
Deep learning is adaptive to intrinsic dimensionality of model smoothness in anisotropic Besov space

Taiji Suzuki

Department of Mathematical Informatics, The University of Tokyo, Tokyo, Japan
RIKEN Center for Advanced Intelligence Project, Tokyo, Japan
taiji@mist.i.u-tokyo.ac.jp

Atsushi Nitanda

Kyushu Institute of Technology, Fukuoka, Japan
RIKEN Center for Advanced Intelligence Project, Tokyo, Japan
nitanda@ai.kyutech.ac.jp

Abstract

Deep learning has exhibited superior performance for various tasks, especially for high-dimensional datasets, such as images. To understand this property, we investigate the approximation and estimation ability of deep learning on *anisotropic Besov spaces*. The anisotropic Besov space is characterized by direction-dependent smoothness and includes several function classes that have been investigated thus far. We demonstrate that the approximation error and estimation error of deep learning only depend on the average value of the smoothness parameters in all directions. Consequently, the curse of dimensionality can be avoided if the smoothness of the target function is highly anisotropic. Unlike existing studies, our analysis does not require a low-dimensional structure of the input data. We also investigate the minimax optimality of deep learning and compare its performance with that of the kernel method (more generally, linear estimators). The results show that deep learning has better dependence on the input dimensionality if the target function possesses anisotropic smoothness, and it achieves an adaptive rate for functions with spatially inhomogeneous smoothness.

1 Introduction

Based on the recent literature pertaining to machine learning, deep learning has exhibited superior performance in several tasks such as image recognition (Krizhevsky et al., 2012), natural language processing (Devlin et al., 2018), and image synthesis (Radford et al., 2015). In particular, its superiority is remarkable for complicated and high-dimensional data like images. This is mainly due to its high flexibility and superior feature-extraction ability for effectively extracting the intrinsic structure of data. Its theoretical analysis also has been extensively developed considering several aspects such as expressive ability, optimization, and generalization error.

Amongst representation ability analysis of deep neural networks such as universal approximation ability (Cybenko, 1989; Hornik, 1991; Sonoda & Murata, 2017), approximation theory of deep neural networks on typical function classes such as Hölder, Sobolev, and Besov spaces have been extensively studied. In particular, analyses of deep neural networks with the ReLU activation (Nair & Hinton, 2010; Glorot et al., 2011) have been recently developed. Schmidt-Hieber (2020) showed that the deep learning with ReLU activations can achieve the minimax optimal estimation accuracy to estimate composite functions in Hölder spaces by using the approximation theory of Yarotsky (2017). Suzuki (2019) generalized this analysis to those on the *Besov space* and the *mixed smooth*

Table 1: Relationship between existing research and our work. α indicates the smoothness of the target function, d is the dimensionality of input x , D is the dimensionality of a low-dimensional structure on which the data are distributed, and $\tilde{\alpha}$ is the average smoothness of an anisotropic Besov space (Eq. (1)).

Function class	Hölder	Besov	mixed smooth Besov	Hölder on a low-dimensional set	anisotropic Besov
Author	Schmidt-Hieber (2020)	Suzuki (2019)	Suzuki (2019)	Nakada & Imaizumi (2020); Schmidt-Hieber (2019); Chen et al. (2019)	This work
Estimation error	$\mathcal{O}(n^{-\frac{2}{2+d}})$	$\mathcal{O}(n^{-\frac{2}{2+d}})$	$\mathcal{O}\left(n^{-\frac{2}{2+1}} \log(n)^{\frac{2(d-1)(\alpha+)}{1+2}}\right)$	$\mathcal{O}(n^{-\frac{2}{2+D}})$	$\mathcal{O}(n^{-\frac{2}{2+1}})$

Besov space by utilizing the techniques developed in approximation theories (Temlyakov, 1993; DeVore, 1998). It was shown that deep learning can achieve an *adaptive approximation* error rate that is faster than that of (non-adaptive) linear approximation methods (DeVore & Popov, 1988; DeVore et al., 1993; Düng, 2011b), and it outperforms any linear estimators (including kernel ridge regression) in terms of the minimax optimal rate.

From these analyses, one can see that the approximation errors and estimation errors are strongly influenced by two factors, i.e., the *smoothness* of the target function and the *dimensionality* of the input (see Table 1). In particular, they suffer from the *curse of dimensionality*, which is unavoidable. However, these analyses are about the worst case errors and do not exploit specific intrinsic properties of the true distributions. For example, practically encountered data usually possess low intrinsic dimensionality, i.e., data are distributed on a low dimensional sub-manifold of the input space (Tenenbaum et al., 2000; Belkin & Niyogi, 2003). Recently, Nakada & Imaizumi (2020); Schmidt-Hieber (2019); Chen et al. (2019); Chen et al. (2019) have shown that deep ReLU network has adaptivity to the intrinsic dimensionality of data and can avoid curse of dimensionality if the intrinsic dimensionality is small. However, one drawback is that they assumed *exact* low dimensionality of the input data. This could be a strong assumption because practically observed data are always noisy, and injecting noise immediately destroys the low-dimensional structure. Therefore, we consider another direction in this paper. In terms of curse of dimensionality, Suzuki (2019) showed that deep learning can alleviate the curse of dimensionality to estimate functions in a so called mixed smooth Besov space (m-Besov). However, m-Besov space assumes strong smoothness toward *all* directions uniformly and does not include the ordinary Besov space as a special case. Moreover, the convergence rate includes heavy poly-log term which is not negligible (see Table 1).

In practice, one of the typically expected properties of a true function on high-dimensional data is that it is invariant against perturbations of an input in some specific directions (Figure 1). For example, in image-recognition tasks, the target function must be invariant against the spatial shift of an input image, which is utilized by data-augmentation techniques (Simard et al., 2003; Krizhevsky et al., 2012). In this paper, we investigate the approximation and estimation abilities of deep learning on *anisotropic Besov spaces* (Nikol’skii, 1975; Vybiral, 2006; Triebel, 2011) (also called dominated mixed-smooth Besov spaces). An anisotropic Besov space is a set of functions that have “direction-dependent” smoothness, whereas ordinary function spaces such as Hölder, Sobolev, and Besov spaces assume isotropic smoothness that is uniform in all directions. We consider a composition of functions included in an anisotropic Besov space, including several existing settings as special cases; it includes analyses of the Hölder space Schmidt-Hieber (2020) and Besov space Suzuki (2019), as well as the low-dimensional sub-manifold setting (Nakada & Imaizumi, 2020; Schmidt-Hieber, 2019; Chen et al., 2019; Chen et al., 2019)¹. By considering such a space, we can show that deep learning can alleviate curse of dimensionality if the smoothness in each direction is highly anisotropic. Interestingly, any linear estimator (including kernel ridge regression) has worse dependence on the dimensionality than deep learning. Our contributions can be summarized as follows:

- We consider a situation in which the target function is included in a class of anisotropic Besov spaces and show that deep learning can avoid the curse of dimensionality *even if the input data*

¹We would like to remark that the analysis of Nakada & Imaizumi (2020) does not require smoothness of the embedded manifold that is not covered in this paper.

do not lie on a low-dimensional manifold. Moreover, deep learning can achieve the optimal adaptive approximation error rate and minimax optimal estimation error rate.

- We compare deep learning with general linear estimators (including kernel methods) and show that deep learning has better dependence on the input dimensionality than linear estimators.

2 Problem setting and the model

In this section, we describe the problem setting considered in this work. We consider the following nonparametric regression model:

$$y_i = f^0(x_i) + \varepsilon_i \quad (i = 1; \dots; n);$$

where x_i is generated from a probability distribution P_X on $[0; 1]^d$, $\varepsilon_i \sim N(0; \sigma^2)$, and the data $D_n = (x_i; y_i)_{i=1}^n$ are independently identically distributed. f^0 is the true function that we want to estimate. We are interested in the mean squared estimation error of an estimator \hat{f} : $E_{D_n}[k\hat{f} - f^0]_{L^2(P_X)}^2$; where $E_{D_n}[\cdot]$ indicates the expectation with respect to the training data D_n .

We consider a least-squares estimator in the deep neural network model as \hat{f} (see Eq. (5)) and discuss its optimality. More specifically, we investigate how the ‘‘intrinsic dimensionality’’ of data affects the estimation accuracy of deep learning. For this purpose, we consider an *anisotropic Besov space* as a model of the target function.

2.1 Anisotropic Besov space

In this section, we introduce the anisotropic Besov which was investigated as the model of the true function in this paper. Throughout this paper, we set the domain of the input to $\mathcal{X} = [0; 1]^d$. For a function $f : \mathcal{X} \rightarrow \mathbb{R}$, let $kfk_p := kfk_{L^p(\mathcal{X})} := (\int_{\mathcal{X}} |f|^p dx)^{1/p}$ for $0 < p < \infty$. For $p = \infty$, we define $kfk_{\infty} := kfk_{L^{\infty}(\mathcal{X})} := \sup_{x \in \mathcal{X}} |f(x)|$. For $\alpha \in \mathbb{R}_+$, let $j_{\alpha} = \sum_{j=1}^d j_{\alpha}^2$.

For a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we define the r th difference of f in the direction $h \in \mathbb{R}^d$ as

$$\tau_h^r(f)(x) := \tau_h^{r-1}(f)(x+h) - \tau_h^{r-1}(f)(x); \quad \tau_h^0(f)(x) := f(x);$$

for $x \in \mathcal{X}$ with $x+rh \in \mathcal{X}$, otherwise, let $\tau_h^r(f)(x) = 0$.

Definition 1. For a function $f \in L^p(\mathcal{X})$ where $p \in (0; \infty]$, the r -th modulus of smoothness of f is defined by $W_{r,p}(f; t) = \sup_{h \in \mathbb{R}^d; |h|_1 \leq t} k\tau_h^r(f)k_p$; for $t = (t_1; \dots; t_d)$; $t_i > 0$.

With this modulus of smoothness, we define the anisotropic Besov space $B_{p,q}(\mathcal{X})$ for $\mathcal{X} = ([0; 1; \dots; 1]^d) \subset \mathbb{R}_+^d$ as follows.

Definition 2 (Anisotropic Besov space ($B_{p,q}(\mathcal{X})$)). For $0 < p, q \leq \infty$, $\mathcal{X} = ([0; 1; \dots; 1]^d) \subset \mathbb{R}_+^d$, $r := \max_i b_i \in \mathbb{C} + 1$, let the seminorm $j_{B_{p,q}}$ be

$$j_{B_{p,q}} := \begin{cases} \left(\sum_{k=0}^{\infty} [2^k W_{r,p}(f; (2^{-k}; \dots; 2^{-k}))]^q \right)^{1/q} & (q < \infty); \\ \sup_k 2^k W_{r,p}(f; (2^{-k}; \dots; 2^{-k})) & (q = \infty); \end{cases}$$

The norm of the anisotropic Besov space $B_{p,q}(\mathcal{X})$ is defined by $kfk_{B_{p,q}} := kfk_p + j_{B_{p,q}}$, and $B_{p,q}(\mathcal{X}) = \{f \in L^p(\mathcal{X}) \mid kfk_{B_{p,q}} < \infty\}$.

Roughly speaking $j_{B_{p,q}}$ represents the smoothness in each direction. If b_i is large, then a function in $B_{p,q}$ is smooth to the i th coordinate direction, otherwise, it is non-smooth to that direction. p is also an important quantity that controls the *spatial inhomogeneity* of the smoothness. If $b_1 = b_2 = \dots = b_d$, then the definition is equivalent to the usual Besov space (DeVore & Popov, 1988; DeVore et al., 1993). Suzuki (2019) analyzed curse of dimensionality of deep learning through a so-called *mixed smooth Besov* (m-Besov) space which imposes a stronger condition toward all directions uniformly.

²We let $\mathbb{N} := \{1; 2; 3; \dots\}$, $\mathbb{Z}_+ := \{0; 1; 2; 3; \dots\}$, $\mathbb{Z}_+^d := \{(z_1; \dots; z_d) \mid z_i \in \mathbb{Z}_+\}$, $\mathbb{R}_+ := \{x \in \mathbb{R} \mid x > 0\}$, and $\mathbb{R}_{++} := \{x \in \mathbb{R} \mid x > 0\}$. We let $[N] := \{1; \dots; N\}$ for $N \in \mathbb{N}$.

Particularly, it imposes stronger smoothness toward non-coordinate axis directions. Moreover, m -Besov space does *not* include the vanilla Besov space as a special case and thus cannot capture the situation that we consider in this paper.

Throughout this paper, for given $\alpha = (\alpha_1; \dots; \alpha_d) \in \mathbb{R}_{++}^d$, we write $\alpha_- := \min_i \alpha_i$ (smallest smoothness) and $\alpha^+ := \max_i \alpha_i$ (largest smoothness). The approximation error of a function in anisotropic Besov spaces is characterized by the harmonic mean of $(\alpha_j)_{j=1}^d$, which corresponds to the average smoothness, and thus we define

$$\tilde{\alpha} := \left(\sum_{j=1}^d \alpha_j^{-1} \right)^{-1}. \quad (1)$$

The Besov space is closely related to other function spaces such as Hölder space. Let $\mathcal{C}^{\alpha}(x) = \frac{\partial^{\alpha} f}{\partial x_1^{\alpha_1} \dots \partial x_d^{\alpha_d}}(x)$:

Definition 3 (Hölder space $\mathcal{C}^{\alpha}(\cdot)$). For a smoothness parameter $\alpha \in \mathbb{R}_{++}^d$ with $\alpha \in \mathbb{N}$, consider an m times differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ where $m = \lfloor \alpha^+ \rfloor$ (the largest integer less than α^+), and let the norm of the Hölder space $\mathcal{C}^{\alpha}(\cdot)$ be $\|f\|_{\mathcal{C}^{\alpha}} := \max_{|j| \leq m} \|f^{(j)}\|_{\infty} + \max_{|j|=m} \sup_{x, y \in \mathbb{R}^d} \frac{|f^{(j)}(x) - f^{(j)}(y)|}{\|x - y\|_{\alpha}^m}$. Then, (α) -Hölder space $\mathcal{C}^{\alpha}(\cdot)$ is defined as $\mathcal{C}^{\alpha}(\cdot) = \{f \mid \|f\|_{\mathcal{C}^{\alpha}} < \infty\}$.

