$$

where  $|T| \leq |\Lambda|d(\Lambda)$  was used.

To estimate  $\text{size}(f_\Lambda)$ , we use (C.13) and find  $\text{size}(f_\Lambda) \leq C_q(\text{size}(f_\Lambda^{(1)}) + \text{size}(f_\Lambda^{(2)}))$ . We calculate

$$\begin{aligned}
 \text{size}\left(f_\Lambda^{(2)}\right) &= \text{size}\left(\left\{\text{Id}_{\mathbb{R}} \circ \tilde{L}_{\nu_j}(y_j)\right\}_{(j,\nu_j) \in T}\right) \\
 &= \sum_{(j,\nu_j) \in T} \text{size}\left(\text{Id}_{\mathbb{R}} \circ \tilde{L}_{\nu_j}(y_j)\right) \\
 &\leq C_q m(\Lambda) d(\Lambda) \max_{(j,\nu_j) \in T} \left(\text{size}(\text{Id}_{\mathbb{R}}) + \text{size}\left(\tilde{L}_{\nu_j}(y_j)\right)\right) \\
 &\leq C_q m(\Lambda) d(\Lambda) \max_{(j,\nu_j) \in T} \text{size}\left(\tilde{L}_{\nu_j}(y_j)\right) \\
 &\leq C_q d(\Lambda) m(\Lambda)^2.
 \end{aligned} \tag{D.16}$$

In (D.16) we used  $|T| \leq m(\Lambda)d$  at the first inequality and

$$\text{size}(\text{Id}_{\mathbb{R}}) \leq C_q \max_{(j,\nu_j) \in T} \text{depth}\left(\tilde{L}_{\nu_j}(y_j)\right) \leq C_q \max_{(j,\nu_j) \in T} \text{size}\left(\tilde{L}_{\nu_j}(y_j)\right), \tag{D.17}$$

which follows from (C.6). At the third inequality in (D.16) we used the size bound for the univariate Legendre polynomials from (C.41).

For  $\text{size}(f_\Lambda^{(1)})$  it holds

$$\begin{aligned}
 \text{size}\left(f_\Lambda^{(1)}\right) &= \sum_{\nu \in \Lambda} \text{size}\left(\text{Id}_{\mathbb{R}} \circ \tilde{\Pi}\right) \\
 &\leq C_q \sum_{\nu \in \Lambda} \left(\text{size}(\text{Id}_{\mathbb{R}}) + \text{size}\left(\tilde{\Pi}\right)\right) \\
 &\leq C_q |\Lambda| \max_{\nu \in \Lambda} \text{size}\left(\tilde{\Pi}\right) \\
 &\leq C_q |\Lambda| \max_{\nu \in \Lambda} |\text{supp } \nu| \\
 &\leq C_q |\Lambda| d(\Lambda).
 \end{aligned} \tag{D.18}$$

In (D.18) we used the size bound for  $\tilde{\Pi}$  from (C.33) and the argument from (D.17). Combining (D.16) and (D.18) shows the size bound for  $f_\Lambda$ .

The network  $f_\Lambda$  consists of sparse concatenations and parallelizations of the networks  $\tilde{\Pi}$  and  $\tilde{L}_j$ . Because we have  $\text{mpar}(\tilde{\Pi}) \leq C_q$  and  $\text{mpar}(\tilde{L}_j) \leq C_q$ , the NN calculus rules (C.17) and (C.1) yield  $\text{mpar}(f_\Lambda) \leq C_q$ . This finishes the proof.  $\blacksquare$

### D.3 Proof of Theorem 18

The following two theorems are similar to Herrmann et al. (2024, Theorem 5) and will be required for the proof of Theorem 18.

**Theorem 52** For  $N, q \in \mathbb{N}$ , consider the sparse FrameNet class  $\mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$ . Let Assumption 2 be satisfied with  $r > 1$  and  $t > 0$ . Fix  $\tau > 0$  (arbitrary small). Then there exists a constant  $C > 0$  independent of  $N$ , such that there exists  $\Gamma_N \in \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$  with

$$\sup_{a \in C_R^r(\mathcal{X})} \|\Gamma_N(a) - G_0(a)\|_{\mathcal{Y}} \leq CN^{-\min\{r-1, t\} + \tau}. \quad (\text{D.19})$$

**Theorem 53** Consider the setting of Theorem 52. Let  $\Psi_{\mathcal{X}}$  be a Riesz basis. Additionally, let  $\gamma$  be as in (3.9). Fix  $\tau > 0$  (arbitrary small). Then there exists a constant  $C > 0$  independent of  $N$ , such that there exists  $\Gamma_N \in \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$  with

$$\|\Gamma_N - G_0\|_{L^2(C_R^r(\mathcal{X}), \gamma; \mathcal{Y})} \leq CN^{-\min\{r-\frac{1}{2}, t\} + \tau}. \quad (\text{D.20})$$

We first show that Theorems 52 and 53 imply Theorem 18.

**Proof** [Proof of Theorem 18] First consider the setting of Theorem 52. Let  $\tau > 0$ . Then there exists a constant  $C$  independent of  $N$  and a FrameNet  $\Gamma_N \in \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$  such that for all  $N \in \mathbb{N}$

$$\|\Gamma_N - G_0\|_{\infty, \text{supp}(\gamma)}^2 \leq \sup_{a \in C_R^r(\mathcal{X})} \|\Gamma_N(a) - G_0(a)\|_{\mathcal{Y}}^2 \leq CN^{-2\min\{r-1, t\} + \tau},$$

where we used (D.19) with  $\tau/2$  and  $\text{supp}(\gamma) \subseteq C_R^r(\mathcal{X})$  by Assumption 2.

Now consider the setting of Theorem 53. Let  $\tau > 0$ . Then there exists a constant  $C$  independent of  $N$  and a FrameNet  $\mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$  with

$$\|\Gamma_N - G_0\|_{L^2(\gamma)}^2 \leq \|\Gamma_N - G_0\|_{L^2(C_R^r(\mathcal{X}), \gamma; \mathcal{Y})}^2 \leq CN^{-2\min\{r-\frac{1}{2}, t\} + \tau},$$

where we used  $\text{supp}(\gamma) \subset C_R^r(\mathcal{X})$  (Assumption 2) and (D.20) with  $\tau/2$ . ■

We are left to prove Theorems 52 and 53. We need some auxiliary results.

### D.3.1 AUXILIARY RESULTS

For  $r > 1$ ,  $R > 0$ ,  $U = [-1, 1]^{\mathbb{N}}$  and  $\sigma_R^r$  from (3.7) we define

$$u : U \rightarrow \mathcal{Y}, \quad u(\mathbf{y}) := (G_0 \circ \sigma_R^r)(\mathbf{y}).$$

For the proofs of Theorems 52 and 53 we do a  $\mathcal{Y}$ -valued tensorized Legendre expansion of  $u$  in the frame  $(\eta_j L_{\nu}(\mathbf{y}))_{j, \nu}$  of  $L^2(U, \pi; \mathcal{Y})$ , which reads

$$u(\mathbf{y}) = G_0(\sigma_R^r(\mathbf{y})) = \sum_{j \in \mathbb{N}} \sum_{\nu \in \mathcal{F}} c_{\nu, j} \eta_j L_{\nu}(\mathbf{y}) \quad (\text{D.21})$$

with Legendre coefficients

$$c_{\nu, j} := \int_U L_{\nu}(\mathbf{y}) \langle u(\mathbf{y}), \tilde{\eta}_j \rangle_{\mathcal{Y}} d\pi(\mathbf{y}). \quad (\text{D.22})$$

Our aim is to construct the network  $\Gamma_N$  out of the tensorized Legendre polynomials with the ‘‘most important’’ contributions to the expansion. This contribution is quantified via

the Legendre coefficients  $c_{\nu,j}$  in (D.22). We therefore have to examine bounds on  $c_{\nu,j}$  and analyze their respective structure. Therefore consider the following order relation on multi-indices in  $\mathcal{F}$  from (3.11). For  $\mu, \nu \in \mathcal{F}$  we write  $\mu \leq \nu$  if and only if  $\mu_j \leq \nu_j$  for all  $j \in \mathbb{N}$ . We call a set  $\Lambda \subset \mathcal{F}$  *downward closed* if and only if  $\nu \in \mathcal{F}$  implies  $\mu \in \mathcal{F}$  for all  $\mu \leq \nu$ . Furthermore, for  $\nu \in \mathcal{F}$ , define

$$\omega_\nu := \prod_{j=1}^{\infty} (1 + 2\nu_j).$$

The following theorem is a special case of Zech (2018, Theorem 2.2.10). The formulation is similar to Herrmann et al. (2024, Theorem 4).

**Theorem 54 (Herrmann et al. 2024, Theorem 4)** *Let Assumption 2 be satisfied with  $r > 1$  and  $t > 0$ . Fix  $\tau > 0$ ,  $p \in (\frac{1}{r}, 1]$  and  $t' \in [0, t]$ . Consider  $\mathcal{F}$  from (3.11), and let  $\pi = \otimes_{j \in \mathbb{N}} \frac{\lambda}{2}$  be the infinite product (probability) measure on  $U = [-1, 1]^{\mathbb{N}}$ , where  $\lambda$  denotes the Lebesgue measure on  $[-1, 1]$ . Then there exists  $C > 0$  and a sequence  $(a_\nu)_{\nu \in \mathcal{F}} \in \ell^p(\mathcal{F})$  of positive numbers such that*

(i) for each  $\nu \in \mathcal{F}$

$$\omega_\nu^\tau \left\| \int_U L_\nu(\mathbf{y}) u(\mathbf{y}) \, d\pi(\mathbf{y}) \right\|_{\mathcal{Y}^{t'}} \leq C a_\nu,$$

(ii) there exists an enumeration  $(\nu_i)_{i \in \mathbb{N}}$  of  $\mathcal{F}$  such that  $(a_{\nu_i})_{i \in \mathbb{N}}$  is monotonically decreasing, the set  $\Lambda_N := \{\nu_i : i \leq N\} \subseteq \mathcal{F}$  is downward closed for each  $N \in \mathbb{N}$ , and additionally

$$m(\Lambda_N) := \max_{i=1, \dots, N} |\nu_i| = \mathcal{O}(\log(|\Lambda_N|)), \quad (\text{D.23})$$

$$d(\Lambda_N) := \max_{i=1, \dots, N} |\text{supp } \nu_i| = o(\log(|\Lambda_N|)) \quad (\text{D.24})$$

for  $N \rightarrow \infty$ ,

(iii) the following expansion holds with absolute and uniform convergence:

$$\forall \mathbf{y} \in U : \quad u(\mathbf{y}) = \sum_{\nu \in \mathcal{F}} L_\nu(\mathbf{y}) \int_U L_\nu(\mathbf{x}) u(\mathbf{x}) \, d\pi(\mathbf{x}) \in \mathcal{Y}^{t'}.$$

The following proposition reformulates Theorem 54 (i) into a bound for  $c_{\nu,j}$ . It was shown in Herrmann et al. (2024, Proposition 2). Recall that  $\theta_j$  denote the weights to define the spaces  $\mathcal{Y}^{t'}$ ,  $t' > 0$ , see Definition 14.

**Proposition 55 (Herrmann et al. 2024, Proposition 2)** *Consider the setting of Theorem 54. Then for each  $\nu \in \mathcal{F}$*

$$\omega_\nu^{2\tau} \sum_{j \in \mathbb{N}} \theta_j^{-2t'} c_{\nu,j}^2 \leq C^2 a_\nu^2.$$

Proposition 55 gives decay of the coefficients  $c_{\nu,j}$  in both  $j$  and  $\nu$ . Since  $\theta_j = \mathcal{O}(j^{-1+\tau})$  for all  $\tau > 0$  we have  $c_{\nu,j}^2 = \mathcal{O}(j^{-1-2t'+\tilde{\tau}})$  for  $\tilde{\tau} < 2\tau t'$  and every  $\nu \in \Lambda_N$ . Furthermore, since  $(a_\nu)_{\nu \in \mathcal{F}} \in \ell^p(\mathcal{F})$  the Legendre coefficients  $c_{\nu,j}$  decay algebraically in  $\nu$ . We continue with a technical lemma, which was shown in Herrmann et al. (2024, Lemma 4).

**Lemma 56 (Herrmann et al. 2024, Lemma 4)** *Let  $\alpha > 1$ ,  $\beta > 0$  and assume two sequences  $(a_i)_{i \in \mathbb{N}}$  and  $(d_j)_{j \in \mathbb{N}}$  in  $\mathbb{R}$  with  $a_i \lesssim i^{-\alpha}$  and  $d_j \lesssim j^{-\beta}$  for all  $i, j \in \mathbb{N}$ . Additionally assume that  $(d_j)_{j \in \mathbb{N}}$  is monotonically decreasing. Suppose that there exists a constant  $C < \infty$  such that the sequence  $(c_{i,j})_{i,j \in \mathbb{N}}$  satisfies*

$$\forall i \in \mathbb{N} : \sum_{j \in \mathbb{N}} c_{i,j}^2 d_j^{-2} \leq C^2 a_i^2.$$

Then for every  $\tau > 0$

(i) for all  $N \in \mathbb{N}$  there exists  $(m_i)_{i \in \mathbb{N}} \subseteq \mathbb{N}_0^{\mathbb{N}}$  monotonically decreasing s.t.  $\sum_{i \in \mathbb{N}} m_i \leq N$  and

$$\sum_{i \in \mathbb{N}} \left( \sum_{j > m_i} c_{i,j}^2 \right)^{\frac{1}{2}} \lesssim N^{-\min\{\alpha-1, \beta\} + \tau},$$

(ii) for all  $N \in \mathbb{N}$  there exists  $(m_i)_{i \in \mathbb{N}} \subseteq \mathbb{N}_0^{\mathbb{N}}$  monotonically decreasing s.t.  $\sum_{i \in \mathbb{N}} m_i \leq N$  and

$$\left( \sum_{i \in \mathbb{N}} \sum_{j > m_i} c_{i,j}^2 \right)^{\frac{1}{2}} \lesssim N^{-\min\{\alpha-\frac{1}{2}, \beta\} + \tau}.$$

In the following, we use Lemma 56 to get a decay property for the Legendre coefficients  $c_{\mathbf{v}_i, j}$  with the enumeration  $\mathbf{v}_i$  of  $\Lambda_N$  from Theorem 54. The sequence  $\mathbf{m} = (m_i)_{i \in \mathbb{N}}$  quantifies which coefficients of the Legendre expansion are ‘‘important’’ and are therefore used to define the surrogate  $\Gamma_N$ .

We first show that Theorem 54 yields sufficient decay on the Legendre coefficients  $c_{\mathbf{v}_i, j}$  s.t. the assumptions of Lemma 56 are satisfied.

**Lemma 57** *Consider the setting of Theorem 54. Let  $\tilde{\tau} > 0$  such that  $1/p > r - \tilde{\tau}/2$ . Then the assumptions of Lemma 56 are fulfilled for  $\alpha = r - \tilde{\tau}/2$ ,  $\beta = t - \tilde{\tau}/2$ ,  $a_i = a_{\mathbf{v}_i}$ ,  $d_j = \theta_j^{t'}$  and  $c_{i,j} = \omega_{\mathbf{v}_i}^{1/2} c_{\mathbf{v}_i, j}$  for  $i, j \in \mathbb{N}$ .*

**Proof** Proposition 55 with  $\tau = \frac{1}{2}$  gives

$$\omega_{\mathbf{v}_i}^{\frac{1}{2}} \left( \sum_{j \in \mathbb{N}} \theta_j^{-2t'} c_{\mathbf{v}_i, j}^2 \right)^{\frac{1}{2}} = \mathcal{O}(a_{\mathbf{v}_i}) = \mathcal{O}\left(i^{-r + \frac{\tilde{\tau}}{2}}\right). \quad (\text{D.25})$$

The last equality in (D.25) holds because  $ia_{\mathbf{v}_i}^p \leq \sum_{j \in \mathbb{N}} a_{\mathbf{v}_j}^p < \infty$  (since  $a_{\mathbf{v}_i}$  is monotonically decreasing) implies  $a_{\mathbf{v}_i} = \mathcal{O}(i^{-1/p}) = \mathcal{O}(i^{-r + \tilde{\tau}/2})$ . Since  $(\theta_j^{t'})_{j \in \mathbb{N}} \in \ell^{1/(t - \tilde{\tau}/2)}$  (see Definition 14) it holds

$$\theta_j^{t'} = \mathcal{O}(j^{-t + \tilde{\tau}/2})$$

with the same argument. ■

## D.3.2 PROOFS OF THEOREMS 52 AND 53

The proof of Theorems 52 and 53 is similar to Herrmann et al. (2024, Sections 4.2-4.4, Proofs of Theorems 1,2 and 5).

**Proof** [Proof of Theorem 52] Let  $(a_{\nu})_{\nu \in \mathcal{F}}$  be the enumeration  $(\nu_{\mathbf{i}})_{\mathbf{i} \in \mathbb{N}}$  from Theorem 54, where we use the case  $\tau = \frac{1}{2}$ . Therefore  $(a_{\nu_{\mathbf{i}}})_{\mathbf{i} \in \mathbb{N}}$  is monotonically decreasing and belongs to  $\ell^p$  with  $p \in (\frac{1}{r}, 1]$ . We further fix  $\tilde{\tau} > 0$  and demand  $\frac{1}{p} > r - \frac{\tilde{\tau}}{2}$ . Fix  $\tilde{N} \in \mathbb{N}$  and set  $\Lambda_{\tilde{N}} := \{\nu_j : j \leq \tilde{N}\} \subset \mathcal{F}$ , which is downward closed by Theorem 54. Now we approximate the tensorized Legendre polynomials  $L_{\nu}$  on the index set  $\Lambda_{\tilde{N}}$ . Let  $\rho \in (0, \frac{1}{2})$ . In the ReLU case, Proposition 17 gives a NN  $f_{\Lambda_{\tilde{N}}, \rho}$  with outputs  $\{\tilde{L}_{\nu, \rho}\}_{\nu \in \Lambda_{\tilde{N}}}$  s.t.

$$\sup_{\mathbf{y} \in U} \max_{\nu \in \Lambda_{\tilde{N}}} |L_{\nu}(\mathbf{y}) - \tilde{L}_{\nu, \rho}(\mathbf{y})| \leq \rho.$$

Using  $|\Lambda_{\tilde{N}}| = \tilde{N}$ , (D.23) and (D.24), it holds for  $\tilde{N} \geq 2$

$$\begin{aligned} \text{depth}(f_{\Lambda_{\tilde{N}}, \rho}) &= \mathcal{O}\left(\log(\tilde{N})^2 \log(\log(\tilde{N}))^2 + \log(\tilde{N}) \log(\rho^{-1})\right), \\ \text{width}(f_{\Lambda_{\tilde{N}}, \rho}) &= \mathcal{O}(\tilde{N} \log(\tilde{N})), \\ \text{size}(f_{\Lambda_{\tilde{N}}, \rho}) &= \mathcal{O}\left(\tilde{N} \log(\tilde{N})^2 \log(\log(\tilde{N})) + \tilde{N} \log(\tilde{N}) \log(\rho^{-1})\right), \\ \text{mpar}(f_{\Lambda_{\tilde{N}}, \rho}) &= 1. \end{aligned}$$

The constants hidden in  $\mathcal{O}(\cdot)$  are independent of  $\tilde{N}$  and  $\rho$ . For  $\tilde{N} \in \mathbb{N}$ , set the accuracy  $\rho := \tilde{N}^{-\min\{r-\frac{1}{2}, t\}}$ . Then it holds

$$\begin{aligned} \text{depth}(f_{\Lambda_{\tilde{N}}, \rho}) &= \mathcal{O}\left(\log(\tilde{N})^2 \log(\log(\tilde{N}))^2\right), \\ \text{size}(f_{\Lambda_{\tilde{N}}, \rho}) &= \mathcal{O}\left(\tilde{N} \log(\tilde{N})^2 \log(\log(\tilde{N}))\right). \end{aligned}$$

Proposition 51 shows that the ReLU bounds also hold for the RePU-case.

By Lemma 57 the assumptions of Lemma 56 are satisfied. Applying Lemma 56 (i) with  $\alpha := r - \tilde{\tau}/2$  and  $\beta := t - \tilde{\tau}/2$  gives a sequence  $(m_i)_{i \in \mathbb{N}} \subset \mathbb{N}_0^{\mathbb{N}}$  such that  $\sum_{i \in \mathbb{N}} m_i \leq \tilde{N}$  and

$$\sum_{i \in \mathbb{N}} \omega_{\nu_{\mathbf{i}}}^{\frac{1}{2}} \left( \sum_{j > m_i} c_{\nu_{\mathbf{i}}, j}^2 \right)^{\frac{1}{2}} \leq C \tilde{N}^{-\min\{r-1, t\} + \tilde{\tau}}. \quad (\text{D.26})$$

We now define

$$\left( \tilde{\gamma}_{\tilde{N}, j} \right) := \sum_{\{i \in \mathbb{N} : m_i \geq j\}} \tilde{L}_{\nu_{\mathbf{i}}, \rho}(\mathbf{y}) c_{\nu_{\mathbf{i}}, j} \quad (\text{D.27})$$

for  $j \in \mathbb{N}$ , where empty sums are set to zero. Recall the uniform distribution  $\pi$  on  $U = [-1, 1]^{\mathbb{N}}$ , see Example 3. With  $\tilde{\gamma}_{\tilde{N}} = (\tilde{\gamma}_{\tilde{N},j})_{j \in \mathbb{N}}$  it holds

$$\begin{aligned}
 \|G_0 \circ \sigma_R^r(\mathbf{y}) - \mathcal{D}_{\mathbf{y}} \circ \tilde{\gamma}_{\tilde{N}}(\mathbf{y})\|_{\mathbf{y}} &= \left\| \sum_{i,j \in \mathbb{N}} c_{\nu_i,j} L_{\nu_i}(\mathbf{y}) \eta_j - \sum_{i \in \mathbb{N}} \sum_{j \leq m_i} c_{\nu_i,j} \tilde{L}_{\nu_i,\rho}(\mathbf{y}) \eta_j \right\|_{\mathbf{y}} \\
 &\leq \left\| \sum_{i \in \mathbb{N}} L_{\nu_i}(\mathbf{y}) \sum_{j > m_i} c_{\nu_i,j} \eta_j \right\|_{\mathbf{y}} + \left\| \sum_{i \in \mathbb{N}} (L_{\nu_i}(\mathbf{y}) - \tilde{L}_{\nu_i,\rho}(\mathbf{y})) \sum_{j \leq m_i} c_{\nu_i,j} \eta_j \right\|_{\mathbf{y}} \\
 &\leq \Lambda_{\Psi_{\mathbf{y}}} \sum_{i \in \mathbb{N}} \underbrace{\|L_{\nu_i}\|_{\infty,\pi}}_{\leq \omega_{\nu_i}^{\frac{1}{2}}} \left( \sum_{j > m_i} c_{\nu_i,j}^2 \right)^{\frac{1}{2}} + \Lambda_{\Psi_{\mathbf{y}}} \rho \sum_{i \in \mathbb{N}} \left( \sum_{j \leq m_i} c_{\nu_i,j}^2 \right)^{\frac{1}{2}} \\
 &\leq \tilde{C} \Lambda_{\Psi_{\mathbf{y}}} \tilde{N}^{-\min\{r-1,t\}+\tilde{\tau}} + \tilde{C} \Lambda_{\Psi_{\mathbf{y}}} \rho \leq \tilde{C} \tilde{N}^{-\min\{r-1,t\}+\tilde{\tau}}
 \end{aligned} \tag{D.28}$$

for all  $\mathbf{y} \in U$ . In (D.28) we used the definition of  $\mathcal{D}_{\mathbf{y}}$ , (D.27) and (D.21) at the first equality. Furthermore, we used (D.26), the definition of  $\rho$  and

$$\begin{aligned}
 \sum_{i \in \mathbb{N}} \left( \sum_{j \leq m_i} c_{\nu_i,j}^2 \right)^{\frac{1}{2}} &\leq \sum_{i \in \mathbb{N}} \left( \sum_{j \in \mathbb{N}} c_{\nu_i,j}^2 \right)^{\frac{1}{2}} \\
 &\leq \tilde{C} \sum_{i \in \mathbb{N}} \omega_{\nu_i}^{\frac{1}{2}} \left( \sum_{j \in \mathbb{N}} \theta_j^{-2t'} c_{\nu_i,j}^2 \right)^{\frac{1}{2}} \\
 &\leq \tilde{C} \sum_{i \in \mathbb{N}} a_{\nu_i} \leq \tilde{C} \sum_{i \in \mathbb{N}} i^{-r+\tilde{\tau}/2} \leq \tilde{C}
 \end{aligned} \tag{D.29}$$

at the second-to-last inequality. We changed the constants  $\tilde{C}$  from line to line in (D.28) and (D.29). The last line of (D.28) shows why the RePU-case does not improve the approximation property qualitatively. In the RePU-case, Proposition 51 gives a  $\sigma_q$ -NN  $f_{\Lambda}$  exactly realizing the tensorized Legendre polynomials, i.e. the case  $\rho = 0$  from above. Therefore the second summand in the last line of (D.28) vanishes. This does not improve the approximation rate due to the first summand. This part depends on the summability properties of the Legendre coefficients  $c_{\nu_i,j}$  following Assumption 2 and is therefore independent of the activation function  $\sigma$ .