Let $\mathcal{C}^0(\cdot)$ be the set of continuous functions equipped with L^{∞} -norm: $\mathcal{C}^0(\cdot) := \{f \mid f \text{ is continuous and } \|f\|_{\infty} < \infty\}$. These function spaces are closely related to each other.

Proposition 1 (Triebel (2011)). There exist the following relations between the spaces:

1. For $\alpha = (\alpha_0; \dots; \alpha_0) \in \mathbb{R}_{++}^d$ with $\alpha_0 \in \mathbb{N}$, it holds that $\mathcal{C}^{\alpha}(\cdot) = B_{1,1}^{\alpha}(\cdot)$:
2. For $0 < p_1, p_2, q \leq \infty$, $p_1 \leq p_2$ and $\alpha \in \mathbb{R}_{++}^d$ with $\tilde{\alpha} > (1/p_1 - 1/p_2)_+$ ³, it holds that⁴ $B_{p_1, q}^{\alpha}(\cdot) \hookrightarrow B_{p_2, q}^{\alpha}(\cdot)$ for $\alpha = 1 - (1/p_1 - 1/p_2)_+ = \tilde{\alpha}$.
3. For $0 < p, q_1, q_2 \leq \infty$, $q_1 < q_2$, and $\alpha \in \mathbb{R}_{++}^d$, it holds that $B_{p, q_1}^{\alpha} \hookrightarrow B_{p, q_2}^{\alpha}$. In particular, with properties 1 and 2, if $\tilde{\alpha} > 1/p$, it holds that $B_{p, q}^{\alpha}(\cdot) \hookrightarrow \mathcal{C}^{\alpha}(\cdot)$ where $\alpha = 1 - 1/p = \tilde{\alpha}$.
4. For $0 < p, q \leq \infty$ and $\alpha \in \mathbb{R}_{++}^d$, if $\tilde{\alpha} > 1/p$, then $B_{p, q}^{\alpha}(\cdot) \hookrightarrow \mathcal{C}^0(\cdot)$.

This result is basically proven by Triebel (2011). For completeness, we provide its derivation in Appendix D. If the average smoothness $\tilde{\alpha}$ is sufficiently large ($\tilde{\alpha} > 1/p$), then the functions in $B_{p, q}^{\alpha}$ are continuous; however, if it is small ($\tilde{\alpha} < 1/p$), then they are no longer continuous. Small p indicates spatially inhomogeneous smoothness; thus, spikes and jumps appear (see Donoho & Johnstone (1998) for this perspective, from the viewpoint of wavelet analysis).

2.2 Model of the true function

As a model of the true function f^0 , we consider two types of models: *Affine composition model* and *deep composition model*. For a Banach space H , we let $U(H)$ be the unit ball of H .

(a) Affine composition model: The first model we introduced is a very naive model which is just a composition of an affine transformation and a function in the anisotropic Besov space:

$$H_a := \{f \circ h \mid h \in U(B_{p, q}^{\alpha}([0, 1]^d)); A \in \mathbb{R}^{d \times d}; b \in \mathbb{R}^d \text{ s.t. } Ax + b \in [0, 1]^d (\forall x \in \mathbb{R}^d)\};$$

where we assume $d \leq d$. Here, we assumed that the affine transformation has an appropriate scaling such that $Ax + b$ is included in the domain of h for all $x \in \mathbb{R}^d$. This is a quite naive model but provides an instructive example to understand how the estimation error of deep learning behaves under the anisotropic setting.

(b) Deep composition model: The *deep composition model* generalizes the affine composition model to a composition of nonlinear functions. Let $m_1 = d, m_{L+1} = 1, m_{\cdot}$ be the dimension of the

³Here, we let $(x)_+ := \max\{x, 0\}$.

⁴The symbol \hookrightarrow means continuous embedding.

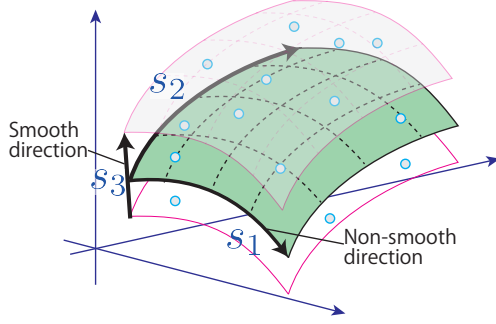


Figure 1: Near low dimensional data distribution with anisotropic smoothness of the target function. The target function has less smoothness (S_1, S_2) toward the first two coordinates on the manifold while it is almost constant toward the third coordinate (large S_3).

l th layer, and let $(\cdot) \in \mathbb{R}_{++}^m$ be the smoothness parameter in the l th layer. The deep composition model is defined as

$$H_{\text{deep}} := f \circ h_H \quad h_1(x) \in [0; 1]^{m_1}, \dots, [0; 1]^{m_{l+1}}; h_{l;k} \in U(B_{p,q}^{(\cdot)}([0; 1]^{m_l})) \quad (8k \in [m_{l+1}])g$$

Here, the interval $[0; 1]$ can be replaced by another compact interval, such as $[a; b]$, but this difference can be absorbed by changing a scaling factor. The assumption $kh_{l;k} \in B_{p,q}^{(\cdot)}([0; 1]^{m_l})$ can also be relaxed, but we do not pursue that direction due to presentation simplicity. This model includes the affine composition model as a special case. However, it requires a stronger assumption to properly evaluate the estimation error on this model.

Examples The model we have introduced includes some instructive examples as listed below:

(a) **Linear projection** Schmidt-Hieber (2020) analyzed estimation of the following model by deep learning: $f^0(x) = g(w^T x)$ where $g \in \mathcal{C}([0; 1])$ and $w \in \mathbb{R}^d$. In this example, the function f^0 varies along only one direction, w . Apparently, this is an example of the affine composition model.

(b) **Distribution on low dimensional smooth manifold** Assume that the input x is distributed on a low-dimensional smooth manifold embedded in \mathbb{R}^d , and the smoothness of the true function f^0 is anisotropic along a coordinate direction on the manifold. We suppose that the low dimensional manifold is d -dimensional and $d \ll d$. In this situation, the true function can be written as $f^0(x) = h(\pi(x))$ where $\pi: \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a map that returns the coordinate of x on the manifold and h is an element in an anisotropic Besov space on \mathbb{R}^d . This situation appears if data is distributed on a low-dimensional sub-manifold of \mathbb{R}^d and the target function is invariant against noise injection to some direction on the manifold at each input point x (Figure 1 illustrates this situation). One typical example of this situation is a function invariant with data augmentation (Simard et al., 2003; Krizhevsky et al., 2012). Even if the noise injection destroys low dimensionality of the data distribution (i.e., $d = d$), an anisotropic smoothness of the target function eases the curse of dimensionality as analyzed below, which is quite different from existing works (Yang & Dunson, 2016; Bickel & Li, 2007; Nakada & Imaizumi, 2020; Schmidt-Hieber, 2019; Chen et al., 2019; Chen et al., 2019).

Related work Here, we introduce some more related work and discuss their relation to our analysis. The statistical analysis on an anisotropic Besov space can be back to Ibragimov & Khas'minskii (1984) who considered density estimation, where the density is assumed to be included in an anisotropic Sobolev space with $p \geq 2$, and derived the minimax optimal rate $n^{-r/(2r+1)}$ with respect to L^r -norm. Nyssbaum (1983, 1987) analyzed a nonparametric regression problem on an anisotropic Besov space. Following these results, several studies have been conducted in the literature pertaining to nonparametric statistics, such as nonlinear kernel estimator Kerkycharian et al. (2001), adaptive confidence band construction Hoffman & Lepski (2002), optimal aggregation Gaifas & Lecue (2011), Gaussian process estimator Bhattacharya et al. (2011, 2014), and kernel ridge regression Hang & Steinwart (2018). Basically, these studies investigated estimation problems in which the target function is in anisotropic Besov spaces, but the composition models considered in this paper have not been analyzed. Hoffman & Lepski (2002); Bhattacharya et al. (2011) considered a dimension reduction model; that is, the target function is dependent on only a few variables of x , but they did not deal with more general models, such as the affine/deep composition models. The nonparametric regression problems where the input data are distributed on a low-dimensional smooth manifold has been studied as a ‘‘manifold regression’’ Yang & Dunson (2016); Bickel & Li

(2007); Yang & Tokdar (2015). Such a model can be considered as a specific example of the deep composition model. In this sense, our analysis is a significant extension of these analyses.

3 Approximation error analysis

Here, we consider approximating the true function f^0 via a deep neural network and derive the approximation error. As the activation function, we consider the ReLU activation denoted by $\sigma(x) = \max\{x, 0\}g(x \in \mathbb{R})$. For a vector x , $\sigma(x)$ is operated in an element-wise manner. The model of neural networks with height L , width W , sparsity constraint S , and norm constraint B as $(L; W; S; B) := f(W^{(L)}(\cdot) + b^{(L)}) \circ \dots \circ (W^{(1)}x + b^{(1)})$, $W^{(l)} \in \mathbb{R}^{W \times W}$, $b^{(l)} \in \mathbb{R}^W$, $W^{(1)} \in \mathbb{R}^{W \times d}$, $b^{(1)} \in \mathbb{R}^W$; $W^{(l)} \in \mathbb{R}^{W \times W}$; $b^{(l)} \in \mathbb{R}^W$ ($1 \leq l \leq L$); $\sum_{l=1}^L (kW^{(l)}k_0 + kb^{(l)}k_0) \leq S$; $\max_l kW^{(l)}k_1 \leq B$; where k_0 is the ℓ_0 -norm of the matrix (the number of non-zero elements of the matrix), and k_1 is the ℓ_1 -norm of the matrix (maximum of the absolute values of the elements). The sparsity constraint and norm bounds are required to obtain the near-optimal rate of the estimation error. To evaluate the accuracy of the deep neural network model in approximating target functions, we define the worst-case approximation error as

$$R_r(F; H) := \sup_{f \in F} \inf_{f \in \mathcal{H}} \|f - f\|_{L^r(\cdot)};$$

where F is the set of functions used for approximation, and H is the set of target functions.

Proposition 2 (Approximation ability for anisotropic Besov space). *Suppose that $0 < p, q, r \leq 1$ and $\tilde{\gamma} \in \mathbb{R}_+^d$ satisfy the following condition: $\tilde{\gamma} > (1-p, 1-r)_+$. Assume that $m \in \mathbb{N}$ satisfies $0 < \tilde{\gamma} < \min(m, m-1+1-p)$. Let $\tilde{\gamma} = (1-p, 1-r)_+$, $\tilde{\gamma} = (\tilde{\gamma})_+ = (2, \dots)$ and $W_0(d) := 6dm(m+2) + 2d$. Then, for $N \in \mathbb{N}$, we can bound the approximation error as*

$$R_r((L_1; W_1; S_1; B_1); U(B_{p,q}(\tilde{\gamma}))) \leq N^{-\tilde{\gamma}};$$

by setting

$$L_1(d) := 3 + 2d \log_2 \left(\frac{3^{d-m}}{c_{(d,m)}} \right) + 5ed \log_2(d-m)e; \quad W_1(d) := NW_0; \quad (2)$$

$$S_1(d) := [(L-1)W_0^2 + 1]N; \quad B_1(d) := O(N^{d(1-\tilde{\gamma})(1-p-\tilde{\gamma})_+}); \quad (3)$$

for $\tilde{\gamma} = N^{-\tilde{\gamma}} \log(N)^{-1}$ and a constant $c_{(d,m)}$ depending only on d and m .

The proof of this proposition is provided in Appendix B. The rate $N^{-\tilde{\gamma}}$ is the optimal adaptive approximation error rate that can be achieved by a model with N parameters (the difference between adaptive and non-adaptive methods is explained in the discussion below). Note that this is an approximation error in an oracle setting and no sample complexity appears here. We notice that we can avoid the *curse of dimensionality* if the average smoothness $\tilde{\gamma}$ is small. This means that if the target function is non-smooth in only a few directions and smooth in other directions, we can avoid the curse of dimensionality. In contrast, if we consider an isotropic Besov space where $\tilde{\gamma}_1 = \tilde{\gamma} = d/2$, then $\tilde{\gamma} = d/2$, which directly depends on the dimensionality d , and we need an exponentially large number of parameters in this situation to achieve ϵ -accuracy. Therefore, the anisotropic smoothness has a significant impact on the approximation error rate. The assumption $\tilde{\gamma} > (1-p, 1-r)_+$ ensures the L_r -integrability of the target function, and the inequality (without equality) admits a near-optimal wavelet approximation of the target function in terms of L_r -norm.

Using this evaluation as a basic tool, we can obtain the approximation error for the deep composition models. We can also obtain the approximation error for the affine composition models, but it is almost identical to Proposition 2. Therefore, we defer it to Appendix A.

Theorem 1 (Deep composition model). *Assume that $\tilde{\gamma}^{(l)} > 1-p$ for all $l = 1, \dots, H$. Then, the estimation error on the deep composition model is bounded as*

$$R_1((L; W; S; B); H_{\text{deep}}) \leq \max_{l \in [H]} N^{-\tilde{\gamma}^{(l)}}; \quad (4)$$

where $\tilde{\gamma}^{(l)} = \tilde{\gamma}^{(l)} \prod_{k=1}^H [(L_k - 1-p)^{\wedge 1}]$, and $L = \sum_{l=1}^H (L_1(m_l) + 1)$; $W = \max_l (W_1(m_l) - m_{l+1})$; $S = \sum_{l=1}^H (S_1(m_l) + 3m_{l+1})$; $B = \max_l B_1(m_l)$.

The proof can be found in Appendix B.1. Since the model is more general than the vanilla anisotropic Besov space, we require a stronger assumption $\tilde{(\cdot)} > 1=\rho$ on $\tilde{(\cdot)}$ than the condition in Proposition 2. This is because we need to bound the Hölder smoothness of the remaining layers to bound the influence of the approximation error in the internal layers to the entire function. Hölder smoothness is ensured according to the embedding property under this condition (Proposition 1). This Hölder smoothness assumption affects the approximation error rate. The convergence rate $\tilde{(\cdot)}$ in Eq. (4) is different from that in Eq. (8). This is because the approximation error in the internal layers are propagated through the remaining layers with Hölder smoothness and its amplitude is controlled by the Hölder smoothness.