Now we argue similarly to Herrmann et al. (2024, Proof of Theorem 1). Consider the scaling  $S_r$  from (3.5). It holds

$$S_r \circ \mathcal{E}_{\mathcal{X}}(a) \in U \quad \forall a \in C_R^r(\mathcal{X}), \tag{D.30}$$

because of (3.8) and (3.5). We define  $\tilde{\Gamma}_{\tilde{N}} := \mathcal{D}_y \circ \tilde{\gamma}_{\tilde{N}} \circ S_r \circ \mathcal{E}_X$  and calculate

$$\begin{aligned}
 & \sup_{a \in C_R^r(X)} \left\| G_0(a) - \tilde{\Gamma}_{\tilde{N}}(a) \right\|_y = \sup_{a \in C_R^r(X)} \left\| G_0(a) - \mathcal{D}_y \circ \tilde{\gamma}_{\tilde{N}} \circ S_r \circ \mathcal{E}_X(a) \right\|_y \\
 &= \sup_{\{\mathbf{y} \in U: \sigma_R^r(\mathbf{y}) \in C_R^r(X)\}} \left\| G_0 \circ \sigma_R^r(\mathbf{y}) - \mathcal{D}_y \circ \tilde{\gamma}_{\tilde{N}} \circ S_r \circ \mathcal{E}_X \circ \sigma_R^r(\mathbf{y}) \right\|_y \\
 &\leq \sup_{\mathbf{y} \in U} \left\| G_0 \circ \sigma_R^r(\mathbf{y}) - \mathcal{D}_y \circ \tilde{\gamma}_{\tilde{N}}(\mathbf{y}) \right\|_y \leq \tilde{C} \tilde{N}^{-\min\{r-1, t\} + \tilde{\tau}}, \tag{D.31}
 \end{aligned}$$

where we used (D.30) and (D.28).

In order to finish the proof of Theorem 52, we relate  $\tilde{N}$  to  $N$  and show  $\Gamma_N := \tilde{\Gamma}_{\tilde{N}} \in \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$ , i.e. we show that the approximation networks we constructed have the desired sparse structure. We simultaneously prove the ReLU- and RePU-case.

In order to analyze the NNs  $\tilde{\gamma}_{\tilde{N}}$  from (D.27), we specify its structure. We set  $n_j = |\{m_i \geq j\}|$  and define

$$\tilde{\gamma}_{\tilde{N}} = \left( \left\{ \Sigma_{n_j} \left( \left\{ \text{Id}_{\mathbb{R}} \circ SM_{c_{\nu_i, j}} \circ \tilde{L}_{\nu_i, \rho} \right\}_{m_i \geq j} \right) \right\}_{j \in \mathbb{N}} \right). \tag{D.32}$$

The round brackets in (D.32) denote a parallelization. The networks  $SM_{c_{\nu_i, j}}$  denote the scalar multiplication networks from Definition 42. Furthermore,  $\Sigma_{n_j}$  denotes the summation network from Definition 41 and we use the identity networks  $\text{Id}_{\mathbb{R}}$  from Lemma 37 or 38 to synchronize the depth. Using the respective bounds for the summation and scalar multiplication networks and the NN calculus for parallelization and sparse concatenation we get

$$\begin{aligned}
 \text{depth}(\tilde{\gamma}_{\tilde{N}}) &\leq 2 + \max_{i, j \in \mathbb{N}^2, m_i \geq j} \left( \text{depth}(\tilde{L}_{\nu_i, \rho}) + \text{depth}(SM_{c_{\nu_i, j}}) \right) + \max_{j \in \mathbb{N}} \text{depth}(\Sigma_{n_j}) \\
 &\leq 3 + \mathcal{O}(\log(\tilde{N})^2 \log(\log(\tilde{N}))) + \max_{i, j \in \mathbb{N}^2, m_i \geq j} C_q \log(|c_{\nu_i, j}|) + 0 \\
 &= \mathcal{O}(\log(\tilde{N})^2 \log(\log(\tilde{N}))), \tag{D.33}
 \end{aligned}$$

$$\begin{aligned}
 \text{width}(\tilde{\gamma}_{\tilde{N}}) &\leq \max \left\{ \sum_{j \in \mathbb{N}} \sum_{m_i \geq j} \text{width}(\tilde{L}_{\nu_i, \rho}), \sum_{j \in \mathbb{N}} \sum_{m_i \geq j} \text{width}(SM_{c_{\nu_i, j}}), \right. \\
 &\quad \left. \sum_{j \in \mathbb{N}} \sum_{m_i \geq j} \text{width}(\text{Id}_{\mathbb{R}}), \sum_{j \in \mathbb{N}} \text{width}(\Sigma_{n_j}) \right\} \\
 &\leq C_q \max \left\{ \mathcal{O}(\tilde{N} \log(\tilde{N})), \sum_{j \in \mathbb{N}} \sum_{m_i \geq j} C_q, \sum_{j \in \mathbb{N}} n_j \right\} = \mathcal{O}(\tilde{N} \log(\tilde{N})),
 \end{aligned}$$

as well as

$$\begin{aligned}
 \text{size}(\tilde{\gamma}_{\tilde{N}}) &\leq C_q \sum_{j \in \mathbb{N}} \sum_{m_i \geq j} \left( \text{size}(\tilde{L}_{\mathbf{v}_i, \rho}) + \text{size}(SM_{c_{\mathbf{v}_i, j}}) + \text{size}(\text{Id}_{\mathbb{R}}) \right) + C_q \sum_{j \in \mathbb{N}} \text{size}(\Sigma_{n_j}) \\
 &= \mathcal{O}\left(\tilde{N} \log(\tilde{N})^2 \log(\log(\tilde{N}))\right) + \sum_{j \in \mathbb{N}} \left[ \left( \sum_{m_i \geq j} C_q \log(|c_{\mathbf{v}_i, j}|) \right) + C_q n_j \right] \\
 &\leq \mathcal{O}\left(\tilde{N} \log(\tilde{N})^2 \log(\log(\tilde{N}))\right) + C_q \sum_{j \in \mathbb{N}} \sum_{m_i \geq j} 1 \\
 &= \mathcal{O}\left(\tilde{N} \log(\tilde{N})^2 \log(\log(\tilde{N}))\right), \tag{D.34}
 \end{aligned}$$

$$\text{mpar}(\tilde{\gamma}_{\tilde{N}}) \leq C_q.$$

In (D.33)–(D.34) we used  $\log(|c_{\mathbf{v}_i, j}|) \leq C_q$  for all  $i, j \in \mathbb{N}$  independent of  $n$ . Furthermore, we used

$$\sum_{j \in \mathbb{N}} n_j = \sum_{j \in \mathbb{N}} \sum_{m_i \geq j} 1 = \sum_{i \in \mathbb{N}} \sum_{j \leq m_i} 1 = \sum_{i \in \mathbb{N}} m_i \leq \tilde{N}.$$

To get rid of the logarithmic terms, we define  $N = N(\tilde{N}) := \max\{1, \tilde{N} \log(\tilde{N})^3\}$  and obtain a NN  $\gamma_N = \tilde{\gamma}_{\tilde{N}}$  with

$$\begin{aligned}
 \text{depth}(\gamma_N) &= \mathcal{O}(\log(N)), \\
 \text{width}(\gamma_N) &= \mathcal{O}(N), \\
 \text{size}(\gamma_N) &= N, \\
 \text{mpar}(\gamma_N) &\leq C_q
 \end{aligned}$$

and error less than

$$\tilde{C} \tilde{N}^{-\min\{r-1, t\} + \tilde{\tau}} = \tilde{C} \tilde{N}^{-\kappa} \leq \tilde{C} (3\kappa/\tilde{\tau})^{3\kappa} N^{-\kappa + \tilde{\tau}} := C N^{-\min\{r-1, t\} + \tau}. \tag{D.35}$$

Per definition  $\text{depth}(\gamma_N) = \mathcal{O}(\log(N))$  and  $\text{width}(\gamma_N) = \mathcal{O}(N)$  yields constants  $\tilde{C}_L, \tilde{C}_p$  and  $N_1, N_2 \in \mathbb{N}$  s.t.

$$\begin{aligned}
 \text{depth}(\gamma_N) &\leq \tilde{C}_L \max\{1, \log(N)\}, & N \geq N_1, \\
 \text{width}(\gamma_N) &\leq \tilde{C}_p N, & N \geq N_2.
 \end{aligned}$$

Setting  $C_L = \max\{\tilde{C}_L, \max_{N=2, \dots, N_1-1} \text{depth}(\gamma_N)/\log(2)\}$  and  $C_p = \max\{\tilde{C}_p, \max_{N=1, \dots, N_2-1} \text{width}(\gamma_N)\}$  shows

$$\begin{aligned}
 \text{depth}(\gamma_N) &\leq C_L \log(N), & N \in \mathbb{N}, \\
 \text{width}(\gamma_N) &\leq C_p N, & N \in \mathbb{N}.
 \end{aligned}$$

In order to show  $\Gamma_N := \mathcal{D}_y \circ \gamma_N \circ S_r \circ \mathcal{E}_x \in \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$ , we are left to show that the maximum Euclidean norm  $\|\cdot\|_2$  of  $\gamma_N$  in  $U$  is independent of  $N$ . It holds for all  $\mathbf{y} \in U$  that

$\mathcal{D}_X \circ S_r^{-1}(\mathbf{y}) \in C_R^r(\mathcal{X})$ . We get

$$\begin{aligned}
 \sup_{\mathbf{y} \in U} \|\gamma_N(\mathbf{y})\|_2 &= \sup_{\mathbf{y} \in U} \|\mathcal{E}_y \circ \Gamma_N \circ \mathcal{D}_X \circ S_r^{-1}(\mathbf{y})\|_2 \\
 &\leq \Lambda_{\Psi_y} \sup_{a \in C_R^r(\mathcal{X})} \|\Gamma_N(a)\|_y \\
 &= \Lambda_{\Psi_y} \sup_{a \in C_R^r(\mathcal{X})} \|\Gamma_N(a) - G_0(a) + G_0(a)\|_y \\
 &\leq \Lambda_{\Psi_y} \sup_{a \in C_R^r(\mathcal{X})} (\|\Gamma_N(a) - G_0(a)\|_y + c\|G_0(a)\|_{y_t}) \\
 &\leq \Lambda_{\Psi_y} (C + cC_{G_0}) =: B,
 \end{aligned} \tag{D.36}$$

where  $\Lambda_{\Psi_y}$  denotes the upper frame bound of  $\Psi_y$  and  $c = \theta_0^t$ , see Definition 14. In (D.36) we used Assumption 2 and the approximation error from (D.35). Thus  $\Gamma_N \in \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$  for all  $N \in \mathbb{N}$ , where we set  $C_s = 1$ . Using  $\text{supp}(\gamma) \subset C_R^r(\mathcal{X})$ , see Assumption 2, in (D.31) finalizes the proof of Theorem 52.  $\blacksquare$

**Proof** [Proof of Theorem 53] By Lemma 57 the assumptions of Lemma 56 are satisfied. Applying Lemma 56 (ii) with  $\alpha := r - \tilde{\tau}/2$  and  $\beta := t - \tilde{\tau}/2$  gives a sequence  $(m_i)_{i \in \mathbb{N}} \subset \mathbb{N}_0^{\mathbb{N}}$  such that  $\sum_{i \in \mathbb{N}} m_i \leq \tilde{N}$  and

$$\left( \sum_{i \in \mathbb{N}} \omega_{\nu_i} \sum_{j > m_i} c_{\nu_i, j}^2 \right)^{\frac{1}{2}} \leq \tilde{C} \tilde{N}^{-\min\{r-\frac{1}{2}, t\} + \tilde{\tau}}. \tag{D.37}$$

Define  $\tilde{\gamma}_{\tilde{N}} = (\tilde{\gamma}_{\tilde{N}, j})_{j \in \mathbb{N}}$  for all  $y \in U$  with  $\tilde{\gamma}_{\tilde{N}, j}$  as in (D.27). Then it holds

$$\begin{aligned}
 \|G_0 \circ \sigma_R^r - \mathcal{D}_y \circ \tilde{\gamma}_{\tilde{N}}\|_{L^2(U, \pi; \mathcal{Y})} &\leq \left\| \sum_{i \in \mathbb{N}} \sum_{j > m_i} c_{\nu_i, j} L_{\nu_i} \eta_j \right\|_{L^2(U, \pi; \mathcal{Y})} \\
 &\quad + \left\| \sum_{i \in \mathbb{N}} \sum_{j \leq m_i} c_{\nu_i, j} \eta_j (L_{\nu_i} - \tilde{L}_{\nu_i, \rho}) \right\|_{L^2(U, \pi; \mathcal{Y})} \\
 &\leq \Lambda_{\Psi_y} \left( \sum_{i \in \mathbb{N}} \underbrace{\|L_{\nu_i}\|_{\infty, \pi}^2}_{\leq \omega_{\nu_i}} \sum_{j > m_i} c_{\nu_i, j}^2 \right)^{\frac{1}{2}} + \Lambda_{\Psi_y} \rho \left( \sum_{i \in \mathbb{N}} \sum_{j \leq m_i} c_{\nu_i, j}^2 \right)^{\frac{1}{2}} \\
 &\leq \tilde{C} \Lambda_{\Psi_y} \tilde{N}^{-\min\{r-\frac{1}{2}, t\} + \tilde{\tau}} + \tilde{C} \Lambda_{\Psi_y} \rho \leq \tilde{C} \tilde{N}^{-\min\{r-\frac{1}{2}, t\} + \tilde{\tau}}.
 \end{aligned} \tag{D.38}$$

In (D.38) we used the definition of  $\mathcal{D}_y$ , (D.27) and (D.21) at the first inequality. Additionally we used that  $(L_{\nu} \eta_j)_{\nu, j}$  is a frame of  $L^2(U, \pi; \mathcal{Y})$  at the second inequality. Finally we used (D.37), the definition of  $\rho$  and an argument similar to (D.29) at the second-to-last inequality. Note that again we changed the constants  $\tilde{C}$  from line to line in (D.38).

Since  $\Psi_{\mathcal{X}}$  is a Riesz basis, we have (see Section 3.1.2 and 3.5)

$$C_R^r(\mathcal{X}) = \{\sigma_R^r(\mathbf{y}), \mathbf{y} \in U\} \quad \text{and} \quad \mathcal{E}_X \circ \sigma_R^r(\mathbf{y}) = S_r^{-1}(\mathbf{y}). \tag{D.39}$$

With  $\tilde{\Gamma}_{\tilde{N}} := \mathcal{D}_y \circ \tilde{\gamma}_{\tilde{N}} \circ S_r \circ \mathcal{E}_x$  we calculate

$$\begin{aligned} \left\| \tilde{\Gamma}_{\tilde{N}} - G_0 \right\|_{L^2(C_R^r(x), (\sigma_R^r)_{\#} \pi; y)} &= \left\| \mathcal{D}_y \circ \tilde{\gamma}_{\tilde{N}} \circ S_r \circ \mathcal{E}_x - G_0 \right\|_{L^2(C_R^r(x), (\sigma_R^r)_{\#} \pi; y)} \\ &= \left\| \mathcal{D}_y \circ \tilde{\gamma}_{\tilde{N}} \circ S_r \circ \mathcal{E}_x \circ \sigma_R^r - G_0 \circ \sigma_R^r \right\|_{L^2(U, \pi; y)} \\ &= \left\| \mathcal{D}_y \circ \tilde{\gamma}_{\tilde{N}} - G_0 \circ \sigma_R^r \right\|_{L^2(U, \pi; y)} \leq \tilde{C} \tilde{N}^{-\min\{r-\frac{1}{2}, t\} + \tau}, \end{aligned}$$

where we used (D.39) and (D.38). Defining  $N = N(\tilde{N}) := \max\{1, \tilde{N} \log(\tilde{N})^3\}$  we can proceed similar to the proof of Theorem 52 from (D.31) on. The reason this works is that the NNs  $\tilde{\gamma}_{\tilde{N}}$  are defined in the same way in the  $L^2$ - and the  $L^\infty$ -case (only the sequence  $\mathbf{m}$  changes, but not its properties). This shows  $\Gamma_N := \tilde{\Gamma}_{\tilde{N}} \in \mathbf{G}_{\text{FN}}^{\text{sp}}(\sigma_q, N)$  for all  $N \in \mathbb{N}$  and thus finishes the proof of Theorem 53.  $\blacksquare$

#### D.4 Proof of Lemma 20

The arguments in the following proof are based on entropy bounds for feedforward neural network classes, first established in Schmidt-Hieber (2020a, Proof of Lemma 5).

Define the supremum norm  $\|\cdot\|_{\infty, \infty}$  on  $\mathbf{g}_{\text{FN}}$  as

$$\|g\|_{\infty, \infty} := \sup_{\mathbf{y} \in \mathbb{R}^{p_0}} \|g(\mathbf{y})\|_{\infty}, \quad g \in \mathbf{g}_{\text{FN}}, \quad (\text{D.40})$$

where  $\|\cdot\|_{\infty}$  denotes the maximum norm in  $\mathbb{R}^n$ . Then Petersen et al. (2021, Proposition 3.5) shows that  $(\mathbf{g}_{\text{FN}}, \|\cdot\|_{\infty, \infty})$  is compact. Since the map  $i : \mathbf{g}_{\text{FN}} \rightarrow \mathbf{G}_{\text{FN}}, g \rightarrow G = \mathcal{D}_y \circ g \circ \mathcal{E}_x$  is linear, also  $(\mathbf{G}_{\text{FN}}, \|\cdot\|_{\infty, \text{supp}(\gamma)})$  and hence  $(\mathbf{G}_{\text{FN}}, \|\cdot\|_n)$  is compact. We now show the entropy bounds for  $\mathbf{G}_{\text{FN}}$ .

**Step 1.** Recall  $\text{depth}(g) \leq L$ ,  $\text{depth}(g) \leq p$ ,  $\text{size}(g) \leq s$  and  $\text{mpar}(g) \leq M$  for  $g \in \mathbf{g}_{\text{FN}}$ . We first estimate the entropy  $H(\mathbf{G}_{\text{FN}}, \|\cdot\|_{\infty, \text{supp}(\gamma)}, \delta)$  against the respective entropy of  $\mathbf{g}_{\text{FN}}$ . For  $G, G' \in \mathbf{G}_{\text{FN}}$  and  $g, g' \in \mathbf{g}_{\text{FN}}$  with  $G = \mathcal{D}_y \circ g \circ S_r \circ \mathcal{E}_x$ ,  $G' = \mathcal{D}_y \circ g' \circ S_r \circ \mathcal{E}_x$ , it holds

$$\begin{aligned} \|G - G'\|_{\infty, \sigma_R^r(U)} &= \sup_{x \in \sigma_R^r(U)} \|\mathcal{D}_y \circ g \circ S_r \circ \mathcal{E}_x(x) - \mathcal{D}_y \circ g' \circ S_r \circ \mathcal{E}_x(x)\|_y \\ &\leq \Lambda_{\Psi_y} \sup_{\mathbf{y} \in U} \|g(\mathbf{y}) - g'(\mathbf{y})\|_2 \leq \Lambda_{\Psi_y} \sqrt{p} \|g - g'\|_{\infty, \infty}, \end{aligned} \quad (\text{D.41})$$

where we used  $\sigma_R^r = \mathcal{D}_x \circ S_r^{-1}$  and  $\|\cdot\|_{\infty, \infty}$  from (D.40). Furthermore, we used  $\|g(u)\|_2 \leq \sqrt{p} \|g\|_{\infty}$  for all  $g \in \mathbf{g}_{\text{FN}}$ , since NNs  $g \in \mathbf{g}_{\text{FN}}$  have  $\text{width}(g) \leq p$ . Then (D.41) yields

$$H(\mathbf{G}_{\text{FN}}, \|\cdot\|_{\infty, \sigma_R^r(U)}, \delta) \leq H\left(\mathbf{g}_{\text{FN}}, \|\cdot\|_{\infty, \infty}, \frac{\delta}{\Lambda_{\Psi_y} \sqrt{p}}\right). \quad (\text{D.42})$$

**Step 2.** It remains to bound  $H(\mathbf{g}_{\text{FN}}, \|\cdot\|_{\infty, \infty}, \delta) = \log(N(\mathbf{g}_{\text{FN}}, \|\cdot\|_{\infty, \infty}, \delta))$ . To this end we follow the proof and notation of Schmidt-Hieber (2020a, Lemma 5). For  $l = 1, \dots, L+1$ , define the matrices  $W_l = (w_{i,j}^l)_{i,j} \in \mathbb{R}^{p^{l-1} \times p^l}$  and the vectors  $B_l = (b_j^l)_j \in \mathbb{R}^{p^l}$ . Furthermore, define

$$\begin{aligned} \sigma^{B_l} : \mathbb{R}^{p^l} &\rightarrow \mathbb{R}^{p^l}, \quad \sigma^{B_l}(x) = \sigma_1(x + B_l) = \max\{0, x + B_l\}, \quad l = 1, \dots, L, \\ \sigma^{B_{L+1}} : \mathbb{R}^{p^{L+1}} &\rightarrow \mathbb{R}^{p^{L+1}}, \quad \sigma^{B_{L+1}}(x) = x + B_{L+1}. \end{aligned} \quad (\text{D.43})$$

Then we can write a NN  $g \in \mathbf{g}_{\text{FN}}$  as a functional composition of  $\sigma^{B_l}$  and  $W_l$ , i.e.

$$g : \mathbb{R}^{p_0} \rightarrow \mathbb{R}^{p_{L+1}}, \quad g(x) = \sigma^{B_{L+1}} W_{L+1} \sigma^{B_L} \dots W_2 \sigma^{B_1} W_1 x.$$

For  $k \in \{1, \dots, L+1\}$  we define the functions

$$\begin{aligned} A_k^+ g : \mathbb{R}^{p_0} &\rightarrow \mathbb{R}^{p_k}, & A_k^+ g(x) &= \sigma^{B_k} W_k \dots \sigma^{B_1} W_1 x, \\ A_k^- g : \mathbb{R}^{p_{k-1}} &\rightarrow \mathbb{R}^{p_{L+1}}, & A_k^- g(x) &= \sigma^{B_{L+1}} W_{L+1} \dots \sigma^{B_k} W_k x. \end{aligned} \quad (\text{D.44})$$

Furthermore, set  $A_0^+ g = \text{Id}_{\mathbb{R}^{p_0}}$  and  $A_{L+2}^- g = \text{Id}_{\mathbb{R}^{p_{L+1}}}$ . For all  $1 \leq l \leq L+1$  holds

$$\begin{aligned} \|\sigma^{B_l}(x)\|_\infty &\leq \|x\|_\infty + M \\ \|W^l(x)\|_\infty &\leq \|W^l\|_\infty \|x\|_\infty \leq Mp \|x\|_\infty. \end{aligned}$$

We claim that for  $k \in \{1, \dots, L+1\}$

$$\sup_{x \in [-1, 1]^{p_0}} \|A_k^+ g(x)\|_\infty \leq (M(p+1))^k$$

and proceed by induction. The case  $k=0$  is trivial. To go from  $k-1$  to  $k$  we compute

$$\begin{aligned} \sup_{x \in [-1, 1]^{p_0}} \|A_k^+ x\|_\infty &= \sup_{x \in [-1, 1]^{p_0}} \|\sigma^{B_k} W_k (\sigma^{B_{k-1}} W_{k-1} \dots \sigma^{B_1} W_1 x)\|_\infty \\ &\leq \sup_{x \in [-(M(p+1))^{k-1}, (M(p+1))^{k-1}]^{p_{k-1}}} \|\sigma^{B_k} W^k x\|_\infty \\ &\leq (Mp(M(p+1))^{k-1} + M) \leq (M(p+1))^k, \end{aligned} \quad (\text{D.45})$$

as claimed.