Approximation error by non-adaptive method The approximation error obtained in the previous section is called an adaptive error rate in the literature regarding approximation theory (DeVore, 1998). If we fix N bases beforehand and approximate the target function by a linear combination of the N bases (which is called the non-adaptive method), then we *cannot* achieve the adaptive error rate obtained in the previous section. Roughly speaking, the approximation error of non-adaptive methods is lower bounded by $N^{-\left(\frac{1}{p} - \frac{1}{\min\{2, r\}g}\right)_+}$ (Myronyuk, 2015, 2016, 2017), which is slower than the approximation error rate of deep neural networks especially for small ρ .

4 Estimation error analysis

In this section, we analyze the accuracy of deep learning in estimating a function in compositions of anisotropic Besov spaces. We consider a least-squares estimator in the deep neural network model:

$$\hat{f} = \operatorname{argmin}_{f: f \in \mathcal{L}(L; W; S; B)} \sum_{i=1}^n (y_i - f(x_i))^2 \quad (5)$$

where f is the *clipping* of f defined by $f = \min\{f, F\} \max\{f, -F\}$ for a constant $F > 0$ which is realized by ReLU units. The network parameters $(L; W; S; B)$ should be specified appropriately as indicated in Theorems 2 and 3. In practice, these parameters can be specified by cross validation. Indeed, we can theoretically show that cross validation can provide the appropriate choice of these parameters in compensation to an additional $\log(n)$ -factor in the estimation error bound. This estimator can be seen as a sparsely regularized estimator because there are constraints on S . In terms of optimization, this requires a combinatorial optimization, but we do not pursue the computational aspect. The estimation error that we derive in this section can involve the optimization error, but for simplicity, we only demonstrate the estimation error of the *ideal* situation where there is no optimization error.

Affine composition model The following theorem provides an upper bound of the estimation error for the affine composition model.

Theorem 2. *Assume the same condition as in Theorem 6; in particular, suppose $0 < p; q \leq 1$ and $\tilde{(\cdot)} > (1-\rho - 1/2)_+$ for $\tilde{(\cdot)} \in \mathbb{R}_{++}^d$. Moreover, we assume that the distribution P_X has a density p_X such that $k p_X k_1 \leq R$ for a constant $R > 0$. If $f^0 \in H_a \setminus L^1(\cdot)$, and $k f^0 k_1 \leq F$ for $F \geq 1$; then, letting $(W; L; S; B) = (L_1(d); W_1(d); S_1(d); (dC + 1)B_1(d))$ as in Theorem 6 with $N = n^{\frac{1}{2^{\tilde{(\cdot)}} + 1}}$, we obtain*

$$E_{D_n}[k f^0 - \hat{f} k_{L^2(P_X)}^2] \leq n^{-\frac{2^{\tilde{(\cdot)}}}{2^{\tilde{(\cdot)}} + 1}} \log(n)^3;$$

where $E_{D_n}[\cdot]$ indicates the expectation with respect to the training data D_n .

The proof is provided in Appendix C. We will show that the convergence rate $n^{-\frac{2^{\tilde{(\cdot)}}}{2^{\tilde{(\cdot)}} + 1}}$ is min-max optimal in Section 5 (see also Kerkycharian & Picard (1992); Donoho et al. (1996); Donoho & Johnstone (1998); Giné & Nickl (2015) for ordinary Besov spaces). The L^1 -norm constraint $k f^0 k_1 \leq F$ is used to derive a uniform bound on the discrepancy between the population and the empirical L^2 -norm. Without this condition, the convergence rate could be slower.

Deep composition model For the deep composition model, we obtain the following convergence rate. This is an extension of Theorem 2 but requires a stronger assumption on the smoothness.

Theorem 3. *Suppose that $0 < p; q \leq 1$ and $\tilde{(\cdot)} > 1=\rho$ for all $\tilde{(\cdot)} \in [H]$. If $f^0 \in H_{\text{deep}} \setminus L^1(\cdot)$, and $k f^0 k_1 \leq F$ for $F \geq 1$, then we obtain*

$$E_{D_n}[k f^0 - \hat{f} k_{L^2(P_X)}^2] \leq \max_{\tilde{(\cdot)} \in [H]} n^{-\frac{2^{\tilde{(\cdot)}}}{2^{\tilde{(\cdot)}} + 1}} \log(n)^3;$$

where $\tilde{\gamma}(\cdot)$ is defined in Theorem 1, and $(L; W; S; B)$ is as given in Theorem 1 with $N = \max_{\ell \in [L]} n^{2 - \frac{1}{\tilde{\gamma}(\ell) + 1}}$.

The proof is provided in Appendix C. We will show that this is also minimax optimal in Theorem 4. Because of the Hölder continuity, the convergence rate becomes slower than the affine composition model (that is, $\tilde{\gamma}(\cdot) = \tilde{\gamma}(\cdot)$). However, this slower rate is unavoidable in terms of the minimax optimal rate. Schmidt-Hieber (2020) analyzed the same situation for the Hölder class which corresponds to $\tilde{\gamma}(\cdot) = \frac{(\cdot)}{d}(\beta)$ and $p = q = 1$. Our analysis far extends their analysis to the setting of anisotropic Besov spaces in which the parameters $(\cdot); p; q$ have much more freedom.

From these two bounds (Theorems 2 and 3), we can see that as the smoothness $\tilde{\gamma}$ becomes large, the convergence rates faster. If the target function is included in the isotropic Besov space with smoothness $\gamma = \frac{(\cdot)}{d}(\beta)$, then the estimation error becomes

$$\text{(Isotropic Besov)} \quad n^{-2 - (2 - \gamma)}$$

In the exponent, the dimensionality d appears, which causes the curse of dimensionality. In contrast, if the target function is in the anisotropic Besov space, and the smoothness in each direction is sufficiently imbalanced such that $\tilde{\gamma}$ does not depend on d , our obtained rate

$$\text{(Anisotropic Besov)} \quad n^{-2 - (\tilde{\gamma} + 1)}$$

avoids the curse of dimensionality. For high-dimensional settings, there would be several redundant directions in which the true function does not change. Deep learning is adaptive to this redundancy and achieves a better estimation error. However, in Section 6, we prove that linear estimators are affected by the dimensionality more strongly than deep learning. This indicates the superiority of deep learning.

5 Minimax optimal rate

Here, we show that the estimation error rate, that we have presented, of deep learning achieves the *minimax optimal rate*. Roughly speaking the minimax optimal risk on a model F of the true function is the smallest worst case error over all estimators:

$$R(F) := \inf_{\hat{f}} \sup_{f \in F} E_{D_n} [k\hat{f} - f^0 k_{L^2(P_X)}^2];$$

where \hat{f} runs over all estimators. The convergence rate of the minimax optimal risk is referred to as minimax optimal rate. We obtain the following minimax optimal rate for anisotropic Besov spaces.

Theorem 4. (a) Affine composition model: For $0 < p; q \leq 1$ and $\gamma \in \mathbb{R}_{++}^d$, assume that $\tilde{\gamma} > \max\{1-p, 1-2; 1-p, 1\}; 0 < \gamma$: Then, the minimax optimal risk of the affine composition model is lower bounded as $R(H_a) \geq n^{-\frac{2}{\tilde{\gamma} + 1}}$. **(b) Deep composition model:** For $0 < p; q \leq 1$ and $(\cdot) \in \mathbb{R}_{++}^d$ ($\cdot = 1; \dots; H$), assume that $\tilde{\gamma}(\cdot) > 1-p$: Let $\epsilon > 0$ be arbitrarily small for $q < 1$, and let $\gamma = 0$ for $q = 0$. Let $\tilde{\gamma}(\cdot) = \tilde{\gamma}(\cdot) \prod_{k=1}^H [(\cdot - 1 - p + \epsilon)^k + 1]$, and $\tilde{\gamma} := \min_{\cdot} \tilde{\gamma}(\cdot)$. Then, the minimax optimal risk of the deep composition model is lower bounded as $R(H_{\text{deep}}) \geq n^{-\frac{2}{\tilde{\gamma} + 1}}$.

The proof is provided in Appendix E (see also Ibragimov & Khas'minskii (1984); Nyssbaum (1987)). From this theorem, we can see that the estimation error of deep learning shown in Theorems 2 and 3 indeed achieve the minimax optimal rate up to a poly-log(n) factor.

6 Suboptimality of linear estimators

In this section, we give the minimax optimal rate in the class of *linear estimators*. The linear estimator is a class of estimators that can be written as

$$\hat{f}(x) = \sum_{i=1}^n y_i \phi_i(x; X^n);$$

where $X^n = (x_1; \dots; x_n)$ and $\phi_i(x; X^n)$ ($i = 1; \dots; n$) are (measurable) functions that only depend on x and X^n . This is linearly dependent on $Y^n = (y_1; \dots; y_n)$. We notice that the kernel

ridge regression is included in this class because it can be written as $\hat{f}(x) = k_{X, X^n}(k_{X^n, X^n} + \lambda)^{-1} Y^n$, which linearly depends on Y^n . This class includes other important estimators, such as the Nadaraya–Watson estimator, the k -nearest neighbor estimator, and the sieve estimator. We compare deep learning with the linear estimators in terms of minimax risk. For this purpose, we define the minimax risk of the class of linear estimators:

$$R^{(\text{lin})}(F) := \inf_{\hat{f}: \text{linear}} \sup_{f \in \mathcal{F}} E_{D_n}[k f^0 - \hat{f}]_{L^2(P_X)}^2;$$

where \hat{f} runs over all *linear estimators*. We can see that linear estimators suffer from the sub-optimal rate because of the following two points: (i) they do not have adaptivity, and (ii) they significantly suffer from the curse of dimensionality.

Theorem 5. (i) Suppose that the input distribution P_X is the uniform distribution on $\mathbb{R}^d = [0; 1]^d$ and assume that $\rho > 1 = \rho$ and $1 - \rho; q \leq 1$. Then, the minimax optimal rate of the linear estimators is lower bounded as

$$R^{(\text{lin})}(U(B_{p,q})) \asymp n^{-\frac{2-\nu}{2+1-\nu}}; \quad (6)$$

where $\nu = 2(1-\rho - 1=2)_+$.

(ii) In addition to the above conditions, we assume that the true function is included in the affine composition model with $d = d, \rho = 1 = \rho = \rho_d$ and $0 < \rho \leq 2$. Let $a_d = 1 + \rho$ with arbitrary small $\rho > 0$ when $d < d=2$, and let $a_d = 0$ when $d = d=2$. Then, the minimax rate of the linear estimators on the affine composition model is lower bounded by

$$R^{(\text{lin})}(H_a) \asymp n^{-\frac{2(\rho - d = \rho + d = 2 + a_d)}{2(\rho - d = \rho + d = 2 + a_d) + d}}; \quad (7)$$

The proof is provided in Appendix F. (i) The lower bound (7) reveals the suboptimality of linear estimators in terms of input dimensionality. Actually, if we consider a particular case where $d = 1, \rho = 1$ and $d = d$, then the obtained minimax rate of linear estimators and the estimation error rate of deep learning can be summarized as

$$\text{linear} : n^{-\frac{2+d}{2+2d}}; \quad \text{deep} : n^{-\frac{2}{2+1}};$$

by Theorem 2 when $\rho > 1$ (which can be checked by noticing $d = \rho = \rho = 1$ in this situation). We can see that the dependence on the dimensionality of linear estimators is significantly worse than that of deep learning. This indicates poor adaptivity of linear estimators to the intrinsic dimensionality of data. Actually, as d becomes large, the rate for the linear estimator approaches to $1 = \rho = \rho$ but that for the deep learning is not affected by d and still faster than $1 = \rho = \rho$. To show the theorem, we used the ‘‘convex-hull argument’’ developed by Hayakawa & Suzuki (2019); Donoho & Johnstone (1998). We combined this technique with the so-called Irie-Miyake’s integral representation (Irie & Miyake, 1988; Hornik et al., 1990). Note that this difference appears because there is an affine transformation in the first layer of the affine composition model. Deep learning is flexible against such a coordinate transform so that it can find directions to which the target function is smooth. In contrast, kernel methods do not have such adaptivity because there is no feature extraction layer. (ii) The lower bound (6) states that when $\rho < 2$ (that is, $\nu > 0$), the minimax rate of the linear estimators is outperformed by that of deep learning (Theorem 2). This is due to the ‘‘adaptivity’’ of deep learning. When ρ is small, the smoothness of the target function is less homogeneous, and it requires an adaptive approximation scheme to achieve the best estimation error. Linear estimators do not have adaptivity and thus fail to achieve the minimax optimal rate. Our bound (6) extends the result by Zhang et al. (2002) to a multivariate anisotropic Besov space while Zhang et al. (2002) investigated the univariate space ($d = 1$).

7 Conclusion

We investigated the approximation error and estimation error of deep learning in the anisotropic Besov spaces. It was proved that the convergence rate is determined by the average of the anisotropic smoothness, which results in milder dependence on the input dimensionality. If the smoothness is

highly anisotropic, deep learning can avoid overfitting. We also compared the error rate of deep learning with that of linear estimators and showed that deep learning has better dependence on the input dimensionality. Moreover, it was shown that deep learning can achieve the adaptive rate and outperform non-adaptive approximation methods and linear estimators if the homogeneity ρ of smoothness is small. These analyses strongly support the practical success of deep learning from a theoretical perspective.

Limitations of this work Our work does not cover the optimization aspect of deep learning. It is assumed that the regularized least squares (5) can be executed. It would be nice to combine our study with recent developments of non-convex optimization techniques (Vempala & Wibisono, 2019; Suzuki & Akiyama, 2021).

Potential negative societal impact Since this is purely theoretical result, it is not expected that there is a direct negative societal impact. However, revealing detailed properties of the deep learning could promote an opportunity to pervert deep learning.

Acknowledgment

TS was partially supported by JSPS KAKENHI (18H03201), Japan Digital Design and JST CREST. AN was partially supported by JSPS Kakenhi (19K20337) and JST-PRESTO.