Moreover, for  $l = 1, \dots, L+1$ ,  $W_l : (\mathbb{R}^{p_{l-1}}, \|\cdot\|_\infty) \rightarrow (\mathbb{R}^{p_l}, \|\cdot\|_\infty)$  is Lipschitz with constant  $Mp$  and  $\sigma^{B_l} : (\mathbb{R}^{p_l}, \|\cdot\|_\infty) \rightarrow (\mathbb{R}^{p_l}, \|\cdot\|_\infty)$  is Lipschitz with constant 1. Thus we can estimate the Lipschitz constant of  $A_k^- g$  for  $k = 1, \dots, L+1$ . It holds

$$\begin{aligned} \|A_k^- g(x) - A_k^- g(y)\|_\infty &= \|\sigma^{B_{L+1}} W_{L+1} \dots \sigma^{B_k} W_k x - \sigma^{B_{L+1}} W_{L+1} \dots \sigma^{B_k} W_k y\|_\infty \\ &\leq Mp \|\sigma^{B_L} W_L \dots \sigma^{B_k} W_k x - \sigma^{B_L} W_L \dots \sigma^{B_k} W_k y\|_\infty \\ &\leq \dots \leq (Mp)^{L+2-k} \|x - y\|_\infty \quad \text{for } x, y \in \mathbb{R}^{p_{k-1}}. \end{aligned} \quad (\text{D.46})$$

Now let  $g, g^* \in \mathbf{g}_{\text{FN}}$  be two NN such that  $|w_{i,j}^l - w_{i,j}^{l,*}| < \varepsilon$  and  $|b_i^l - b_i^{l,*}| < \varepsilon$  for all  $i \leq p_{l+1}$ ,  $j \leq p_l$ ,  $l \leq L+1$ . Then

$$\begin{aligned} \|g - g^*\|_{\infty, \infty} &\leq \sum_{k=1}^{L+1} \left\| A_{k+1}^- g \sigma^{B_k} W_k A_{k-1}^+ g^* - A_{k+1}^- g \sigma^{B_k^*} W_k^* A_{k-1}^+ g^* \right\|_{\infty, \infty} \\ &\leq \sum_{k=1}^{L+1} (Mp)^{L+1-k} \left\| \sigma^{B_k} W_k A_{k-1}^+ g^* - \sigma^{B_k^*} W_k^* A_{k-1}^+ g^* \right\|_{\infty, \infty} \\ &\leq \sum_{k=1}^{L+1} (Mp)^{L+1-k} \left( \|(W_k - W_k^*) A_{k-1}^+ g^*\|_{\infty, \infty} + \|B_k - B_k^*\|_\infty \right) \\ &\leq \varepsilon \sum_{k=1}^{L+1} (Mp)^{L+1-k} \left( pM^{k-1} (p+1)^{k-1} + 1 \right) \\ &< \varepsilon (L+1) M^L (p+1)^{L+1}, \end{aligned} \quad (\text{D.47})$$

where we used (D.45), (D.46) and  $M \geq 1$ . The total number of weight and biases is less than  $(L+1)(p^2+p)$ . Therefore there are at most

$$\binom{(L+1)(p^2+p)}{s} \leq ((L+1)(p^2+p))^s$$

combinations to pick  $s$  nonzero parameters. Since all parameters are bounded by  $M$ , we choose  $\varepsilon = \delta/((L+1)M^L(p+1)^{L+1})$  and obtain the covering bound for all  $\delta > 0$

$$\begin{aligned} N(\mathbf{g}_{\text{FN}}, \|\cdot\|_{\infty, \infty}, \delta) &\leq \max \left\{ 1, \sum_{s^*=1}^s (2M\varepsilon^{-1}(L+1)(p^2+p))^{s^*} \right\} \\ &\leq \max \left\{ 1, \sum_{s^*=1}^s (2\delta^{-1}(L+1)M^{L+1}(p+1)^{L+1}(L+1)(p^2+p))^{s^*} \right\} \\ &\leq \max \left\{ 1, \sum_{s^*=1}^s (2\delta^{-1}(L+1)^2M^{L+1}(p+1)^{L+3})^{s^*} \right\} \\ &\leq (2^{L+6}L^2M^{L+1}p^{L+3} \max\{1, \delta^{-1}\})^{s+1}, \end{aligned} \quad (\text{D.48})$$

where we used  $L \geq 1$  and  $p \geq 1$  at the last inequality. Eqs. (D.48) and (D.42) show (3.16).

Applying (3.16) to the sparse FrameNet class  $\mathbf{G}_{\text{FN}}^{\text{SP}}(\sigma_1, N)$  gives

$$\begin{aligned} &H(\mathbf{G}_{\text{FN}}^{\text{SP}}(\sigma_1, N), \|\cdot\|_{\infty, \sigma_R^r(U)}, \delta) \\ &\leq (\text{size}_N + 1) \log \left( 2^{\text{depth}_N+6} \Lambda_{\Psi_y} \text{depth}_N^2 M^{\text{depth}_N+1} \text{width}_N^{\text{depth}_N+4} \max\{1, \delta^{-1}\} \right) \\ &\leq (C_s N + 1) \\ &\quad \times \log \left( 2^{C_L \log(N)+6} \Lambda_{\Psi_y} (C_L \log(N))^2 M^{C_L \log(N)+1} (C_p N)^{C_L \log(N)+4} \max\{1, \delta^{-1}\} \right) \\ &\leq C_H^{\text{SP}} N (1 + \log(N)^2 + \log(\max\{1, \delta^{-1}\})), \quad N \in \mathbb{N}, \quad \delta > 0, \end{aligned} \quad (\text{D.49})$$

where we defined

$$\begin{aligned} C_H^{\text{SP}} &= 2C_s \left( (C_L + 6) \log(2) + \log(\Lambda_{\Psi_y}) + C_L^2 + \right. \\ &\quad \left. + (C_L + 1) \log(M) + (C_L + 4)(\log(C_p) + 1) \right). \end{aligned}$$

Applying (3.16) to the fully connected FrameNet class  $\mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_1, N)$  gives

$$\begin{aligned} &H(\mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_1, N), \|\cdot\|_{\infty, \sigma_R^r(U)}, \delta) \\ &\leq (s^{\text{FC}}(N) + 1) \log \left( 2^{\text{depth}_N+6} \Lambda_{\Psi_y} \text{depth}_N^2 M^{\text{depth}_N+1} \text{width}_N^{\text{depth}_N+4} \max\{1, \delta^{-1}\} \right) \\ &\leq ((\text{depth}_N + 1) (\text{width}_N^2 + \text{width}_N) + 1) \\ &\quad \times \log \left( 2^{C_L \log(N)+6} \Lambda_{\Psi_y} (C_L \log(N))^2 M^{C_L \log(N)+1} (C_p N)^{C_L \log(N)+4} \max\{1, \delta^{-1}\} \right) \\ &\leq ((C_L \log(N) + 1) (C_p^2 N^2 + C_p N) + 1) \\ &\quad \times \log \left( 2^{C_L \log(N)+6} \Lambda_{\Psi_y} (C_L \log(N))^2 M^{C_L \log(N)+1} (C_p N)^{C_L \log(N)+4} \max\{1, \delta^{-1}\} \right) \\ &\leq C_H^{\text{FC}} N^2 (1 + \log(N)^3 + \log(\max\{1, \delta^{-1}\})), \quad N \in \mathbb{N}, \end{aligned} \quad (\text{D.50})$$

where we defined

$$C_H^{\text{FC}} = 8C_L C_p^2 \left( (C_L + 6) \log(2) + \log(\Lambda_{\Psi_y}) + C_L^2 + (C_L + 1) \log(M) + (C_L + 4)(\log(C_p) + 1) \right).$$

Equations (D.49) and (D.50) finish the proof of Lemma 20.

## D.5 Proof of Lemma 22

The following proof is a modification of Schmidt-Hieber (2020a, Proof of Lemma 5) to the case where the activation function is not globally, but only locally Lipschitz continuous. The compactness of  $(\mathbf{G}, \|\cdot\|_{\infty, \text{supp}(\gamma)})$  follows similarly to the ReLU case since Petersen et al. (2021, Proposition 3.5) holds for any continuous activation function.

Let  $q \in \mathbb{N}$ ,  $q \geq 2$  and let  $\sigma_q : \mathbb{R} \rightarrow \mathbb{R}$ ,  $\sigma_q(x) = \max\{0, x\}^q$  denote the RePU activation function. Recall  $\text{depth}(g) \leq L$ ,  $\text{width}(g) \leq p$ ,  $\text{size}(g) \leq s$  and  $\text{mpar}(g) \leq M$  for  $g \in \mathbf{g}_{\text{FN}}$ . We argue analogously to the ReLU-case in Lemma 20 and bound the entropy of the NN class  $\mathbf{g}_{\text{FN}}(\sigma_q, L, p, s, M, B)$ . Recall the definitions of  $\sigma^{B_l}$ ,  $A_k^+ g$  and  $A_k^- g$  from (D.43)-(D.44). Similar to (D.45) it holds that

$$\begin{aligned} \|A_k^+ g\|_{\infty, \infty} &= \sup_{x \in [-1, 1]^{p_0}} \|A_k^+ g(x)\|_{\infty} \\ &\leq \sup_{x \in [-M(p+1), M(p+1)]^{p_1}} \|\sigma^{B_k} W_k \sigma^{B_{k-1}} \dots W_2 \sigma_q x\|_{\infty} \\ &\leq \sup_{x \in [-M^q(p+1)^q, M^q(p+1)^q]^{p_1}} \|\sigma^{B_k} W_k \sigma^{B_{k-1}} \dots W_2 x\|_{\infty} \\ &\leq \dots \leq (M(p+1))^{\sum_{j=1}^k q^j} \leq (M(p+1))^{q^{k+1}}, \end{aligned}$$

where we used  $M \geq 1$ .

In the RePU-case,  $A_k^- g$  is only locally Lipschitz: Since  $|\sigma_q(x)'| \leq q|x|^{q-1}$  it holds

$$|\sigma_q(x) - \sigma_q(y)| \leq q \max\{|x|, |y|\}^{q-1} |x - y| \quad \forall x, y \in \mathbb{R}.$$

Therefore for  $k = 1, \dots, L+1$  and  $x, y \in \mathbb{R}^{p_{k-1}}$ ,  $\|x\|_{\infty}, \|y\|_{\infty} \leq C$ , we get

$$\begin{aligned} &\|A_k^- g(x) - A_k^- g(y)\|_{\infty} \\ &= \|\sigma^{B_{L+1}} W_{L+1} \sigma^{B_L} \dots W_k x - \sigma^{B_{L+1}} W_{L+1} \sigma^{B_L} \dots W_k y\|_{\infty} \\ &\leq M p \|\sigma^{B_L} W_L \sigma^{B_{L-1}} \dots W_k x - \sigma^{B_L} W_L \sigma^{B_{L-1}} \dots W_k y\|_{\infty} \\ &\leq M p q \left( \sup_{\|x\|_{\infty} \leq C} \|W_L \sigma^{B_{L-1}} \dots W_k x\|_{\infty} \right)^{q-1} \\ &\quad \times \|W_L \sigma^{B_{L-1}} \dots W_k x - W_L \sigma^{B_{L-1}} \dots W_k y\|_{\infty} \\ &\leq (M p q)^{L+2-k} \left( \sup_{\|x\|_{\infty} \leq C} \|W_L \sigma^{B_{L-1}} \dots W_k x\|_{\infty} \right)^{q-1} \\ &\quad \times \left( \sup_{\|x\|_{\infty} \leq C} \|W_{L-1} \sigma^{B_{L-2}} \dots W_k x\|_{\infty} \right)^{q-1} \times \dots \times \left( \sup_{\|x\|_{\infty} \leq C} \|W_k x\|_{\infty} \right)^{q-1} \|x - y\|_{\infty}. \end{aligned}$$

Using

$$\begin{aligned}
 \sup_{\|x\|_\infty \leq C} \|W_j \sigma^{B_{j-1}} \dots W_k x\|_\infty &\leq \sup_{\|x\|_\infty \leq (M(p+1)C)^q} \|W_j \sigma^{B_{j-1}} \dots W_{k+1} x\|_\infty \\
 &\leq \sup_{\|x\|_\infty \leq (M(p+1))^{q+q^2} C^{q^2}} \|W_j \sigma^{B_{j-1}} \dots W_{k+2} x\|_\infty \\
 &\leq \dots \leq \sup_{\|x\|_\infty \leq (M(p+1))^{\sum_{l=1}^{j-k} q^l} C^{q^{j-k}}} \|W_j x\|_\infty \\
 &\leq (M(p+1))^{\sum_{l=0}^{j-k} q^l} C^{q^{j-k}} \leq (M(p+1)C)^{q^{j-k+1}}
 \end{aligned}$$

for  $j = k, \dots, L$  and  $C, M \geq 1$ , we get

$$\begin{aligned}
 \|A_k^- g(x) - A_k^- g(y)\|_\infty &\leq (Mpq)^{L+2-k} \prod_{j=k}^L \left( M(p+1) \hat{C} \right)^{q^{j-k+1}} \|x - y\|_\infty \quad (\text{D.51}) \\
 &\leq (Mpq)^{L+2-k} \left( M(p+1) \hat{C} \right)^{q^{L+2-k}} \|x - y\|_\infty, \quad x, y \in \mathbb{R}^{p_{k-1}}.
 \end{aligned}$$

Now we proceed similarly to (D.47). Let  $g, g^* \in \mathbf{g}_{\text{FN}}$  be two NN such that  $|w_{i,j}^l - w_{i,j}^{l,*}| < \varepsilon$  and  $|b_i^l - b_i^{l,*}| < \varepsilon$  for all  $i \leq p_{l+1}$ ,  $j \leq p_l$ ,  $l \leq L+1$ . Then with  $A_0^+ g = \text{Id}_{\mathbb{R}^{p_0}}$  and  $A_{L+2}^- g = \text{Id}_{\mathbb{R}^{p_{L+1}}}$  we estimate

$$\begin{aligned}
 &\|g - g^*\|_{\infty, \infty} \\
 &\leq \sum_{k=1}^{L+1} \left\| A_{k+1}^- g \sigma^{B_k} W_k A_{k-1}^+ g^* - A_{k+1}^- g \sigma^{B_k^*} W_k^* A_{k-1}^+ g^* \right\|_{\infty, \infty} \\
 &\leq \sum_{k=1}^{L+1} (Mpq)^{L+1-k} \left( M(p+1) (M(p+1))^{q^{k+1}} \right)^{q^{L+1-k}} \\
 &\quad \left\| \sigma^{B_k} W_k A_{k-1}^+ g^* - \sigma^{B_k^*} W_k^* A_{k-1}^+ g^* \right\|_{\infty, \infty} \\
 &\leq \sum_{k=1}^{L+1} (Mpq)^{L+1-k} \left( M(p+1) (M(p+1))^{q^{k+1}} \right)^{q^{L+1-k}} \\
 &\quad \left( \|(W_k - W_k^*) A_{k-1}^+ g^*\|_{\infty, \infty} + \|B_k - B_k^*\|_{\infty, \infty} \right) q \left( M(p+1) \|A_{k-1}^+ g^*\|_{\infty, \infty} \right)^{q-1} \\
 &\leq 2\varepsilon \sum_{k=1}^{L+1} (Mpq)^{L+2-k} \left( M(p+1) (M(p+1))^{q^{k+1}} \right)^{q^{L+1-k}} \|A_{k-1}^+ g^*\|_{\infty, \infty} \\
 &\quad \left( M(p+1) \|A_{k-1}^+ g^*\|_{\infty, \infty} \right)^{q-1} \\
 &\leq 2\varepsilon (L+1) (M(p+1)q)^{L+q} \left( M(p+1) (M(p+1))^{q^{L+2}} \right)^{q^L} \left( (M(p+1))^{q^{L+1}} \right)^q \\
 &< \varepsilon L q^{L+q} (2pM)^{4q^{2L+2}} (2M\sqrt{p}(p^2+p)(L+1))^{-1}. \quad (\text{D.52})
 \end{aligned}$$

In (D.52) we used the Lipschitz bound (D.51) with

$$C = \max \left\{ 1, \|\sigma^{B_k} W_k A_{k-1}^+ g^*\|_{\infty, \infty}, \|\sigma^{B_k^*} W_k^* A_{k-1}^+ g^*\|_{\infty, \infty} \right\} \leq (M(p+1))^{q^{k+1}},$$

and  $p \geq 1$ ,  $L \geq 1$ ,  $q \geq 2$  at the last inequality.

As in the proof of Lemma 20, there are

$$\binom{(L+1)(p^2+p)}{s} \leq ((L+1)(p^2+p))^s$$

combinations to pick  $s$  nonzero weights and biases. Since all parameters are bounded by  $M$ , we choose

$$\varepsilon = \frac{2M\sqrt{p}(p^2+p)(L+1)\delta}{Lq^{L+q}(2pM)^{4q^{2L+2}}}$$

and obtain the covering bound

$$\begin{aligned} N(\mathbf{g}_{\text{FN}}, \|\cdot\|_{\infty, \infty}, \delta) &\leq \max \left\{ 1, \sum_{s^*=1}^s (2M\epsilon^{-1}(L+1)(p^2+p))^{s^*} \right\} \\ &\leq \max \left\{ 1, \sum_{s^*=1}^s \left( Lq^{L+q} (2pM)^{4q^{2L+2}} (\sqrt{p}\delta)^{-1} \right)^{s^*} \right\} \\ &\leq \left( Lq^{L+q} (2pM)^{4q^{2L+2}} \sqrt{p}^{-1} \max\{1, \delta^{-1}\} \right)^{s+1}. \end{aligned} \quad (\text{D.53})$$

Eqs. (D.53) and (D.42) show (3.17).

Applying (3.17) to the sparse FrameNet class  $\mathbf{G}_{\text{FN}}^{\text{SP}}(\sigma_q, N)$  gives

$$\begin{aligned} &H(\mathbf{G}_{\text{FN}}^{\text{SP}}(\sigma_q, N), \|\cdot\|_{\infty, \sigma_R^r(U)}, \delta) \\ &\leq (s^{\text{SP}}(N) + 1) \log \left( \Lambda_{\Psi_y} \text{depth}_N q^{\text{depth}_N+q} (2\text{width}_N M)^{4q^{2\text{depth}_N+2}} \max\{1, \delta^{-1}\} \right) \\ &\leq (C_s N + 1) \log \left( \Lambda_{\Psi_y} C_L \log(N) q^{C_L \log(N)+q} (2C_p N M)^{4q^{2C_L \log(N)+2}} \max\{1, \delta^{-1}\} \right) \\ &\leq C_H^{\text{SP}} N^{1+2C_L \log(q)} (1 + \log(N) + \log(\max\{1, \delta^{-1}\})), \end{aligned} \quad (\text{D.54})$$

where we set

$$C_H^{\text{SP}} = 2C_s (\log(\Lambda_{\Psi_y}) + C_L + (C_L + q) \log(q) + 4q^2 (\log(2C_p M) + 1)).$$

Applying (3.17) to the fully connected FrameNet class  $\mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_q, N)$  gives the entropy bound

$$\begin{aligned} &H(\mathbf{G}_{\text{FN}}^{\text{full}}(\sigma_q, N), \|\cdot\|_{\infty, \sigma_R^r(U)}, \delta) \\ &\leq (s^{\text{FC}}(N) + 1) \log \left( \Lambda_{\Psi_y} \text{depth}_N q^{\text{depth}_N+q} (2\text{width}_N M)^{4q^{2\text{depth}_N+2}} \max\{1, \delta^{-1}\} \right) \\ &\leq ((\text{depth}_N + 1) (\text{width}_N^2 + \text{width}_N) + 1) \\ &\quad \times \log \left( \Lambda_{\Psi_y} \text{depth}_N q^{\text{depth}_N+q} (2\text{width}_N M)^{4q^{2\text{depth}_N+2}} \max\{1, \delta^{-1}\} \right) \\ &\leq ((C_L \log(N) + 1) (C_p^2 N^2 + C_p N) + 1) \\ &\quad \times \log \left( \Lambda_{\Psi_y} C_L \log(N) q^{C_L \log(N)+q} (2C_p N M)^{4q^{2C_L \log(N)+2}} \max\{1, \delta^{-1}\} \right) \\ &\leq C_H^{\text{FC}} N^{2+2C_L \log(q)} (1 + \log(N)^2 + \log(\max\{1, \delta^{-1}\})), \end{aligned} \quad (\text{D.55})$$

where we set

$$C_H^{\text{FC}} = 8C_s C_p^2 \left( \log(\Lambda_{\Psi_y}) + C_L + (C_L + q) \log(q) + 4q^2 (\log(2C_p M) + 1) \right).$$

Equations (D.54) and (D.55) finish the proof of Lemma 22.

## Appendix E. Proofs of Section 4

### E.1 Proof of Theorem 27

In Herrmann et al. (2024, Proof of Proposition 3, Step 1), the holomorphy in Assumption 2 is verified for  $\mathcal{X}, \mathcal{Y}$  in (4.5) with  $r_0 > d/2$  and  $t \in [0, (1 + r_0 - d/2 - t_0)/d]$ . Moreover,  $\gamma = (\sigma_R^r)_{\#}\pi$  in particular shows  $\text{supp}(\gamma) \subseteq C_R^r(\mathcal{X})$  and hence verifies the second part of Assumption 2. Substituting  $\mathfrak{s} = r_0 + rd$ , i.e.  $r = \frac{\mathfrak{s} - r_0}{d}$ , and taking  $t = (1 + r_0 - d/2 - t_0)/d - \tau$  with some small  $\tau$ , Theorem 23 (i) then gives

$$\mathbb{E}_{G_0} [\|\hat{G}_n - G_0\|_{L^2(\gamma)}^2] \leq Cn^{-\frac{\kappa}{\kappa+1} + \tau},$$

where

$$\kappa = 2 \min \left\{ \frac{\mathfrak{s} - r_0}{d} - \frac{1}{2}, \frac{1 + r_0 - \frac{d}{2} - t_0}{d} \right\} - \tau$$

for all  $r_0 > d/2$  and  $t_0 \in [0, 1]$ .

From here on the proof is essentially the same as Herrmann et al. (2024, Proof of Proposition 3, Step 2); the only difference is that while Herrmann et al. (2024) uses the uniform bound in Theorem 18 (i), we require the  $L^2$ -bound in Theorem 18 (ii). For completeness, we repeat the argument. The constraint  $r > 1$  implies  $\mathfrak{s} > r_0 + d$  on  $\mathfrak{s}$ . We now choose  $r_0 > \frac{d}{2}$  in order to maximize the convergence rate. Solving

$$\frac{\mathfrak{s} - r_0}{d} - \frac{1}{2} = \frac{1 + r_0 - \frac{d}{2} - t_0}{d}$$

for  $r_0$  gives

$$r_0 = \frac{\mathfrak{s} + t_0 - 1}{2}. \tag{E.1}$$

The constraint  $r_0 > \frac{d}{2}$  implies the constraint  $\mathfrak{s} > d + 1 - t_0$ .

We look at two cases separately. First, if  $\mathfrak{s} \in (\frac{3d}{2}, 2d + 1 - t_0]$ , we set  $r_0 := \frac{d}{2} + \tau_2$ , where we choose  $\tau_2 > 0$  s.t.  $\tau_2 < \mathfrak{s} - 3d/2$  which guarantees  $\mathfrak{s} > r_0 + d$ . For  $\tau < \tau_2/d$ , we obtain the convergence rate

$$\kappa = 2 \min \left\{ \frac{\mathfrak{s} - \frac{d}{2} - \tau_2}{d} - \frac{1}{2}, \frac{1 + \frac{d}{2} + \tau_2 - \frac{d}{2} - t_0}{d} \right\} - \tau \geq 2 \min \left\{ \frac{\mathfrak{s}}{d} - 1 - \frac{2\tau_2}{d}, \frac{1 - t_0}{d} \right\}.$$

In the case  $\mathfrak{s} > 2d + 1 - t_0$ , define  $r_0$  as in (E.1). The constraint  $\mathfrak{s} > r_0 + d$  amounts to

$$\mathfrak{s} > \frac{\mathfrak{s} + t_0 - 1}{2} + d \quad \Leftrightarrow \quad \mathfrak{s} > 2d + t_0 - 1,$$

which holds since  $\mathfrak{s} > 2d + 1 - t_0 \geq 2d + t_0 - 1$  for all  $t_0 \in [0, 1]$ . In this case we get the convergence rate

$$\kappa = 2 \frac{\mathfrak{s} - r_0}{d} - 1 - \tau = \frac{\mathfrak{s} + 1 - t_0}{d} - 1 - \tau.$$

Choosing  $\tau_1 > 8\tau_2/d > 8\tau$  shows (4.6) and finishes the proof of Theorem 27.

## References

- Milton Abramowitz and Irene A Stegun. *Handbook of mathematical functions with formulas, graphs, and mathematical tables*, volume 55. US Government printing office, 1948.
- Ben Adcock, Nick Dexter, and Sebastian Moraga. Optimal deep learning of holomorphic operators between banach spaces. *Advances in Neural Information Processing Systems*, 37:27725–27789, 2024.
- Ben Adcock, Simone Brugiapaglia, Nick Dexter, and Sebastian Moraga. Near-optimal learning of banach-valued, high-dimensional functions via deep neural networks. *Neural Networks*, 181:106761, 2025. ISSN 0893-6080. doi: <https://doi.org/10.1016/j.neunet.2024.106761>. URL <https://www.sciencedirect.com/science/article/pii/S0893608024006853>.
- Sergios Agapiou and Sven Wang. Laplace priors and spatial inhomogeneity in bayesian inverse problems. *Bernoulli*, 30(2):878–910, 2024.
- Anima Anandkumar, Kamyar Azizzadenesheli, Kaushik Bhattacharya, Nikola Kovachki, Zongyi Li, Burigede Liu, and Andrew Stuart. Neural operator: Graph kernel network for partial differential equations. In *ICLR 2020 Workshop on Integration of Deep Neural Models and Differential Equations*, 2019. URL <https://openreview.net/forum?id=fg2ZFmXF03>.
- Ivo Babuška, Fabio Nobile, and Raúl Tempone. A stochastic collocation method for elliptic partial differential equations with random input data. *SIAM Rev.*, 52(2):317–355, 2010. ISSN 0036-1445,1095-7200. doi: 10.1137/100786356. URL <https://doi.org/10.1137/100786356>.
- Markus Bachmayr, Albert Cohen, Ronald DeVore, and Giovanni Migliorati. Sparse polynomial approximation of parametric elliptic PDEs. Part II: Lognormal coefficients. *ESAIM Math. Model. Numer. Anal.*, 51(1):341–363, 2017. ISSN 2822-7840,2804-7214. doi: 10.1051/m2an/2016051. URL <https://doi.org/10.1051/m2an/2016051>.
- Andrew Barron, Lucien Birgé, and Pascal Massart. Risk bounds for model selection via penalization. *Probability theory and related fields*, 113:301–413, 1999.
- Sebastian Becker, Arnulf Jentzen, Marvin S Müller, and Philippe von Wurstemberger. Learning the random variables in monte carlo simulations with stochastic gradient descent: Machine learning for parametric pdes and financial derivative pricing. *Mathematical Finance*, 34(1):90–150, 2024.
- Kaushik Bhattacharya, Bamdad Hosseini, Nikola B. Kovachki, and Andrew M. Stuart. Model reduction and neural networks for parametric PDEs. *SMAI J. Comput. Math.*, 7: 121–157, 2021.
- Lucien Birgé and Pascal Massart. Rates of convergence for minimum contrast estimators. *Probability Theory and Related Fields*, 97:113–150, 1993.