References

- M. Belkin and P. Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*, 15(6):1373–1396, 2003.
- A. Bhattacharya, D. Pati, and D. B. Dunson. Adaptive dimension reduction with a gaussian process prior. *arXiv preprint arXiv:1111.1044*, 1445, 2011.
- A. Bhattacharya, D. Pati, and D. Dunson. Anisotropic function estimation using multi-bandwidth gaussian processes. *Annals of statistics*, 42(1):352, 2014.
- P. J. Bickel and B. Li. Local polynomial regression on unknown manifolds. In *Complex datasets and inverse problems*, pp. 177–186. Institute of Mathematical Statistics, 2007.
- M. Chen, H. Jiang, W. Liao, and T. Zhao. Nonparametric Regression on Low-Dimensional Manifolds using Deep ReLU Networks. *arXiv e-prints*, art. arXiv:1908.01842, Aug 2019.
- M. Chen, H. Jiang, W. Liao, and T. Zhao. Efficient approximation of deep relu networks for functions on low dimensional manifolds. In *Advances in Neural Information Processing Systems*, pp. 8172–8182, 2019.
- G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals, and Systems (MCSS)*, 2(4):303–314, 1989.
- J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv e-prints*, art. arXiv:1810.04805, Oct 2018.
- R. A. DeVore. Nonlinear approximation. *Acta Numerica*, 7:51–150, 1998.
- R. A. DeVore and V. A. Popov. Interpolation of Besov spaces. *Transactions of the American Mathematical Society*, 305(1):397–414, 1988.
- R. A. DeVore, G. Kyriazis, D. Leviatan, and V. M. Tikhomirov. Wavelet compression and nonlinear-widths. *Advances in Computational Mathematics*, 1(2):197–214, 1993.
- D. L. Donoho and I. M. Johnstone. Minimax estimation via wavelet shrinkage. *The Annals of Statistics*, 26(3):879–921, 1998.
- D. L. Donoho, I. M. Johnstone, G. Kerkycharian, and D. Picard. Density estimation by wavelet thresholding. *The Annals of Statistics*, 24(2):508–539, 1996.

- D. Dũng. B-spline quasi-interpolant representations and sampling recovery of functions with mixed smoothness. *Journal of Complexity*, 27(6):541–567, 2011a.
- D. Dũng. Optimal adaptive sampling recovery. *Advances in Computational Mathematics*, 34(1): 1–41, 2011b.
- S. Gaïffas and G. Lecue. Hyper-sparse optimal aggregation. *Journal of Machine Learning Research*, 12(Jun):1813–1833, 2011.
- E. Giné and R. Nickl. *Mathematical Foundations of Infinite-Dimensional Statistical Models*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2015.
- X. Glorot, A. Bordes, and Y. Bengio. Deep sparse rectifier neural networks. In *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pp. 315–323, 2011.
- H. Hang and I. Steinwart. Optimal learning with anisotropic gaussian svms. *arXiv preprint arXiv:1810.02321*, 2018.
- S. Hayakawa and T. Suzuki. On the minimax optimality and superiority of deep neural network learning over sparse parameter spaces. *arXiv preprint arXiv:1905.09195*, 2019.
- M. Hoffman and O. Lepski. Random rates in anisotropic regression (with a discussion and a rejoinder by the authors). *The Annals of Statistics*, 30(2):325–396, 04 2002.
- K. Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2): 251–257, 1991.
- K. Hornik, M. Stinchcombe, and H. White. Universal approximation of an unknown mapping and its derivatives using multilayer feedforward networks. *Neural Networks*, 3(5):551–560, 1990.
- I. Ibragimov and R. Khas'minskii. More on the estimation of distribution densities. *Journal of Soviet Mathematics*, 25(3):1155–1165, 1984.
- B. Irie and S. Miyake. Capabilities of three-layered perceptrons. In *IEEE 1988 International Conference on Neural Networks*, pp. 641–648, 1988.
- G. Kerkycharian and D. Picard. Density estimation in Besov spaces. *Statistics & Probability Letters*, 13:15–24, 1992.
- G. Kerkycharian, O. Lepski, and D. Picard. Nonlinear estimation in anisotropic multi-index de-noising. *Probability Theory and Related Fields*, 121(2):137–170, Oct 2001.
- A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- C. Leisner. Nonlinear wavelet approximation in anisotropic besov spaces. *Indiana University mathematics journal*, pp. 437–455, 2003.
- S. Mallat. *A Wavelet Tour of Signal Processing*. Academic Press, 1999.
- V. Myronyuk. Trigonometric approximations and kolmogorov widths of anisotropic besov classes of periodic functions of several variables. *Ukrainian Mathematical Journal*, 66(8), 2015.
- V. V. Myronyuk. Kolmogorov widths of the anisotropic besov classes of periodic functions of many variables. *Ukrainian Mathematical Journal*, 68(5):718–727, Oct 2016.
- V. V. Myronyuk. Widths of the anisotropic besov classes of periodic functions of several variables. *Ukrainian Mathematical Journal*, 68(8):1238–1251, Jan 2017. ISSN 1573-9376. doi: 10.1007/s11253-017-1290-1. URL <https://doi.org/10.1007/s11253-017-1290-1>.
- V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on Machine Learning*, pp. 807–814, 2010.

- R. Nakada and M. Imaizumi. Adaptive approximation and generalization of deep neural network with intrinsic dimensionality. *Journal of Machine Learning Research*, 21(174):1–38, 2020. URL <http://jmlr.org/papers/v21/20-002.html>.
- S. M. Nikol'skii. *Approximation of functions of several variables and imbedding theorems*, volume 205. Springer-Verlag Berlin Heidelberg, 1975.
- M. Nyssbaum. Optimal filtration of a function of many variables in white gaussian noise. *Problems of Information Transmission*, 19:23–29, 1983.
- M. Nyssbaum. Nonparametric estimation of a regression function that is smooth in a domain in \mathbb{R}^k . *Theory of Probability & Its Applications*, 31(1):108–115, 1987.
- A. Radford, L. Metz, and S. Chintala. Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. *arXiv e-prints*, art. arXiv:1511.06434, Nov 2015.
- G. Raskutti, M. J. Wainwright, and B. Yu. Minimax-optimal rates for sparse additive models over kernel classes via convex programming. *The Journal of Machine Learning Research*, 13(1):389–427, 2012.
- J. Schmidt-Hieber. Nonparametric regression using deep neural networks with ReLU activation function. *ArXiv preprint arXiv:1708.06633(v3)*, 2018.
- J. Schmidt-Hieber. Deep ReLU network approximation of functions on a manifold. *arXiv preprint arXiv:1908.00695*, 2019.
- J. Schmidt-Hieber. Nonparametric regression using deep neural networks with relu activation function. *The Annals of Statistics*, 48(4):1875–1897, 2020.
- P. Y. Simard, D. Steinkraus, and J. C. Platt. Best practices for convolutional neural networks applied to visual document analysis. In *Proceedings of the Seventh International Conference on Document Analysis and Recognition-Volume 2*, pp. 958. IEEE Computer Society, 2003.
- S. Sonoda and N. Murata. Neural network with unbounded activation functions is universal approximator. *Applied and Computational Harmonic Analysis*, 43(2):233–268, 2017.
- T. Suzuki. Fast generalization error bound of deep learning from a kernel perspective. In *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pp. 1397–1406. PMLR, 2018.
- T. Suzuki. Adaptivity of deep ReLU network for learning in Besov and mixed smooth Besov spaces: optimal rate and curse of dimensionality. In *International Conference on Learning Representations (ICLR2019)*, 2019. URL <https://openreview.net/forum?id=H1ebTsActm>.
- T. Suzuki and S. Akiyama. Benefit of deep learning with non-convex noisy gradient descent: Provable excess risk bound and superiority to kernel methods. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=2m0g1wEafh>.
- V. Temlyakov. *Approximation of Periodic Functions*. Nova Science Publishers, 1993.
- J. B. Tenenbaum, V. De Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *science*, 290(5500):2319–2323, 2000.
- H. Triebel. Entropy numbers in function spaces with mixed integrability. *Revista matemática complutense*, 24(1):169–188, 2011.
- A. B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer Series in Statistics. Springer New York, 2008.
- A. W. van der Vaart and J. A. Wellner. *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer, New York, 1996.
- S. Vempala and A. Wibisono. Rapid convergence of the unadjusted langevin algorithm: Isoperimetry suffices. In *Advances in Neural Information Processing Systems*, pp. 8094–8106, 2019.

- J. Vybiral. Function spaces with dominating mixed smoothness. *Dissertationes Math. (Rozprawy Mat.)*, 436:3–73, 2006.
- Y. Yang and A. Barron. Information-theoretic determination of minimax rates of convergence. *The Annals of Statistics*, 27(5):1564–1599, 1999.
- Y. Yang and D. B. Dunson. Bayesian manifold regression. *The Annals of Statistics*, 44(2):876–905, 2016.
- Y. Yang and S. T. Tokdar. Minimax-optimal nonparametric regression in high dimensions. *The Annals of Statistics*, 43(2):652–674, 2015.
- D. Yarotsky. Error bounds for approximations with deep relu networks. *Neural Networks*, 94:103–114, 2017.
- S. Zhang, M.-Y. Wong, and Z. Zheng. Wavelet threshold estimation of a regression function with random design. *Journal of Multivariate Analysis*, 80(2):256–284, 2002.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes]
 - (c) Did you discuss any potential negative societal impacts of your work? [No]
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes] Each theorem explicitly describes assumptions on which the theorem is established.
 - (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix in the supplementary material. All proofs are included there.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [N/A]
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A]
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [N/A]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [N/A]
 - (b) Did you mention the license of the assets? [N/A]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

—Appendix—

A Approximation error of Affine composition model

Theorem 6 (Affine composition model). *Assume that the distribution of $x = Ax + b \in \mathbb{R}^d$ has a bounded density function on $[0; 1]^d$ when x obeys the uniform distribution on $[0; 1]^d$, and each element in A and b is bounded by a constant C . Assume that $0 < p; q; r \leq 1$ and $\tilde{\gamma} \in \mathbb{R}_{++}^d$ satisfy $\tilde{\gamma} > (1-p \ 1-r)_+$. Then, it holds that*

$$R_r(\mathcal{L}_1(d); W_1(d); S_1(d); (dC + 1)B_1(d); H_a) \leq N_{\tilde{\gamma}}; \quad (8)$$

where $\mathcal{L}_1(\cdot); W_1(\cdot); S_1(\cdot); B_1(\cdot)$ are defined in Eq. (3).

The assumption $\tilde{\gamma} > (1-p \ 1-r)_+$ ensures the L_r -integrability of the target function, and the inequality (without equality) admits a near-optimal wavelet approximation of the target function in terms of L_r -norm. From this theorem, the approximation error is almost identical to that for $B_{p,q}(\cdot)$ (Proposition 2).

B Proofs of approximation error bounds

To show the approximation accuracy, a key step is to show that the ReLU neural network can approximate the cardinal B-spline with high accuracy. Let $N(x) = 1$ ($x \in [0; 1]$); 0 (otherwise), then the cardinal B-spline of order m is defined by taking $m + 1$ -times convolution of N :

$$N_m(x) = \underbrace{(N * N * \dots * N)}_{m+1 \text{ times}}(x);$$

where $f * g(x) := \int f(x-t)g(t)dt$. It is known that N_m is a piece-wise polynomial of order m . For $k \in \mathbb{Z}_+^d$ and $j = (j_1; \dots; j_d) \in \mathbb{Z}_+^d$, let

$$M_{k;j}^d(x) = \prod_{i=1}^d N_m(2^{bk} \cdot_i^c x_i - j_i);$$

where $\cdot_i^c \in \mathbb{R}_{++}^d$ is a given smoothness parameter (we omit the dependency on \cdot_i^c from the notation which would be obvious from the context). Here, k controls the spatial ‘‘resolution’’ and j specifies the location on which the basis is put. Basically, we approximate a function f in an anisotropic Besov space via super-position of $M_{k;j}^m(x)$, which is closely related to wavelet analysis (Mallat, 1999). The following is a key lemma that was proven by Suzuki (2019).

Lemma 1 (Approximation of cardinal B-spline basis by the ReLU activation). *There exists a constant $c_{(d;m)}$ depending only on d and m such that, for all $\epsilon > 0$, there exists a neural network $M \in \mathcal{N}(L_0; W_0; S_0; B_0)$ with $L_0 := 3 + 2 \lceil \log_2 \left(\frac{3^d \cdot m}{c_{(d;m)}} \right) + 5 \rceil$, $d \log_2(d - m) \leq W_0 \leq 6dm(m+2) + 2d$, $S_0 := L_0 W_0^2$ and $B_0 := 2(m+1)^m$ that satisfies*

$$\|kM_{0;0}^d - M_{k;L^1(\mathbb{R}^d)}\| \leq \epsilon;$$

and $M(x) = 0$ for all $x \notin [0; m+1]^d$.

Let

$$k_k \cdot_j^c := \sum_{j=1}^d b_k \cdot_j^c$$

for a $k \in \mathbb{Z}$. For order $m \in \mathbb{N}$ of the cardinal B-spline bases, let

$$J_i(k) = f \cdot m; \ m + 1; \dots; \cdot 2^{bk} \cdot_i^c \ 1; 2^{bk} \cdot_i^c g$$

and

$$J(k) := J_1(k) \ J_2(k) \ \dots \ J_d(k);$$

and the quasi-norm of the coefficient $(k_j)_{k,j}$ for $k \in \mathbb{Z}_+$ and $j \in J(k)$ be

$$k(k_j)_{k,j} k_{b_{p,q}} := \left\{ \sum_{k=0}^1 \left[2^{k \lfloor \frac{d}{2} \rfloor} \left(\sum_{i=1}^d b_{i,c=k} \right)^{1-p} \left(\sum_{j \in J(k)} k_j j^p \right)^{1-p} \right]^q \right\}^{1-q}.$$

For $p = 1$ or $q = 1$, the definition should be appropriately modified as usual.