- Ismaël Castillo and Richard Nickl. Nonparametric bernstein–von mises theorems in gaussian white noise. *The Annals of Statistics*, 41(4), 2013. ISSN 0090-5364. doi: 10.1214/13-AOS1133.
- Gaëlle Chagny, Anouar Meynaoui, and Angelina Roche. Adaptive nonparametric estimation in the functional linear model with functional output. *Electronic Journal of Statistics*, 19(1):2990–3039, 2025.
- Abdellah Chkifa, Albert Cohen, Ronald DeVore, and Christoph Schwab. Sparse adaptive Taylor approximation algorithms for parametric and stochastic elliptic PDEs. *ESAIM Math. Model. Numer. Anal.*, 47(1):253–280, 2013. ISSN 0764-583X. doi: 10.1051/m2an/2012027. URL <https://doi.org/10.1051/m2an/2012027>.
- Abdellah Chkifa, Albert Cohen, Giovanni Migliorati, Fabio Nobile, and Raul Tempone. Discrete least squares polynomial approximation with random evaluations—application to parametric and stochastic elliptic PDEs. *ESAIM Math. Model. Numer. Anal.*, 49(3): 815–837, 2015a. ISSN 0764-583X. doi: 10.1051/m2an/2014050. URL <https://doi.org/10.1051/m2an/2014050>.
- Abdellah Chkifa, Albert Cohen, and Christoph Schwab. Breaking the curse of dimensionality in sparse polynomial approximation of parametric PDEs. *J. Math. Pures Appl. (9)*, 103(2):400–428, 2015b. ISSN 0021-7824,1776-3371. doi: 10.1016/j.matpur.2014.04.009. URL <https://doi.org/10.1016/j.matpur.2014.04.009>.
- Ole Christensen. *An Introduction to Frames and Riesz Bases [recurso electrónico]*. Applied and Numerical Harmonic Analysis. Birkhauser, Cham, second edition, 2016. ISBN 978-3-319-25613-9.
- Ludovica Cicci, Stefania Fresca, and Andrea Manzoni. Deep-hyromnet: A deep learning-based operator approximation for hyper-reduction of nonlinear parametrized pdes. *Journal of Scientific Computing*, 93(2):57, 2022.
- K. A. Cliffe, M. B. Giles, R. Scheichl, and A. L. Teckentrup. Multilevel Monte Carlo methods and applications to elliptic PDEs with random coefficients. *Comput. Vis. Sci.*, 14(1):3–15, 2011. ISSN 1432-9360,1433-0369. doi: 10.1007/s00791-011-0160-x. URL <https://doi.org/10.1007/s00791-011-0160-x>.
- Albert Cohen and Ronald DeVore. Approximation of high-dimensional parametric PDEs. *Acta Numer.*, 24:1–159, 2015. ISSN 0962-4929. doi: 10.1017/S0962492915000033. URL <https://doi.org/10.1017/S0962492915000033>.
- Albert Cohen, Ronald DeVore, and Christoph Schwab. Convergence rates of best  $N$ -term Galerkin approximations for a class of elliptic sPDEs. *Found. Comput. Math.*, 10(6): 615–646, 2010. ISSN 1615-3375. doi: 10.1007/s10208-010-9072-2. URL <http://dx.doi.org/10.1007/s10208-010-9072-2>.
- Albert Cohen, Ronald Devore, and Christoph Schwab. Analytic regularity and polynomial approximation of parametric and stochastic elliptic PDE’s. *Anal. Appl. (Singap.)*, 9(1):

- 11–47, 2011. ISSN 0219-5305. doi: 10.1142/S0219530511001728. URL <http://dx.doi.org/10.1142/S0219530511001728>.
- Albert Cohen, Giovanni Migliorati, and Fabio Nobile. Discrete least-squares approximations over optimized downward closed polynomial spaces in arbitrary dimension. *Constr. Approx.*, 45(3):497–519, 2017. ISSN 0176-4276. doi: 10.1007/s00365-017-9364-8. URL <https://doi.org/10.1007/s00365-017-9364-8>.
- Albert Cohen, Christoph Schwab, and Jakob Zech. Shape holomorphy of the stationary navier-stokes equations. *SIAM J. Math. Analysis*, 50(2):1720–1752, 2018. doi: <https://doi.org/10.1137/16M1099406>.
- Niccolò Dal Santo, Simone Deparis, and Luca Pegolotti. Data driven approximation of parametrized pdes by reduced basis and neural networks. *Journal of Computational Physics*, 416:109550, 2020.
- Oleg Davydov and Rob Stevenson. Hierarchical riesz bases for  $h^s(\omega)$ ,  $1 < s < 5/2$ . *Constructive Approximation*, 22(3):365–394, 2005. ISSN 1432-0940. doi: 10.1007/s00365-004-0593-2.
- Maarten V. de Hoop, Nikola B. Kovachki, Nicholas H. Nelsen, and Andrew M. Stuart. Convergence rates for learning linear operators from noisy data. *SIAM/ASA Journal on Uncertainty Quantification*, 11(2):480–513, 2023. doi: 10.1137/21M1442942. URL <https://doi.org/10.1137/21M1442942>.
- Tim De Ryck, Samuel Lanthaler, and Siddhartha Mishra. On the approximation of functions by tanh neural networks. *Neural Networks*, 143:732–750, 2021.
- R. DeVore, R. Howard, and C. Micchelli. Optimal nonlinear approximation. *Manuscripta mathematica*, 63(4):469–478, 1989. URL <http://eudml.org/doc/155392>.
- R.A. DeVore and G.G. Lorentz. *Constructive Approximation*. Grundlehren der mathematischen Wissenschaften. Springer Berlin Heidelberg, 1993. ISBN 9783540506270. URL [https://books.google.de/books?id=cDqNW6k7\\_ZwC](https://books.google.de/books?id=cDqNW6k7_ZwC).
- Sjoerd Dirksen. Tail bounds via generic chaining. *Electronic Journal of Probability*, 20, 2015. ISSN 1083-6489. doi: 10.1214/EJP.v20-3760.
- Alireza Doostan and Houman Owhadi. A non-adapted sparse approximation of pdes with stochastic inputs. *Journal of Computational Physics*, 230(8):3015–3034, 2011. ISSN 0021-9991. doi: <https://doi.org/10.1016/j.jcp.2011.01.002>. URL <https://www.sciencedirect.com/science/article/pii/S0021999111000106>.
- Dinh Dung, Van Kien Nguyen, Christoph Schwab, and Jakob Zech. *Analyticity and sparsity in uncertainty quantification for PDEs with Gaussian random field inputs*, volume 2334 of *Lecture Notes in Mathematics*. Springer, Cham, 2023. ISBN 978-3-031-38383-0. doi: 10.1007/978-3-031-38384-7. URL <https://doi.org/10.1007/978-3-031-38384-7>.

- D. E. Edmunds and H. Triebel. *Function spaces, entropy numbers, differential operators*, volume 120 of *Cambridge Tracts in Mathematics*. Cambridge University Press, Cambridge, 1996. ISBN 0-521-56036-5. doi: 10.1017/CBO9780511662201.
- Dennis Elbrachter, Dmytro Perekrestenko, Philipp Grohs, and Helmut Bolcskei. Deep neural network approximation theory. *IEEE Transactions on Information Theory*, 67(5):2581–2623, 2021. ISSN 0018-9448. doi: 10.1109/TIT.2021.3062161.
- Utku Evci, Fabian Pedregosa, Aidan Gomez, and Erich Elsen. The difficulty of training sparse neural networks. *arXiv preprint arXiv:1906.10732*, 2019.
- Evarist Giné and Richard Nickl. *Mathematical foundations of infinite-dimensional statistical models*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, New York, 2016.
- Evarist Giné and Richard Nickl. *Mathematical foundations of infinite-dimensional statistical models*. Cambridge series in statistical and probabilistic mathematics. Cambridge University Press, New York (NY), 2016. ISBN 1107043166.
- Sonja Greven and Fabian Scheipl. A general framework for functional regression modelling. *Statistical Modelling*, 17(1-2):1–35, 2017.
- Helmut Harbrecht, Michael Peters, and Markus Siebenmorgen. Analysis of the domain mapping method for elliptic diffusion problems on random domains. *Numerische Mathematik*, 134(4):823–856, 2016.
- Fernando Henríquez and Christoph Schwab. Shape holomorphy of the calderón projector for the laplacian in  $\mathbb{R}^2$ . *Integral Equations and Operator Theory*, 93(4):43, 2021.
- Lukas Herrmann, Christoph Schwab, and Jakob Zech. Neural and spectral operator surrogates: unified construction and expression rate bounds. *Advances in Computational Mathematics*, 50(4):1–43, 2024. ISSN 1019-7168. doi: 10.1007/s10444-024-10171-2. URL <https://link.springer.com/article/10.1007/s10444-024-10171-2>.
- Jan S. Hesthaven, Gianluigi Rozza, and Benjamin Stamm. *Certified reduced basis methods for parametrized partial differential equations*. SpringerBriefs in Mathematics. Springer, Cham; BCAM Basque Center for Applied Mathematics, Bilbao, 2016. ISBN 978-3-319-22469-5. doi: 10.1007/978-3-319-22470-1. URL <https://doi.org/10.1007/978-3-319-22470-1>. BCAM SpringerBriefs.
- J.S. Hesthaven and S. Ubbiali. Non-intrusive reduced order modeling of nonlinear problems using neural networks. *Journal of Computational Physics*, 363:55–78, 2018. ISSN 0021-9991. doi: <https://doi.org/10.1016/j.jcp.2018.02.037>. URL <https://www.sciencedirect.com/science/article/pii/S0021999118301190>.
- R. Hiptmair, L. Scarabosio, C. Schillings, and Ch. Schwab. Large deformation shape uncertainty quantification in acoustic scattering. *Advances in Computational Mathematics*, 44(5):1475–1518, 10 2018. ISSN 1572-9044. doi: 10.1007/s10444-018-9594-8. URL <https://doi.org/10.1007/s10444-018-9594-8>.

- Viet Ha Hoang and Christoph Schwab. N-term wiener chaos approximation rates for elliptic pdes with lognormal gaussian random inputs. *Mathematical Models and Methods in Applied Sciences*, 24(04):797–826, 2014. doi: 10.1142/S0218202513500681. URL <https://doi.org/10.1142/S0218202513500681>.
- Carlos Jerez-Hanckes, Christoph Schwab, and Jakob Zech. Electromagnetic wave scattering by random surfaces: shape holomorphy. *Math. Models Methods Appl. Sci.*, 27(12):2229–2259, 2017. ISSN 0218-2025. doi: 10.1142/S0218202517500439. URL <https://doi.org/10.1142/S0218202517500439>.
- Padraig Kirwan. *Complexifications of multilinear and polynomial mappings*. PhD thesis, University College Dublin, 1997. Ph.D. thesis, National University of Ireland, Galway.
- Nikola B Kovachki, Samuel Lanthaler, and Hrushikesh Mhaskar. Data complexity estimates for operator learning. *arXiv preprint arXiv:2405.15992*, 2024a.
- Nikola B Kovachki, Samuel Lanthaler, and Andrew M Stuart. Operator learning: Algorithms and analysis. *Handbook of Numerical Analysis*, 25:419–467, 2024b.
- Fabian Kröpfl, Roland Maier, and Daniel Peterseim. Operator compression with deep neural networks. *Advances in Continuous and Discrete Models*, 2022(1):29, 2022.
- Gitta Kutyniok, Philipp Petersen, Mones Raslan, and Reinhold Schneider. A theoretical analysis of deep neural networks and parametric PDEs. *Constr. Approx.*, 55(1):73–125, 2022. ISSN 0176-4276,1432-0940. doi: 10.1007/s00365-021-09551-4. URL <https://doi.org/10.1007/s00365-021-09551-4>.
- Samuel Lanthaler. Operator learning with pca-net: upper and lower complexity bounds. *Journal of Machine Learning Research*, 24(318):1–67, 2023. URL <http://jmlr.org/papers/v24/23-0478.html>.
- Samuel Lanthaler, Siddhartha Mishra, and George E. Karniadakis. Error estimates for deeponets: a deep learning framework in infinite dimensions. *Transactions of Mathematics and Its Applications*, 6(1), 2022. doi: 10.1093/imatrm/tnac001.
- Bo Li. Better approximations of high dimensional smooth functions by deep neural networks with rectified power units. *Communications in Computational Physics*, 27(2):379–411, 2020. ISSN 1815-2406. doi: 10.4208/cicp.OA-2019-0168.
- Bo Li, Shanshan Tang, and Haijun Yu. Powernet: Efficient representations of polynomials and smooth functions by deep neural networks with rectified power units. *arXiv preprint arXiv:1909.05136*, 2019.
- Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar. Fourier neural operator for parametric partial differential equations. *arXiv preprint arXiv:2010.08895*, 2020.
- Hao Liu, Haizhao Yang, Minshuo Chen, Tuo Zhao, and Wenjing Liao. Deep nonparametric estimation of operators between infinite dimensional spaces. *Journal of Machine Learning Research*, 25(24):1–67, 2024.

- Lu Lu, Pengzhan Jin, Guofei Pang, Zhongqiang Zhang, and George Em Karniadakis. Learning nonlinear operators via deepnet based on the universal approximation theorem of operators. *Nature Machine Intelligence*, 3(3):218–229, 3 2021. ISSN 2522-5839. doi: 10.1038/s42256-021-00302-5. URL <https://doi.org/10.1038/s42256-021-00302-5>.
- H. N. Mhaskar and N. Hahm. Neural networks for functional approximation and system identification. *Neural computation*, 9(1):143–159, 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.1.143.
- Hrushikesh N Mhaskar. Neural networks for optimal approximation of smooth and analytic functions. *Neural computation*, 8(1):164–177, 1996.
- Jeffrey S Morris and Raymond J Carroll. Wavelet-based functional mixed models. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 68(2):179–199, 2006.
- Gustavo A. Muñoz, Yannis Sarantopoulos, and Andrew Tonge. Complexifications of real Banach spaces, polynomials and multilinear maps. *Studia Math.*, 134(1):1–33, 1999. ISSN 0039-3223.
- Nicholas H Nelsen and Andrew M Stuart. Operator learning using random features: A tool for scientific computing. *SIAM Review*, 66(3):535–571, 2024.
- R. Nickl, S. van de Geer, and S. Wang. Convergence rates for penalised least squares estimators in PDE-constrained regression problems. *SIAM J. Uncert. Quant.*, 8, 2020.
- Richard Nickl. Donsker-type theorems for nonparametric maximum likelihood estimators. *Probability Theory and Related Fields*, 138(3-4), 2007. ISSN 0178-8051. doi: 10.1007/s00440-006-0031-4.
- Richard Nickl. Bernstein–von mises theorems for statistical inverse problems i: Schrödinger equation. *Journal of the European Mathematical Society*, 22(8):2697–2750, 2020. ISSN 1435-9855. doi: 10.4171/JEMS/975.
- Richard Nickl. *Bayesian Non-linear Statistical Inverse Problems*. EMS press, 2023.
- Richard Nickl and Sven Wang. On polynomial-time computation of high-dimensional posterior measures by Langevin-type algorithms. *Journal of the European Mathematical Society*, 2024.
- F. Nobile, R. Tempone, and C. G. Webster. A sparse grid stochastic collocation method for partial differential equations with random input data. *SIAM Journal on Numerical Analysis*, 46(5):2309–2345, 2008. doi: 10.1137/060663660. URL <https://doi.org/10.1137/060663660>.
- Thomas O’Leary-Roseberry, Peng Chen, Umberto Villa, and Omar Ghattas. Derivative-informed neural operator: An efficient framework for high-dimensional parametric derivative learning. *Journal of Computational Physics*, 496:112555, 2024. ISSN 0021-9991. doi: <https://doi.org/10.1016/j.jcp.2023.112555>. URL <https://www.sciencedirect.com/science/article/pii/S0021999123006502>.

- J. A. A. Opschoor, Ch. Schwab, and J. Zech. Exponential relu dnn expression of holomorphic maps in high dimension. *Constructive Approximation*, 55(1):537–582, 2022a. ISSN 1432-0940. doi: 10.1007/s00365-021-09542-5. URL <https://link.springer.com/article/10.1007/s00365-021-09542-5#citeas>.
- Joost A. A. Opschoor, Philipp C. Petersen, and Christoph Schwab. Deep relu networks and high-order finite element methods. *Analysis and Applications*, 18(05):715–770, 2020. ISSN 0219-5305. doi: gjh43n.
- Joost A. A. Opschoor, Christoph Schwab, and Jakob Zech. Deep learning in high dimension: ReLU neural network expression for Bayesian PDE inversion. In *Optimization and control for partial differential equations—uncertainty quantification, open and closed-loop control, and shape optimization*, volume 29 of *Radon Ser. Comput. Appl. Math.*, pages 419–462. De Gruyter, Berlin, 2022b. ISBN 978-3-11-069596-0. doi: 10.1515/9783110695984-015. URL <https://doi.org/10.1515/9783110695984-015>.
- Houman Owhadi and Gene Ryan Yoo. Kernel flows: From learning kernels from data into the abyss. *Journal of Computational Physics*, 389:22–47, 2019.
- Thomas O’Leary-Roseberry, Umberto Villa, Peng Chen, and Omar Ghattas. Derivative-informed projected neural networks for high-dimensional parametric maps governed by pdes. *Computer Methods in Applied Mechanics and Engineering*, 388:114199, 2022. ISSN 0045-7825. doi: <https://doi.org/10.1016/j.cma.2021.114199>. URL <https://www.sciencedirect.com/science/article/pii/S0045782521005302>.
- Philipp Petersen and Felix Voigtlaender. Optimal approximation of piecewise smooth functions using deep relu neural networks. *Neural networks : the official journal of the International Neural Network Society*, 108:296–330, 2018. doi: 10.1016/j.neunet.2018.08.019.
- Philipp Petersen, Mones Raslan, and Felix Voigtlaender. Topological properties of the set of functions generated by neural networks of fixed size. *Foundations of Computational Mathematics*, 21(2):375–444, 2021. ISSN 1615-3375. doi: 10.1007/s10208-020-09461-0. URL <https://link.springer.com/article/10.1007/s10208-020-09461-0>.
- Allan Pinkus. Approximation theory of the MLP model in neural networks. In *Acta numerica, 1999*, volume 8 of *Acta Numer.*, pages 143–195. Cambridge Univ. Press, Cambridge, 1999. ISBN 0-521-77088-2. doi: 10.1017/S0962492900002919. URL <https://doi.org/10.1017/S0962492900002919>.
- Tomaso Poggio, Hrushikesh Mhaskar, Lorenzo Rosasco, Brando Miranda, and Qianli Liao. Why and when can deep-but not shallow-networks avoid the curse of dimensionality: a review. *International Journal of Automation and Computing*, 14(5):503–519, 2017.
- David Pollard. *Convergence of Stochastic Processes*. Springer New York, New York, NY, 1984. ISBN 978-1-4612-9758-1. doi: 10.1007/978-1-4612-5254-2.
- Alfio Quarteroni, Andrea Manzoni, and Federico Negri. *Reduced basis methods for partial differential equations*, volume 92 of *Unitext*. Springer, Cham, 2016. ISBN 978-3-319-15430-5. doi: 10.1007/978-3-319-15431-2. URL <https://doi.org/10.1007/978-3-319-15431-2>. An introduction, *La Matematica per il 3+2*.

- Bogdan Raonic, Roberto Molinaro, Tobias Rohner, Siddhartha Mishra, and Emmanuel de Bezenac. Convolutional neural operators. In *ICLR 2023 Workshop on Physics for Machine Learning*, 2023.
- Holger Rauhut and Christoph Schwab. Compressive sensing Petrov-Galerkin approximation of high-dimensional parametric operator equations. *Math. Comp.*, 86(304):661–700, 2017. ISSN 0025-5718. doi: 10.1090/mcom/3113. URL <http://dx.doi.org/10.1090/mcom/3113>. Report 2014-14, Seminar for Applied Mathematics, ETH Zürich.
- Markus Reiß. Asymptotic equivalence for nonparametric regression with multivariate and random design. *Ann. Statist.*, 36(4):1957–1982, 2008. ISSN 0090-5364. doi: 10.1214/07-AOS525. URL <http://dx.doi.org/10.1214/07-AOS525>.
- Johannes Schmidt-Hieber. Supplement to “nonparametric regression using deep neural networks with relu activation function”. *The Annals of Statistics*, 48(4), 2020a. ISSN 0090-5364. doi: 10.1214/19-AOS1875SUPP.
- Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with ReLU activation function. *The Annals of Statistics*, 48(4), 2020b. doi: 10.1214/19-aos1875. URL <https://doi.org/10.1214/19-aos1875>.
- C. Schwab and A. M. Stuart. Sparse deterministic approximation of Bayesian inverse problems. *Inverse Problems*, 28(4):045003, 32, 2012. ISSN 0266-5611,1361-6420. doi: 10.1088/0266-5611/28/4/045003. URL <https://doi.org/10.1088/0266-5611/28/4/045003>.
- Christoph Schwab and Jakob Zech. Deep learning in high dimension: neural network expression rates for generalized polynomial chaos expansions in UQ. *Anal. Appl. (Singap.)*, 17(1):19–55, 2019. ISSN 0219-5305,1793-6861. doi: 10.1142/S0219530518500203. URL <https://doi.org/10.1142/S0219530518500203>.
- Christoph Schwab and Jakob Zech. Deep learning in high dimension: neural network expression rates for analytic functions in  $L^2(\mathbb{R}^d, \gamma_d)$ . *SIAM/ASA J. Uncertain. Quantif.*, 11(1):199–234, 2023. ISSN 2166-2525. doi: 10.1137/21M1462738. URL <https://doi.org/10.1137/21M1462738>.
- Christoph Schwab, Andreas Stein, and Jakob Zech. Deep operator network approximation rates for lipschitz operators. *Analysis and Applications*, pages 1–41, 2025. doi: 10.1142/S0219530525500307. URL <https://doi.org/10.1142/S0219530525500307>.
- Euan A Spence and Jared Wunsch. Wavenumber-explicit parametric holomorphy of helmholtz solutions in the context of uncertainty quantification. *SIAM/ASA Journal on Uncertainty Quantification*, 11(2):567–590, 2023.
- Andrew M. Stuart. Inverse problems: a Bayesian perspective. *Acta Numer.*, 19:451–559, 2010. ISSN 0962-4929. doi: 10.1017/S0962492910000061.
- Taiji Suzuki. Adaptivity of deep relu network for learning in besov and mixed smooth besov spaces: optimal rate and curse of dimensionality. In *he 7th International Conference on Learning Representations (ICLR2019)*, volume 7, 2019.

- Michel Talagrand. *The generic chaining: Upper and lower bounds for stochastic processes*. Springer, Berlin and New York, 2005. ISBN 3-540-24518-9. doi: 10.1007/3-540-27499-5.
- Michel Talagrand. *Upper and lower bounds for stochastic processes*, volume 60. Springer, 2014.
- Hans Triebel. *Interpolation theory, function spaces, differential operators*. Barth, Heidelberg and Leipzig, 2., rev. and enl. ed. edition, 1995. ISBN 3335004205.
- Alexandre B Tsybakov. *Introduction to Nonparametric Estimation*. Springer series in statistics. Springer, Dordrecht, 2009. doi: 10.1007/b13794. URL <https://cds.cern.ch/record/1315296>.
- Sara van de Geer. *Empirical Processes in M-Estimation*. Cambridge U. Press, 2000.
- Sara van de Geer. Least squares estimation with complexity penalties. *Mathematical Methods of statistics*, 10(3):355, 2001.
- Sara A. van de Geer. *Applications of empirical process theory*. Cambridge series in statistical and probabilistic mathematics. Cambridge University Press, Cambridge, digitally printed version edition, 2000. ISBN 052165002X.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47 of *Cambridge series in statistical and probabilistic mathematics*. Cambridge University Press, Cambridge, United Kingdom and New York, NY, 2018. ISBN 1108415199.
- Dongbin Xiu and George Em Karniadakis. The wiener–askey polynomial chaos for stochastic differential equations. *SIAM Journal on Scientific Computing*, 24(2):619–644, 2002. doi: 10.1137/S1064827501387826. URL <https://doi.org/10.1137/S1064827501387826>.
- Dmitry Yarotsky. Error bounds for approximations with deep relu networks. *Neural networks : the official journal of the International Neural Network Society*, 94:103–114, 2017. doi: 10.1016/j.neunet.2017.07.002.
- Jakob Zech. *Sparse-Grid Approximation of High-Dimensional Parametric PDEs*. PhD thesis, ETH Zurich, 2018.