Lemma 2. Assume the condition $\tilde{m} > (1-p)(1-r)_+$ in Proposition 2 and $0 < \tilde{m} < \min(m; m-1+1-p)$ where $m \in \mathbb{N}$ is the order of the cardinal B-spline bases. Then, $f \in B_{p,q}$ admits the following decomposition:

$$f = \sum_{k=0}^1 \sum_{j \in J(k)} k_j M_{k,j}^d(x) \quad (9)$$

with convergence in the sense of L^p , and the coefficient (k_j) yields the following norm equivalence

$$k f k_{B_{p,q}} \approx k(k_j)_{k,j} k_{b_{p,q}}. \quad (10)$$

For an integer $K \in \mathbb{N}$, let $N = d^{kK} e$, then for any $f \in B_{p,q}$, there exists f_N that satisfies

$$k f - f_N k_{L^r(\cdot)} \leq N^{-\tilde{m}} k f k_{B_{p,q}};$$

and has the following form:

$$f_N(x) = \sum_{k=0}^K \sum_{j \in J(k)} k_j M_{k,j}^d(x) + \sum_{k=K+1}^K \sum_{i=1}^{n_k} k_{j_i} M_{k,j_i}^d(x); \quad (11)$$

where $K = dK(1+1) e$, $n_k = d^{kK} e - (kK) e$ ($k = K+1; \dots; K$) for $\tilde{m} = (1-p)(1-r)_+$ and $\tilde{m} = (2) e$, and $(j_i)_{i=1}^{n_k} \in J(k)$.

Proof of Lemma 2. Leisner (2003) showed that there exists a bounded linear operator P_k that can be expressed as

$$P_k(f)(x) = \sum_{j \in J(k)} a_{k,j} M_{k,j}^d(x) \quad (12)$$

where k_j is constructed in a certain way, and for every $f \in L^p([0;1]^d)$ with $0 < p < 1$, it holds

$$k f - P_k(f) k_{L^p} \leq C_{W,r,p}(f; (2^{-k} e_1; \dots; 2^{-k} e_d));$$

(See Theorem 3.2.4 of Leisner (2003) and DeVore & Popov (1988)). Let

$$\rho_k(f) := P_k(f) - P_{k-1}(f); \quad P_{-1}(f) = 0.$$

Then, Leisner (2003) showed that when $0 < p; q < 1$ and $0 < \tilde{m} < \min(m; m-1+1-p)$, f belongs to $B_{p,q}$ if and only if f can be decomposed into

$$f = \sum_{k=0}^1 \rho_k(f);$$

with the convergence condition

$$k(\rho_k(f))_{k=0}^1 k_{b_q(L^p)} := \left[\sum_{k=0}^1 (2^{-k} k \rho_k k_{L^p})^q \right]^{1-q} < 1;$$

In particular, it is shown that

$$k f k_{B_{p,q}^s} \approx k(\rho_k(f))_{k=0}^1 k_{b_q^s(L^p)}. \quad (13)$$

Here, each p_k can be expressed as $p_k(x) = \sum_{j \in J(k)} a_{k,j} M_{k,j}^d(x)$ for a coefficient $(a_{k,j})_{k,j}$ which could be different from $(a_{k,j})_{k,j}$ appearing in Eq. (12), and thus $f \in B_{p,q}$ can be decomposed into

$$f = \sum_{k=0}^1 \sum_{j \in J(k)} a_{k,j} M_{k,j}^d(x)$$

with convergence in the sense of L^p . Moreover, it is shown that $\|p_k\|_{L^p} \leq (2^{-kd} \sum_{j \in J(k)} |a_{k,j}|^p)^{1/p}$ and thus

$$\|f\|_{B_{p,q}} \leq \left(\sum_{k=0}^1 \sum_{j \in J(k)} |a_{k,j}|^p \right)^{1/p};$$

This yields the first assertion.

Next, we move to the second assertion. If $p \geq r$, the assertion can be shown in the same manner as Theorem 3.1 of Düng (2011a). More precisely, we can show the assertion in a similar line to the following proof for $p < r$ by setting $K = K$. Thus, we show the assertion only for $p < r$. In this regime, we need to use an adaptive approximation method. In the following, we assume $p < r$. For a given K , by appropriately choosing K later, we set

$$R_K(f)(x) = \sum_{0 \leq k \leq K} p_k + \sum_{k \in \mathbb{Z}_+ : K < k} G_k(p_k);$$

where $G_k(p_k)$ is given as

$$G_k(p_k) = \sum_{i=1}^{n_k} a_{k,i} M_{k,i}^d(x)$$

where $(a_{k,i})_{i=1}^{n_k}$ is the sorted coefficients in decreasing order of absolute value: $|a_{k,1}| \geq |a_{k,2}| \geq \dots \geq |a_{k,n_k}|$. Then, it holds that

$$\|p_k\|_{L^r} \leq \|G_k(p_k)\|_{L^r} \leq \|p_k\|_{L^p} 2^{k(k-p)/r};$$

where $\gamma := (1/p - 1/r)$ (see the proof of Theorem 3.1 of Düng (2011b) and Lemma 5.3 of Düng (2011a)). Moreover, we also have

$$\|p_k\|_{L^r} \leq \|p_k\|_{L^p} 2^{k(k-p)\gamma}$$

for $k \in \mathbb{Z}_+$ with $k > K$.

Here, we define N as

$$N = 2^{K(K-p)\gamma} e;$$

Let $\tilde{K} = \lfloor N \rfloor$,

$$K = \lfloor N \rfloor + 1 = \tilde{K} + 1;$$

and

$$n_k = \left\lceil 2^{k(K-p)\gamma} \left(\sum_{i=1}^{n_k} |a_{k,i}|^p \right)^{1/p} \right\rceil$$

for $k \in \mathbb{Z}_+$ with $K+1 \leq k \leq \tilde{K}$.

Then, by Lemma 5.3 of Düng (2011a), we have

$$\begin{aligned} \|f - R_K(f)\|_{L^r} &\leq \sum_{K < k \leq \tilde{K}} \|p_k\|_{L^r} + \sum_{K < k} \|p_k\|_{L^r} \\ &\leq \sum_{K < k \leq \tilde{K}} [\|p_k\|_{L^p} 2^{k(k-p)\gamma}]^r + \sum_{K < k} [2^{k(k-p)\gamma} \|p_k\|_{L^p}]^r; \end{aligned} \quad (14)$$

(a) Suppose that $q \geq r$ and $r < 1$. Then,

$$\begin{aligned} \|f - R_K(f)\|_{L^q} &= \|f - R_K(f)\|_{L^r}^{q/r} \\ &\leq \left\{ \sum_{K < k \leq \tilde{K}} [2^{k(k-p)\gamma} \|p_k\|_{L^p}]^r + \sum_{K < k} [2^{k(k-p)\gamma} \|p_k\|_{L^p}]^r \right\}^{q/r} \end{aligned} \quad (* \text{ Eq. (14)})$$

$$\begin{aligned}
& \sum_{K < k} [2^{kkk_-} n_k k p_k k_{LP}]^q + \sum_{K < k} [2^{kkk_-} k p_k k_{LP}]^q \\
& N^{-q} \sum_{K < k} [2^{(kkk_-)kKk_-} k p_k k_{LP}]^q \\
& + 2^{-q} \sum_{K < k} [2^{kkk_-} k p_k k_{LP}]^q \\
& \cdot (N^{-2} \sum_{K < k} [2^{(kkk_-)kKk_-} + 2^{(kkk_-)kKk_-}]^q k f k_{MB_{p,q}}^q) \quad (* \text{ Eq. (13)}) \\
& \cdot (N^{-1})^q k f k_{MB_{p,q}}^q
\end{aligned}$$

where we used $2^{kkk_-} = 2^{-k}$ in (i), and $N^{-1} = 2^{kKk_-}$ and $(kkk_-)kKk_- = (kkk_-)kKk_-$ in (ii).

(b) Suppose that $q > r$ and $r < 1$. Then, letting $\theta = q - r (> 1)$ and $\theta^0 = 1 - (1 - \theta) = q - r$ (note that $1 + \frac{1}{\theta} = 1$), we have

$$\begin{aligned}
& k f R_K(f) k_{L^r} \cdot \sum_{K < k} [2^{kkk_-} n_k k p_k k_{LP}]^r + \sum_{K < k} [2^{kkk_-} k p_k k_{LP}]^r \quad (* \text{ Eq. (14)}) \\
& 2^{-r} \sum_{K < k} [2^{(kkk_-)kKk_-} k p_k k_{LP}]^r \\
& + \sum_{K < k} [2^{kkk_-} k p_k k_{LP}]^r (2^{(kkk_-)kKk_-})^r \\
& (2^{-r} kKk_- + 2^{(kkk_-)kKk_-})^r \left\{ \sum_{K < k} [2^{(kkk_-)kKk_-} k p_k k_{LP}]^r \right. \\
& \left. + \sum_{K < k} [2^{kkk_-} k p_k k_{LP}]^r 2^{(kkk_-)kKk_-} \right\} \\
& (2^{-r} kKk_- + 2^{(kkk_-)kKk_-})^r \left\{ \sum_{K < k} [2^{kkk_-} k p_k k_{LP}]^r + \sum_{K < k} [2^{kkk_-} k p_k k_{LP}]^r \right\}^{1-\theta} \\
& \left\{ \sum_{K < k} [2^{(kkk_-)kKk_-}]^{r \theta^0} + \sum_{K < k} [2^{(s)kKk_-} k]^{r \theta^0} \right\}^{1-\theta} \\
& \cdot (2^{-r} kKk_- + 2^{(kkk_-)kKk_-})^r k f k_{B_{p,q}}^r \quad (* \text{ Eq. (13) and } 2^{kkk_-} = 2^{-k}) \\
& \cdot (N^{-1})^r k f k_{B_{p,q}}^r :
\end{aligned}$$

(c) Suppose that $r = 1$. Then, similarly to the analysis in (b), we can evaluate

$$\begin{aligned}
& k f R_K(f) k_{L^r} \\
& \cdot 2^{-1} \sum_{K < k} [2^{(kkk_-)kKk_-} k p_k k_{LP}] \\
& + \sum_{K < k} [2^{kkk_-} k p_k k_{LP}] (2^{(kkk_-)kKk_-}) \\
& \cdot (2^{-1} kKk_- + 2^{(kkk_-)kKk_-}) k f k_{B_{p,q}} \\
& \cdot N^{-1} k f k_{B_{p,q}} :
\end{aligned}$$

This concludes the proof. \square

Proof of Proposition 2. We adopt the proof line employed by [Suzuki \(2019\)](#). Basically, we combine Lemma 1 and Lemma 2. We substitute the approximated cardinal B-spline basis M into the

decomposition of f_N (11). Let the set of indexes $(k;j) \in \mathbb{Z} \times \mathbb{Z}$ that consists f_N given in Eq. (11) be E_N , i.e., $f_N = \sum_{(k;j) \in E_N} M_{k;j}^d$. Accordingly, we set $f := \sum_{(k;j) \in E_N} M_{k;j}^d$. Note that for each x , the number of $(k;j) \in E_N$ that satisfy $M_{k;j}(x) \neq 0$ is bounded by $(m+1)^d(1+K)$, and $\max_{(k;j) \in E_N} |M_{k;j}(x)| \leq 2^K = (\tilde{c}^{-1-p})_+$ by the norm equivalence Eq. (10). For each $x \in \mathbb{R}^d$, it holds that

$$\begin{aligned} |f_N(x) - f(x)| &\leq \sum_{(k;j) \in E_N} |M_{k;j}^d(x) - M_{k;j}(x)| \\ &\leq \sum_{(k;j) \in E_N} |M_{k;j}^d(x)| \\ &\leq (m+1)^d(1+K) 2^K (\tilde{c}^{-1-p})_+ \|f\|_{B_{p,q}^s} \\ &\leq \log(N) N^{(1-p)(\tilde{c}^{-1-p})_+} \|f\|_{B_{p,q}^s}, \end{aligned}$$

where we used the definition of K in the last inequality. This evaluation yields that, for each $f \in U(B_{p,q}(\cdot))$, it holds that

$$\|f\|_{L^r} \leq \|f_N\|_{L^r} + \|f - f_N\|_{L^r} \leq \log(N) N^{(1-p)(\tilde{c}^{-1-p})_+} \|f\|_{B_{p,q}^s} + N^{-\tilde{c}}.$$

By taking \tilde{c} to satisfy $\log(N) N^{(1-p)(\tilde{c}^{-1-p})_+} = N^{-\tilde{c}}$, we obtain the approximation error bound.

As we have seen above $\max_{(k;j) \in E_N} |M_{k;j}^d| \leq 2^K = (\tilde{c}^{-1-p})_+ N^{(1-p)(\tilde{c}^{-1-p})_+}$. The max of the absolute values of parameters used in $M_{k;j}^d$ can be bounded by 2^K (see Suzuki (2019)) which is bounded by $N^{d(1-p)(\tilde{c}^{-1-p})_+}$. Then, we obtain the assertion. \square

B.1 Proof of Theorem 6 and Theorem 1

Proof of Theorem 6. This proof is almost obvious from Proposition 2. We know that, from Proposition 2, for $g \in U(B_{p,q}([0;1]^d))$, there exists $f \in (L_1(d); W_1(d); S_1(d); B_1(d))$ such that

$$\|f - g\|_{L^r} \leq N^{-\tilde{c}}.$$

Because the density of the distribution of $Ax + b$ is bounded above when x obeys the uniform distribution on $[0;1]^d$, this also yields

$$\|f(A+b)\|_{L^r} \leq \|g(A+b)\|_{L^r} \leq N^{-\tilde{c}}.$$

(note that the Lebesgue measure on $[0;1]^d$ corresponds to the uniform distribution on $[0;1]^d$). If f can be written as

$$f(x) = (W^{(L_1)}(x) + b^{(L)}) \cdot (W^{(1)}x + b^{(1)});$$

then we have

$$\|f(A+b)\|_{L^r} = \|W^{(L_1)}(A+b) + b^{(L)}\|_{L^r} \leq \|W^{(1)}A + b^{(1)} + W^{(1)}b\|_{L^r} \leq (L_1(d); W_1(d); S_1(d); (dC+1)B_1(d));$$

\square

Proof of Theorem 1.

$$H_{\text{deep}} := f h_H \quad h_{\cdot;k}(x) \in [0;1]^{m_{\cdot}}; [0;1]^{m_{\cdot+1}}; h_{\cdot;k} \in U(B_{p,q}^{(\cdot)}([0;1]^{m_{\cdot}})) \quad (\delta_k \in [m_{\cdot+1}])g;$$

Since $\tilde{c}^{(\cdot)} > 1-p$, we can show that for each $h_{\cdot;k}$, there exists $f_{\cdot;k} \in (L_1(m_{\cdot}); W_1(m_{\cdot}); S_1(m_{\cdot}); B_1(m_{\cdot}))$ such that

$$\|f_{\cdot;k} - h_{\cdot;k}\|_{L^1} \leq N^{-\tilde{c}}.$$

Moreover, from the proof of Proposition 2, we can share all parameters other than the last layer among $f_{\cdot;k}$ ($k = 1, \dots, m_{\cdot+1}$). If necessary, we may modify $f_{\cdot;k}$ so that $0 \leq f_{\cdot;k}(x) \leq 1$ ($x \in [0;1]^{m_{\cdot}}$).

$[0; 1]^m$) by adding one additional clipping layer which can be realized by ReLU (actually, the clipping operator can be constructed by a linear combination of 2 nodes with ReLU activation as $f(x) = \max\{x; 0\} - \max\{x - 1; 0\} = \min\{x, 1\}$ for $x \in \mathbb{R}$). The approximation error of the whole layer can be evaluated as

$$\begin{aligned} & k h_H \quad h_1 \quad f_H \quad f_1 k_1 \\ & \sum_{l=1}^H k h_H \quad h_{l+1} \quad h_l \quad f_{l-1} \quad f_l \quad h_H \quad h_{l+1} \quad f_l \quad f_{l-1} \quad f_1 k_1 \\ & \sum_{l=1}^H k h_H \quad h_{l+1} \quad h_l \quad h_H \quad h_{l+1} \quad f_l k_1 : \end{aligned}$$

Proposition 1 tells that $h_{\cdot, k} \in C_{(-1, p)}^{(\beta)}$; thus, $h_{\cdot, k}$ is β -Hölder continuous where $\beta := (\beta - 1/p) \wedge 1$. Their composition $h_H \circ h_{H-1} \circ \dots \circ h_{1+1}$ is β -Hölder continuous where $B_\beta = \prod_{l=1}^H \beta_l$. Therefore, we have

$$k h_H \quad h_{l+1} \quad h_l \quad h_H \quad h_{l+1} \quad f_l k_1 \leq k h_{\cdot} \quad f_{\cdot} k_1^{B_\beta};$$

where $k \quad k_1$ for a vector-valued function $g : \mathbb{R}^d \rightarrow \mathbb{R}^{d_0}$ is defined as $\sup_x \|g(x)\|_k$. Summing up this evaluation for $l = 1, \dots, H$ concludes that

$$k h_H \quad h_1 \quad f_H \quad f_1 k_1 \leq \sum_{l=1}^H N^{B_\beta} \cdot \max_{2H} N^{\beta(l)};$$

Consequently, the whole network can be realized as an element of $(L; W; S; B)$ where

$$\begin{aligned} L &= \sum_{l=1}^H (L_1(m_{\cdot}) + 1); \quad W = \max(W_1(m_{\cdot}) - m_{\cdot+1}); \\ S &= \sum_{l=1}^H (S_1(m_{\cdot}) + 3m_{\cdot+1}); \quad B = \max B_1(m_{\cdot}); \end{aligned}$$

□

C Proofs of estimation error bound (Theorem 2 and Theorem 3)

Proof of Theorem 2. We follow the proof strategy from Schmidt-Hieber (2018); Suzuki (2019) which uses Proposition 4. It suffices to the covering number of $\hat{F} = \{f \circ f \circ \dots \circ f \mid f \in (L; W; S; B)g\}$ for $(L; W; S; B)$ given in Theorem 6 where f is the clipped version of a function f . Note that the covering number of \hat{F} is not larger than that of $(L; W; S; B)$. Hence, it is sufficient to evaluate that of $(L; W; S; B)$. From Lemma 6, the covering number of this class is upper bounded by

$$\log N(\cdot; \hat{F}; k \quad k_1) \leq N \log(N) [\log(N)^2 + \log(\cdot)];$$

From Proposition 2, there exists $f \in (L; W; S; B)$ such that

$$k f^0 \leq R_K (f^0) k_2 \leq N^{\tilde{\beta}};$$

Moreover, we notice that $k f \leq f^0 k_{L^2(P_X)} \leq R k f \leq f^0 k_2^2$: for any $f : [0; 1]^d \rightarrow \mathbb{R}$ because the density p_X of P_X is bounded by R . Therefore, by applying Proposition 4 with $\beta = 1/n$, we have

$$E_{D_n} [k \hat{F} \leq f^0 k_{L^2(P_X)}] \leq N^{\tilde{\beta}} + \frac{N \log(N) (\log(N)^2 + \log(n))}{n} + \frac{1}{n};$$

Here, we can minimize the right hand side by setting $N = n^{\frac{1}{2-\tilde{\beta}}}$ up to $\log(n)^3$ -order, and then we obtain the estimation error of the least squares estimator as

$$n^{-\frac{\tilde{\beta}}{2-\tilde{\beta}}} \log(n)^3;$$

This yields the assertion. □

Proof of Theorem 3. The proof is almost identical to the proof of Theorem 2, except that we use Theorem 1 as an approximation error bound. □

D Embedding theorem

Lemma 3. For $0 < p^{(1)}; p^{(2)} \leq 1$, let $(\cdot)^{(1)}; (\cdot)^{(2)} \in \mathbb{R}_{++}^d$ such that they satisfy

$$\frac{1}{p^{(1)}} \leq \frac{1}{p^{(2)}}; \quad (15)$$

$$p^{(2)} \leq p^{(1)};$$

$$p^{(1)} < p^{(2)};$$

for $0 < \dots < 1$. Then, it holds that

$$B_{p^{(1)}; q}^{(1)} \leq B_{p^{(2)}; q}^{(2)}$$

Proof. We show the assertion only for the situation where $p^{(1)} \notin 1$, $p^{(2)} \notin 1$, and $q \notin 1$. The proof for the setting in which $p^{(1)} = 1$, $p^{(2)} = 1$, or $q = 1$ is satisfied is almost identical. Recall the following norm equivalence shown in Lemma 2:

$$kfk_{B_{p; q}^s} \leq k(\cdot; j)_{k; j} k_{b_{p; q}} = \left\{ \sum_{k=0}^1 \left[2^{k[-(\sum_{i=1}^d bk_i^{(1)} c=k)=p]} \left(\sum_{j \in 2J(k)} j_{k; j} j^p \right)^{1=p} \right]^q \right\}^{1=q};$$

when $p; q < 1$. Since $\frac{p^{(1)}}{p^{(2)}} < 1$, it holds that

$$\begin{aligned} \left(\sum_{j \in 2J(k)} j_{k; j} j^{p^{(1)}} \right)^{1=p^{(1)}} &= \left(\sum_{j \in 2J(k)} j_{k; j} j^{p^{(2)} \frac{p^{(1)}}{p^{(2)}}} \right)^{1=p^{(1)}} \\ &= \left(\sum_{j \in 2J(k)} j_{k; j} j^{p^{(2)}} \right)^{\frac{p^{(1)}}{p^{(2)}} \frac{1}{p^{(1)}}} = \left(\sum_{j \in 2J(k)} j_{k; j} j^{p^{(2)}} \right)^{\frac{1}{p^{(2)}}}; \end{aligned}$$

Moreover, we have

$$\begin{aligned} &2^{k[-(\sum_{i=1}^d bk_i^{(1)} c=k)=p^{(1)}]} \\ &\leq 2^{k[-(\sum_{i=1}^d bk_i^{(1)} c=k)=p^{(1)}]} = 2^{k[-(\sum_{i=1}^d bk_i^{(1)} c=k)=p^{(1)}]} \left(\frac{1}{p^{(1)}} \right) \stackrel{(a)}{=} 2^{k[-(\sum_{i=1}^d bk_i^{(1)} c=k)=p^{(1)}]} \left(\frac{1}{p^{(1)}} + \frac{1}{p^{(2)}} - \frac{1}{p^{(2)}} \right) \\ &\stackrel{(b)}{=} 2^{k[-(\sum_{i=1}^d bk_i^{(2)} c=k)=p^{(2)}]} = 2^{k[-(\sum_{i=1}^d bk_i^{(2)} c=k)=p^{(2)}]}, \end{aligned}$$

where we used the condition $p^{(2)} \leq p^{(1)}$ in (a), and we used the condition from Eq. (15) in (b). These relations yield the following evaluation:

$$\begin{aligned} &kfk_{B_{p^{(1)}; q}^{(1)}} \leq k(\cdot; j)_{k; j} k_{b_{p^{(1)}; q}^{(1)}} \\ &= \left\{ \sum_{k=0}^1 \left[2^{k[-(\sum_{i=1}^d bk_i^{(1)} c=k)=p^{(1)}]} \left(\sum_{j \in 2J(k)} j_{k; j} j^{p^{(1)}} \right)^{1=p^{(1)}} \right]^q \right\}^{1=q} \\ &\leq \left\{ \sum_{k=0}^1 \left[2^{k[-(\sum_{i=1}^d bk_i^{(2)} c=k)=p^{(2)}]} \left(\sum_{j \in 2J(k)} j_{k; j} j^{p^{(2)}} \right)^{1=p^{(2)}} \right]^q \right\}^{1=q} \\ &\leq kfk_{B_{p^{(2)}; q}^{(2)}}; \end{aligned}$$

This yields the assertion. \square

By combining Lemma 3 with the relation $B_{1; 1} \leq C$ – (Triebel, 2011), we immediately obtain the following corollary.

Corollary 1. Suppose that $\tilde{p} > p$, then for $\tilde{p} = \frac{p}{\dots}$; it holds that

$$B_{p; q} \leq B_{1; q} \leq B_{1; 1} \leq C;$$

E Minimax optimality

In this section, we demonstrate the proof of Theorem 4. Before this, we prepare the basic notions. The ϵ -covering number $N(\epsilon; C; \hat{d})$ of a metric space C equipped with a metric \hat{d} that is the minimal number of balls with radius ϵ measured by the metric \hat{d} required to cover the set C (van der Vaart & Wellner, 1996). Similarly, the ϵ -packing number $M(\epsilon; C; \hat{d})$ is defined as the largest number of elements $f_1, \dots, f_M \in C$ such that $\hat{d}(f_i; f_j) \geq \epsilon$ for all $i \neq j$.

Raskutti et al. (2012) showed the following inequality in their proof of Theorem 2(b) by utilizing the result by Yang & Barron (1999).

Lemma 4. *Let F be the model of the true function. For a given $n > 0$ and $n' > 0$, let Q be the n' -packing number $M(n'; F; L^2(P_X))$ of F and N be the n -covering number of that. Suppose that they satisfy the following condition:*

$$\frac{n}{2} \leq \frac{n'}{n} \log(N);$$

$$8 \log(N) \leq \log(Q); 4 \log(2) \leq \log(Q); \quad (16)$$

Then, the minimax learning rate is lower bounded as

$$\inf_{\hat{f}} \sup_{f \in 2F} E_{D_n}[k\hat{f} - f]_{L^2(P_X)}^2 \geq \frac{2}{4}.$$

This concludes the assertion.

Now, we are ready to show Theorem 4.

Proof of Theorem 4. Proposition 10 of Triebel (2011) showed that the ϵ -covering number of the unit ball of anisotropic Besov spaces $B_{p,q}(\cdot)$ can be evaluated as

$$\log N(\epsilon; U(B_{p,q}(\cdot)); k_{k_r}) \sim \epsilon^{-1};$$

for $0 < p, q \leq 1, 1 \leq r < \infty$, and \mathbb{R}_+^d that satisfy

$$\epsilon^{-1} \geq \max \left\{ \frac{1}{p}, \frac{1}{r}, \frac{1}{p}, 1, 0 \right\};$$

Affine composition model:

Apparently, $U(B_{p,q}(\cdot))$ is included in H_a . Hence, noting that P_X is the uniform distribution and $k_{k_2} = k_{L^2(P_X)}$, the covering number of H_a can be lower bounded by

$$\log N(H_a; k_{L^2(P_X)}) \leq \epsilon^{-1};$$

From this evaluation, Lemma 4 yields that there exists $C_1 > 0$ independent of n such that

$$\inf_{\hat{f}} \sup_{f \in 2H_{\text{aff}}} E_{D_n}[k\hat{f} - f]_{L^2(P_X)}^2 \leq C_1 n^{-\frac{2}{2^*+1}};$$

To see this, we may just set $n' = n' = n^{-\frac{2}{2^*+1}}$ in Eq. (16) of Lemma 4.

Deep composition model:

Next, we show the minimax rate for the deep composition model. Basically, we follow the same strategy developed by Schmidt-Hieber (2018), but we need to modify some technical details because we are dealing with anisotropic Besov spaces while Schmidt-Hieber (2018) analyzed isotropic Hölder space. Let $\dot{\cdot} := \min_{\cdot \in [H]} \dot{\cdot}(\cdot)$, and $s(\cdot) := (\dot{\cdot}(\cdot) - 1/p + \epsilon)^{\wedge 1}$ where $\epsilon > 0$ can be arbitrary small for $q < 1$ and $\epsilon = 0$ for $q = 1$. Without loss of generality, we may assume that $\dot{\cdot}(\cdot) \geq \frac{(\cdot)}{2}$ for $\cdot \in [H]$. Let us consider a sub-model H_{deep}^0 of H_{deep} defined as

$$H_{\text{deep}}^0 := f g_H \quad g_1 j$$

$$\begin{aligned}
g_{\cdot}(x) &= (x_{\cdot} = 1; \dots; x_{\cdot} = 1); \\
g_{\cdot}(x) &= (g_{\cdot-1}(x); 0; \dots; 0)^{\succ} \text{ where } g_{\cdot-1} \in U(B_{p,q}(\cdot)); \\
g_{\cdot}(x) &= (x_1^{s^{(\cdot)}}; 0; \dots; 0)^{\succ} (\cdot = \cdot + 1; \dots; H)g_{\cdot}
\end{aligned}$$

For $\cdot = \cdot + 1; \dots; H$, through a cumbersome calculation, we can verify that $x_1^{s^{(\cdot)}} \in B_{p,q}^{(\cdot)}([0;1])$ for $x \in [0;1]$, which ensures $g_{\cdot,j}(x) \in B_{p,q}^{(\cdot)}([0;1]^d)$ for $j = 1; \dots; d$. To lower bound the covering number, we concretely construct a subset the cardinality of which can be easily estimated. For that purpose, we use the expansion $f = \sum_{k=0}^{\infty} \sum_{j \in \hat{J}(k)} k_{\cdot,j} M_{k,j}^d(x)$ and the norm equivalence $\|f\|_{B_{p,q}(\cdot)} \approx \|k_{\cdot,j}\|_{k_{\cdot,j} B_{p,q}(\cdot)}$ given in Lemma 2. For a while, we let $\cdot := (\cdot)$ and $B := \prod_{q=\cdot+1}^H s^{(\cdot)}$. We define $k \in \mathbb{N}$ so that k satisfies $2^{k-1} \approx n^{\frac{1}{1+2B}}$. For this choice of k , take a subset $\hat{J}(k) \subseteq J(k)$ such that $j \in \hat{J}(k) \Rightarrow j \in J(k)$ and for each $j; j^0 \in \hat{J}(k)$ with $j \neq j^0$, the supports of $M_{k,j}^d$ and M_{k,j^0}^d are disjoint. Using this index set $\hat{J}(k)$, we consider a set of functions that is given by

$$\hat{H}_{\cdot} := \left\{ f = \sum_{j \in \hat{J}(k)} k_{\cdot,j} M_{k,j}^d(x) \mid k_{\cdot,j} \in \{0, 2^{-k-g}\} \right\}.$$

We can check that $\|f\|_{B_{p,q}(\cdot)} = 2^{k \sum_{j=1}^d j} = 2^{k-1}$ and $\|f\|_{B_{p,q}(\cdot)} = 1$ for all $f \in \hat{H}_{\cdot}$ from the norm equivalence (10). For any $g_w = \sum_{j \in \hat{J}(k)} w_j 2^{-k} M_{k,j}^d(x) \in \hat{H}_{\cdot}$ ($w_j \in \{0, 1\}^{|\hat{J}(k)|}$), we can see that

$$\begin{aligned}
f_w(x) &= g_H \quad g_{\cdot+1} \quad g_w \quad g_{\cdot-1} \quad g_1(x) \\
&= \sum_{j \in \hat{J}(k)} w_j 2^{-k} M_{k,j}^d(x).
\end{aligned}$$

If $w \notin W^0$, then we can see that

$$\begin{aligned}
\|f_w\|_{L^2(P_X)} &\approx \text{Ham}(w; W^0) 2^{-k-1} \\
&\approx \text{Ham}(w; W^0) 2^{-k-(2B+1)};
\end{aligned}$$

where Ham is the Hamming distance because $\|M_{k,j}^d\|_{L^2(P_X)} \approx 2^{-k-1}$.

Then, by the Varshamov–Gilbert bound (see Lemma 2.9 of Tsybakov (2008), for example), there exists a subset $W_k \subseteq \{0, 1\}^{|\hat{J}(k)|}$ such that $|W_k| \geq 2^{|\hat{J}(k)|/8}$ and $\text{Ham}(w; W^0) \leq |\hat{J}(k)|/8$ for all $w; w \in W_k$ with $w \notin W^0$. This yields

$$\|f_w\|_{L^2(P_X)} \approx 2^{-k-1} 2^{-k-(2B+1)} = 2^{-2k-1} \approx n^{-\frac{2B}{2B+1}};$$

where the definition of k is used. This implies that there exists a subset $H_{\text{deep}}^0 \subseteq H_{\text{deep}}^0$ such that

$$\log(N(n; H_{\text{deep}}^0; k_{L^2(P_X)})) \approx n^{\frac{1}{1+2B}}$$

for $n \approx n^{\frac{B}{2B+1}}$. Then, by Lemma 4, we obtain that the minimax optima rate on H_{deep} is lower bounded as

$$\inf_{\hat{f}} \sup_{f \in F} E_{D_n}[\|\hat{f} - f\|_{L^2(P_X)}^2] \approx n^{-\frac{B}{2B+1}}.$$

□

F Minimax optimal rate of linear estimators

Define the convex hull of a function class F as

$$\text{conv}(F) := \left\{ f(x) = \sum_{j=1}^M \lambda_j f_j(x) \mid M = 1; 2; \dots; f_j \in F; \lambda_j \geq 0; \sum_{j=1}^M \lambda_j = 1 \right\};$$

Let $\overline{\text{conv}}(\cdot)$ is the closure of $\overline{\text{conv}}(\cdot)$ with respect to $L_2(P_X)$ -norm.

Proposition 3 (Hayakawa & Suzuki (2019)). *The minimax optimal rate of linear estimators on a target function class F is the same as that on the convex hull of F :*

$$\inf_{\hat{f}: \text{linear } F \rightarrow \mathbb{R}} \sup_{f \in F} E_{D_n}[k f^\circ - \hat{f}]_{L^2(P_X)}^2 = \inf_{\hat{f}: \text{linear } F \rightarrow \overline{\text{conv}}(F)} \sup_{f \in F} E_{D_n}[k f^\circ - \hat{f}]_{L^2(P_X)}^2;$$

See Hayakawa & Suzuki (2019) for the proof of this proposition.

Proof of Theorem 5. We basically follow the strategy developed by Zhang et al. (2002). Let μ be the uniform measure on \mathcal{X} . They essentially showed the following statement in their Theorem 1. Suppose that the space \mathcal{X} has even partition A such that $|A| = 2^K$ for an integer $K \geq \mathbb{N}$, each A has equivalent measure $\mu(A) = 2^{-K}$ for all $A \in \mathcal{A}$, and \mathcal{A} is indeed a partition of \mathcal{X} , i.e., $\bigcup_{A \in \mathcal{A}} A = \mathcal{X}$, $A \cap A' = \emptyset$ for $A, A' \in \mathcal{A}$ and $A \neq A'$. Then, if K is chosen as $n^{-1} \leq 2^K \leq n^{-2}$ for constants $\epsilon_1, \epsilon_2 > 0$ that are independent of n , then there exists an event E such that, for a constant $C^\circ > 0$,

$$P(E) \geq 1 - o(1);$$

We call this property of \mathcal{A} ‘‘Condition A.’’

Here, we consider a set F of functions on \mathcal{X} for which there exists $F > 0$ that satisfies the following conditions:

1. There exists $F > 0$ such that, for any $A \in \mathcal{A}$, there exists $g \in F$ that satisfies $g(x) \geq \frac{1}{2} F$ for all $x \in A$,
2. There exists K° and $C^\circ > 0$ such that $\frac{1}{n} \sum_{i=1}^n g(x_i)^2 \leq C^\circ 2^{-2K^\circ}$ for any $g \in F$ on the event E .

We call this condition of the function class F ‘‘Condition B.’’

Let the minimax optimal rate of linear estimators on the function class F be

$$R = \inf_{\hat{f}: \text{linear } F \rightarrow \mathbb{R}} \sup_{f \in F} E_{D_n}[k \hat{f} - f]_{L^2(P_X)}^2;$$

Then, under Conditions A and B, there exists a constant F_1 such that at least one of the following inequalities holds:

$$\frac{F^2}{4F_1 C^\circ} 2^{K^\circ} \leq R; \tag{17a}$$

$$\frac{F^3}{32} 2^{2K} \leq R; \tag{17b}$$

for sufficiently large n .

(i) *Proof of Eq. (6).*

For given $k \geq \mathbb{N}$ (which will be fixed later), let $\mathcal{C} = 2^{-k} \lfloor \sum_{i=1}^d b_k \lfloor \frac{C}{b_k} \rfloor \rfloor$. Then, from the wavelet expansion of anisotropic Besov space (9),

$$f_w = \sum_{j \in J(k)} w_j M_{k,j}^d(x) \in CU(B_{p,q}(\cdot));$$

where $C > 0$ is a constant and $w = (w_j)_{j \in J(k)}$ is a one-hot vector, i.e., $w_j = 1$ for some $j \in J(k)$ and $w_{j'} = 0$ for all $j' \in J(k)$ with $j' \neq j$. This expansion ensures that, for $K = \sum_{i=1}^d b_k \lfloor \frac{C}{b_k} \rfloor$, there exists a partition \mathcal{A} of \mathcal{X} that satisfies Condition A, and for any $A \in \mathcal{A}$, there exists w such that $f_w(x) \geq \frac{1}{2} C$ for all $x \in A$ and

$$\frac{1}{n} \sum_{i=1}^n f_w(x_i)^2 \leq \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{x_i \in A\}} \leq 2^{-K};$$

on the event E , which ensures that $F = \frac{1}{n} \sum_{j=1}^n w_j$ is a one-hot vector g satisfies Condition B. Hence, by choosing $k \geq N$ so that $2^k \leq n^{\frac{1}{2} + \frac{1}{p} + 1}$ (recall that $K = \sum_{i=1}^d b_k \frac{1}{C}$ by definition), and setting $K = K^0$, then Eq. (17) gives

$$R \leq n^{\frac{2-v}{2-v+1}};$$

for $v = 2(1-p^{-1})$. This yields the assertion because $F \in CU(B_{p,q}(R))$ for a constant C .

(ii) *Proof of Eq. (7).*

Let $\mu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i} = \mu_d = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$. For m such that $\frac{1}{m} < \min\{f, m^{-1} + 1 - pg\}$, let $\mu_d(x) = \prod_{j=1}^d N_m(x_j - (m+1)^{-2}) (x \in \mathbb{R}^d)$.

(ii-a) *Setting of $d = d=2$:*

Let $V_{d,d} := \{U \in \mathbb{R}^d \mid U^T U = I_d\}$ be the Stiefel manifold and let $\nu_{d,d}$ be the invariant measure on the Stiefel manifold (i.e., the uniform distribution). Then, let $\mu_d : \mathbb{R}^d \rightarrow \mathbb{R}$ be

$$\mu_d(x) = \int \mu_d(Ux) d\nu_{d,d}(U) \quad (x \in \mathbb{R}^d);$$

We can see that μ_d is spherically symmetric and there exists $F, C > 0$ such that

$$\mu_d(x) \leq F (8x \in \mathbb{R}^d \text{ s.t. } \|x\| = 1);$$

and

$$\mu_d(x) \begin{cases} C \|x\|^d & (\|x\| = 1); \\ 1 & (\|x\| = 1); \end{cases}$$

The last inequality can be checked by the fact that for a sufficiently large $R > 0$, the measure of the set $\{x \in \mathbb{R}^d \mid \|x\| = R; \mu_d(x) > 0\} \cap \mathbb{R}^d \stackrel{1}{=} \mathbb{R}^{d-1} = \mathbb{R}^d$ (here, \mathbb{R} is the uniform probability measure on the sphere $S_{d-1}(R) = \{x \in \mathbb{R}^d \mid \|x\| = R\}$ and $\int_{S_{d-1}(R)} 1 = 1$).

By the construction of μ_d and the wavelet expansion of anisotropic Besov space (9) with the norm equivalence (10), we have that there exists a constant $c > 0$ such that, for any $k \geq N$ and $b = [\frac{1}{2} - 2^{-k}(\frac{m+1}{2} - \lfloor \frac{m+1}{2} \rfloor)] (1, \dots, 1) \in \mathbb{R}^d$, it holds that

$$c \mu_d(2^k(\cdot - b)) \geq U(B_{p,q}([0; 1]^d));$$

where $\mu_d = 2^{-k(d-p)}$. Here, let $0 < c < 1$ be a constant such that $cU(x - b^0) + b \geq [0; 1]^d$ for any $x, b^0 \in [0; 1]^d$ and any $U \in V_{d,d}$. Then, we have that, for any $b^0 \in [0; 1]^d$,

$$c \mu_d(2^k c U(\cdot - b^0)) = c \mu_d(2^k(\cdot - b)) \quad (cU(\cdot - b^0) + b) \geq H_a;$$

for any $U \in V_{d,d}$. By the convex hull argument (Proposition 3), this yields that

$$R^{\text{lin}}(H_a) = R^{\text{lin}}(\text{conv}(H_a)) \leq R^{\text{lin}}(c \mu_d(2^k c U(\cdot - b^0)) \mid b^0 \in [0; 1]^d);$$

Hence, it suffices to lower bound the far right-hand side of this inequality. We consider a partition A of \mathbb{R}^d , where $A \subseteq \mathbb{R}^d$ has the form $A = [2^{-k}j_1; 2^{-k}(j_1 + 1)] \times \dots \times [2^{-k}j_d; 2^{-k}(j_d + 1)]$ for $0 \leq j_i \leq 2^k - 1$ ($i = 1, \dots, d$). Let $\hat{J}(k) = \{j_1, \dots, j_d \mid 0 \leq j_i \leq 2^k - 1\}$ and $A_j = [2^{-k}j_1; 2^{-k}(j_1 + 1)] \times \dots \times [2^{-k}j_d; 2^{-k}(j_d + 1)] \in A$ for $j \in \hat{J}(k)$. Let $\mu_{A_j} = c \mu_d(2^k c(\cdot - b_{A_j}))$, where $b_{A_j} = (2^{-k}(j_1 + 1), \dots, 2^{-k}(j_d + 1)) \in \mathbb{R}^d$ for $j \in \hat{J}(k)$. We can see that $\mu_{A_j} = 2^{dk}$. Hence, A satisfies Condition A with $K = dk$ if 2^k is in polynomial order with respect to n .

Moreover, there exists $F > 0$ such that $\mu_{A_j}(x) \leq F$ for all $x \in A$. Next, we evaluate $\frac{1}{n} \sum_{i=1}^n \mu_{A_j}(x_i)^2$. On the event E , there exists C^0 such that $\frac{1}{n} \sum_{i=1}^n \mu_{A_j}(x_i)^2 \leq C^0 n^{-2k} = C^0 n^{-2k}$ for all $A^0 \in A$. Here, let

$$\mu_{A_j}(x) := \begin{cases} c C^0 2^k c(x - b_{A_j})^d & (2^k c(x - b_{A_j}) \leq 1); \\ c & (\text{otherwise}); \end{cases}$$

then $\int_{A^0} f_A(x) dx$. Thus, we can upper bound $\frac{1}{n} \sum_{i=1}^n \int_{A^0} f_A(x_i)^2 dx$ as

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \int_{A^0} f_A(x_i)^2 dx &= \frac{1}{n} \sum_{i=1}^n \int_{A^0} f_A(x_i)^2 dx = \frac{1}{n} \sum_{A^0 \ni A} \sum_{x_i \in A} \int_{A^0} f_A(x_i)^2 dx = \frac{1}{n} \sum_{A^0 \ni A} C^0 \frac{n}{2^k} \max_{x \in A} \int_{A^0} f_A(x)^2 dx \\ &= C^0 \sum_{A^0 \ni A} (A^0) \max_{x \in A} \int_{A^0} f_A(x)^2 dx = C^0 \sum_{A^0 \ni A} (A^0) \min_{x \in A} \int_{A^0} f_A(x)^2 dx \frac{\max_{x \in A} \int_{A^0} f_A(x)^2 dx}{\min_{x \in A} \int_{A^0} f_A(x)^2 dx} \\ &= C^0 \sum_{A^0 \ni A} (A^0) \min_{x \in A} \int_{A^0} f_A(x)^2 dx \frac{\max_{x \in A} \int_{A^0} f_A(x)^2 dx}{\min_{x \in A} \int_{A^0} f_A(x)^2 dx} \\ &= C^0 \sum_{A^0 \ni A} (A^0) \min_{x \in A} \int_{A^0} f_A(x)^2 dx \max_{x: k2^k c(x-b_A)^k} \frac{k2^k c(x-b_A)^k}{1 + (k2^k c(x-b_A)^k + c)2^k} \\ &= C^0 \sum_{A^0 \ni A} (A^0) \min_{x \in A} \int_{A^0} f_A(x)^2 dx (1 + c \frac{2^k}{d})^{2d} \\ &= C^0 (1 + c \frac{2^k}{d})^{2d} \int_{A^0} f_A(x)^2 dx \end{aligned}$$

The quantity $\int_{A^0} f_A(x)^2 dx$ on the right-hand side can be evaluated as

$$\begin{aligned} \int_{A^0} f_A(x)^2 dx &= \int_{x: kx-b_A \leq 2^k} f_A(x)^2 dx + \int_{x: c \frac{2^k}{d} < kx-b_A \leq 2^k} f_A(x)^2 dx \\ &= 2^{-kd} + C \frac{2^{2d}}{2^{2kd}} \int_{c \frac{2^k}{d} < r \leq 2^k} r^{-2d} r^{d-1} dr \\ &= 2^{-kd} + 2^{-2kd} \max_{r \in [c \frac{2^k}{d}, 2^k]} r^{k(2d-d)}; 1g \\ &= \max_{r \in [c \frac{2^k}{d}, 2^k]} r^{-kd}; 2^{-2kd}g; \end{aligned}$$

Therefore, we have that, for a constant C^0 , on the event E , we have that

$$\frac{1}{n} \sum_{i=1}^n \int_{A^0} f_A(x_i)^2 dx \leq C^0 (2^{-kd} + 2^{-2kd});$$

Let $F = f_A^j \mathbb{1}_{A \geq Ag}$, then F satisfies Condition B. When $d = d=2$, by choosing k so that $c2^k < n^{\frac{1}{2(d-p)}}$ and $K = K^0 = dk$, then Eq. (17) yields

$$R^{\text{lin}}(F) \leq n^{\frac{2(d-p+d-2)}{2(d-p+d-2)+d}};$$

This concludes the proof.

(ii-b) Setting of $d < d=2$:

Let A be the partition of \mathbb{R}^d as defined in the proof for $d = d=2$, i.e., $jA_j = 2^{dk}$ and each $A \geq A$ can be written as $A = [2^{-k}j_1; 2^{-k}(j_1+1)] \times \dots \times [2^{-k}j_d; 2^{-k}(j_d+1)]$ for $0 \leq j_i \leq 2^k - 1$ ($i = 1; \dots; d$). Pick up $A \geq A$ and let $j \in \hat{J}(k)$ be the index such that $A = A_j$. For a while, we fix A and let $b = b_{A_j} = (2^{-k}(j_1+1); \dots; 2^{-k}(j_d+1))$ accordingly. For $w = (w; b) \in \mathbb{R}^d \times \mathbb{R}^d$, let $A : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be

$$A(x) := 2^k [x_1 - b_1; \dots; x_d - b_d; w^>(x_{d,d} - b_{d,d}) + b];$$

and consider

$$\phi(x) := \phi_d(A(x));$$

We take its convex hull with respect to w . We note that

$$\phi(x) = \left(\prod_{j=1}^{d-1} N_m(2^k(x_j - b_j) \mid (m+1)=2) \right) N_m(2^k[w^>(x_{d,d} - b_{d,d}) + b] \mid (m+1)=2);$$

To analyze its convex hull, it suffices to consider the convex hull of the last term $N_m(2^k[w^{>}(x_{d;d} b_{d;d}) + b])$ ($m+1=2$). Hence, we set $(\cdot) := N_m(\cdot)$ ($m+1=2$) and consider a set of functions

$$\mathbb{F}_{C;}^{(\cdot)} := \{f_X \geq \mathbb{R}^d \text{ } d+1 \text{ } \mathcal{I} \text{ } a(\cdot) (w^> x + b))\} \text{ } j \text{ } j \text{ } a \text{ } j \text{ } 2C; \text{ } k \text{ } w \text{ } k \text{ } 1; \text{ } j \text{ } b \text{ } j \text{ } 2(a; b \geq \mathbb{R}; w \geq \mathbb{R}^d \text{ } d+1) \text{ } g$$

for $C > 0$, $\delta > 0$. We also define the Fourier transform of $\hat{f}(l) := (2\pi)^{-1} \int e^{-il^T x} f(x) dx$ ($l \in \mathbb{R}$). Then, by Lemma 5, we have that, for $h = 2^{-k}$ and $\delta = h^{-1}$,

$$\inf_{g \in \text{conv}(\mathbb{F}_{C;}^{(\cdot)})} \sup_{x \in [0;1]^d} \left| g(x) \exp\left(\frac{kx \cdot ck^2}{2h^2}\right) \right| \\ \frac{4}{j^{2^{\wedge}(1)j}} \left[C_d \text{ } d+1 \text{ } R^{2(d \text{ } d \text{ } 1)} \exp(-R^2/2) + \exp(-R) \right];$$

where $C = \frac{1}{j^{2^{\wedge}(1)j}} = (h^{-1})$ and $R = h^{-1} (2^{\overline{d}} + 1)$. This indicates that, for a fixed $A \geq A$, the convex hull of the set $\bar{f}_A \text{ } j = (w; b) \geq \mathbb{R}^d \text{ } d+1 \text{ } R; \text{ } k \text{ } w \text{ } k \text{ } 1; \text{ } j \text{ } b \text{ } j \text{ } 2; \text{ } j \text{ } a \text{ } j \text{ } g$ where $\delta = 2^{-k} (d=p)$ contains \bar{f}_A which satisfies

$$\left\| \bar{f}_A \text{ } (2C)^{-1} \left(\prod_{j=1}^{d-1} N_m(2^k(x_j \text{ } b_j)) \text{ } (m+1=2) \right) \exp\left(\frac{kx_{d;d} \text{ } b_{d;d} k^2}{2h^2}\right) \right\|_1 \\ = O\left(2^{-k(1+\delta)} (h^{-1})^{2(d \text{ } d \text{ } 1)} \exp(-h^2/2) + \exp(-h) \right);$$

We can see that on the event E , it holds that

$$\frac{1}{n} \sum_{i=1}^n \bar{f}_A(x_i) \leq (A) (2^{-k(1+\delta)})^2 \cdot 2^{-kd} 2^{-2k(d-p+1)} 2^{-2k} = 2^{-2k(d-p+1+d-2)-2k};$$

for all $A \geq A$. Let $F = \bar{f}_A \text{ } j \text{ } A \geq A \text{ } g$, then F satisfies Condition B. Note that, by the definition of $\mathbb{F}_{C;}^{(\cdot)}$ it holds that $\bar{f}_A \in \text{conv}(H_a)$ for all $A \geq A$. Thus

$$R^{\text{lin}}(H_a) = R^{\text{lin}}(\text{conv}(H_a)) = R^{\text{lin}}(F);$$

Therefore, by choosing k such that $2^k \leq n^{2(\frac{1}{d-p+1+d-2} + d+2)}$, and setting $K = K^0 = dk$, then Eq. (17) gives

$$R^{\text{lin}}(F) \leq n^{\frac{2(\frac{d-p+1+d-2}{d-p+1+d-2} + d)}{2}};$$

□

Lemma 5 (Suzuki & Akiyama (2021)). Let $h > 0$ and $R := h^{-1} (2^{\overline{d}} + 1)$. Then, for $C = \frac{1}{j^{2^{\wedge}(1)j}}$, the Gaussian RBF kernel can be approximated by

$$\inf_{g \in \text{conv}(\mathbb{F}_{C;}^{(\cdot)})} \sup_{x \in [0;1]^d} \left| g(x) \exp\left(\frac{kx \cdot ck^2}{2h^2}\right) \right| \\ \frac{4}{j^{2^{\wedge}(1)j}} \left[C_d R^{2(d-2)} \exp(-R^2/2) + \exp(-R) \right]$$

for any $c \in [0;1]^d$, where C_d is a constant depending only on d . In particular, the right hand side is $O(\exp(-n))$ if $R = n$.

G Auxiliary lemmas

The following proposition which were shown in Schmidt-Hieber (2018); Hayakawa & Suzuki (2019); Suzuki (2018) is convenient to show the estimation error rate.

Proposition 4 (Schmidt-Hieber (2018); Hayakawa & Suzuki (2019)). Let F be a set of functions. Let \hat{f} be the least-squares estimator in F :

$$\hat{f} = \operatorname{argmin}_{f \in F} \sum_{i=1}^n (y_i - f(x_i))^2;$$

Assume that $kF^0_{k_1} \subset F$ and all $f \in F$ satisfies $kfk_1 \leq F$ for some $F > 1$. If $\epsilon > 0$ satisfies $N(\cdot; F; k_{k_1}) \leq 3$, then it holds that

$$\mathbb{E}_{D_n} [k\hat{f} - f^0_{k^2_{L^2(P_X)}}] \leq C \left[\inf_{f \in F} kf - f^0_{k^2_{L^2(P_X)}} + (F^2 + \epsilon) \frac{\log N(\cdot; F; k_{k_1})}{n} + (F + \epsilon) \right];$$

where C is a universal constant.

The following lemma provides the covering number of the deep neural network model.

Lemma 6 (Covering number evaluation). The covering number of $(L; W; S; B)$ can be bounded by

$$\log N(\cdot; (L; W; S; B); k_{k_1}) \leq S \log((L(B-1))^L (W+1)^{2L}) + 2SL \log((B-1)(W+1)) + S \log((L-1));$$

Proof of Lemma 6. Given a network $f \in (L; W; S; B)$ expressed as

$$f(x) = (W^{(L)}(x) + b^{(L)}) \dots (W^{(1)}x + b^{(1)});$$

let

$$A_k(f)(x) = (W^{(k-1)}(x) + b^{(k-1)}) \dots (W^{(1)}x + b^{(1)});$$

and

$$B_k(f)(x) = (W^{(L)}(x) + b^{(L)}) \dots (W^{(k)}(x) + b^{(k)});$$

for $k = 2; \dots; L$. Corresponding to the last and first layers, we define $B_{L+1}(f)(x) = x$ and $A_1(f)(x) = x$ respectively. Then, it is easy to see that $f(x) = B_{k+1}(f)(W^{(k)}(x) + b^{(k)}) A_k(f)(x)$. Now, suppose that a pair of different two networks $f, g \in (L; W; S; B)$ given by

$$f(x) = (W^{(L)}(x) + b^{(L)}) \dots (W^{(1)}x + b^{(1)}); \quad g(x) = (W^{(L)^0}(x) + b^{(L)^0}) \dots (W^{(1)^0}x + b^{(1)^0});$$

has parameters with distance $\|W^{(k)} - W^{(k)^0}\|_{k_1} \leq k$ and $\|b^{(k)} - b^{(k)^0}\|_{k_1} \leq k$. Now, note that $kA_k(f)_{k_1} \leq \max_j kW_{j::}^{(k-1)}_{k_1} kA_{k-1}(f)_{k_1} + kb^{(k-1)}_{k_1} \leq WBkA_{k-1}(f)_{k_1} + B(B-1)(W+1)kA_{k-1}(f)_{k_1} \leq (B-1)^{k-1}(W+1)^{k-1}$, and similarly, the Lipschitz continuity of $B_k(f)$ with respect to k_{k_1} -norm is bounded as $(BW)^{L-k+1}$. Then, it holds that

$$\begin{aligned} & \|f(x) - g(x)\| \\ &= \left| \sum_{k=1}^L B_{k+1}(g)(W^{(k)}(x) + b^{(k)}) A_k(f)(x) - B_{k+1}(g)(W^{(k)^0}(x) + b^{(k)^0}) A_k(f)(x) \right| \\ & \leq \sum_{k=1}^L (BW)^{L-k} k(W^{(k)}(x) + b^{(k)}) A_k(f)(x) - (W^{(k)^0}(x) + b^{(k)^0}) A_k(f)(x)_{k_1} \\ & \leq \sum_{k=1}^L (BW)^{L-k} [W(B-1)^{k-1}(W+1)^{k-1} + 1] \\ & \leq \sum_{k=1}^L (BW)^{L-k} (B-1)^{k-1}(W+1)^k \leq L(B-1)^{L-1}(W+1)^L. \end{aligned}$$

Thus, for a fixed sparsity pattern (the locations of non-zero parameters), the covering number is bounded by $(\lfloor L(B-1)^{L-1}(W+1)^L \rfloor)^S$. There are the number of configurations of the sparsity pattern is bounded by $\binom{W+1}{S} (W+1)^{LS}$. Thus, the covering number of the whole space is bounded as

$$(W+1)^{LS} \{ \lfloor L(B-1)^{L-1}(W+1)^L \rfloor \}^S = \lfloor L(B-1)^{L-1}(W+1)^{2L} \rfloor^S;$$

which yields the assertion. \square