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

Anonymous Author(s)

Affiliation

Address

email

Abstract

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

16 1 Introduction

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

24 Amongst representation ability analysis of deep neural networks such as universal approximation
25 ability (Cybenko, 1989; Hornik, 1991; Sonoda & Murata, 2017), approximation theory of deep
26 neural networks on typical function classes such as Hölder, Sobolev, and Besov spaces have been
27 extensively studied. In particular, analyses of deep neural networks with the ReLU activation (Nair
28 & Hinton, 2010; Glorot et al., 2011) have been recently developed. Schmidt-Hieber (2018) showed
29 that the deep learning with ReLU activations can achieve the minimax optimal estimation accuracy
30 to estimate composite functions in Hölder spaces by using the approximation theory of Yarotsky
31 (2017). Suzuki (2019) generalized this analysis to those on the *Besov space* and the *mixed smooth*
32 *Besov space* by utilizing the techniques developed in approximation theories (Temlyakov, 1993;
33 DeVore, 1998). It was shown that deep learning can achieve an *adaptive approximation* error rate
34 that is faster than that of (non-adaptive) linear approximation methods (DeVore & Popov, 1988;
35 DeVore et al., 1993; Dũng, 2011), and it outperforms any linear estimators (including kernel ridge
36 regression) in terms of the minimax optimal rate.

Table 1: Relationship between existing research and our work. β 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 $\tilde{\beta}$ 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 (2018)	Suzuki (2019)	Suzuki (2019)	Nakada & Imaizumi (2020); Schmidt-Hieber (2019); Chen et al. (2019)	This work
Estimation error	$\tilde{O}(n^{-\frac{2\beta}{2\beta+d}})$	$\tilde{O}(n^{-\frac{2\beta}{2\beta+d}})$	$\tilde{O}\left(n^{-\frac{2\beta}{2\beta+1}} \times \log(n)^{\frac{2(d-1)(u+\beta)}{1+2\beta}}\right)$	$\tilde{O}(n^{-\frac{2\beta}{2\beta+D}})$	$\tilde{O}(n^{-\frac{2\beta}{2\beta+1}})$

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

54 In practice, one of the typically expected properties of a true function on high-dimensional data is
55 that it is invariant against perturbations of an input in some specific directions (Figure 1). For ex-
56 ample, in image-recognition tasks, the target function must be invariant against the spatial shift of
57 an input image, which is utilized by data-augmentation techniques (Simard et al., 2003; Krizhevsky
58 et al., 2012). In this paper, we investigate the approximation and estimation abilities of deep learn-
59 ing on *anisotropic Besov spaces* (Nikol’skii, 1975; Vybiral, 2006; Triebel, 2011) (also called domi-
60 nated mixed-smooth Besov spaces). An anisotropic Besov space is a set of functions that have
61 “direction-dependent” smoothness, whereas ordinary function spaces such as Hölder, Sobolev, and
62 Besov spaces assume isotropic smoothness that is uniform in all directions. We consider a com-
63 position of functions included in an anisotropic Besov space, including several existing settings as
64 special cases; it includes analyses of the Hölder space Schmidt-Hieber (2018) and Besov space
65 Suzuki (2019), as well as the low-dimensional sub-manifold setting (Nakada & Imaizumi, 2020;
66 Schmidt-Hieber, 2019; Chen et al., 2019; Chen et al., 2019)¹. By considering such a space, we
67 can show that deep learning can alleviate curse of dimensionality if the smoothness in each direc-
68 tion is highly anisotropic. Interestingly, any linear estimator (including kernel ridge regression) has
69 worse dependence on the dimensionality than deep learning. Our contributions can be summarized
70 as follows:

- 71 • We consider a situation in which the target function is included in a class of anisotropic Besov
72 spaces and show that deep learning can avoid the curse of dimensionality *even if the input data*
73 *do not lie on a low-dimensional manifold*. Moreover, deep learning can achieve the optimal
74 adaptive approximation error rate and minimax optimal estimation error rate.
- 75 • We compare deep learning with general linear estimators (including kernel methods) and show
76 that deep learning has better dependence on the input dimensionality than linear estimators.

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.

77 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:

$$y_i = f^\circ(x_i) + \xi_i \quad (i = 1, \dots, n),$$

78 where x_i is generated from a probability distribution P_X on $[0, 1]^d$, $\xi_i \sim N(0, \sigma^2)$, and the
 79 data $D_n = (x_i, y_i)_{i=1}^n$ are independently identically distributed. f° is the true function that
 80 we want to estimate. We are interested in the mean squared estimation error of an estimator \hat{f} :
 81 $E_{D_n}[\|\hat{f} - f^\circ\|_{L^2(P_X)}^2]$, where $E_{D_n}[\cdot]$ indicates the expectation with respect to the training data D_n .
 82 We consider a least-squares estimator in the deep neural network model as \hat{f} (see Eq. (5)) and
 83 discuss its optimality. More specifically, we investigate how the ‘‘intrinsic dimensionality’’ of data
 84 affects the estimation accuracy of deep learning. For this purpose, we consider an *anisotropic Besov*
 85 *space* as a model of the target function.

86 2.1 Anisotropic Besov space

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

91 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

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

92 for $x \in \Omega$ with $x+rh \in \Omega$, otherwise, let $\Delta_h^r(f)(x) = 0$.

93 **Definition 1.** For a function $f \in L^p(\Omega)$ where $p \in (0, \infty]$, the r -th modulus of smoothness of f is
 94 defined by $w_{r,p}(f, t) = \sup_{h \in \mathbb{R}^d: |h_i| \leq t_i} \|\Delta_h^r(f)\|_p$, for $t = (t_1, \dots, t_d)$, $t_i > 0$.

95 With this modulus of smoothness, we define the anisotropic Besov space $B_{p,q}^\beta(\Omega)$ for $\beta =$
 96 $(\beta_1, \dots, \beta_d)^\top \in \mathbb{R}_{++}^d$ as follows.

97 **Definition 2** (Anisotropic Besov space $(B_{p,q}^\beta(\Omega))$). For $0 < p, q \leq \infty$, $\beta = (\beta_1, \dots, \beta_d)^\top \in \mathbb{R}_{++}^d$,
 98 $r := \max_i \lfloor \beta_i \rfloor + 1$, let the seminorm $|\cdot|_{B_{p,q}^\beta}$ be

$$|f|_{B_{p,q}^\beta} := \begin{cases} \left(\sum_{k=0}^{\infty} [2^k w_{r,p}(f, (2^{-k/\beta_1}, \dots, 2^{-k/\beta_d}))]^q \right)^{1/q} & (q < \infty), \\ \sup_{k \geq 0} 2^k w_{r,p}(f, (2^{-k/\beta_1}, \dots, 2^{-k/\beta_d})) & (q = \infty). \end{cases}$$

99 The norm of the anisotropic Besov space $B_{p,q}^\beta(\Omega)$ is defined by $\|f\|_{B_{p,q}^\beta} := \|f\|_p + |f|_{B_{p,q}^\beta}$, and
 100 $B_{p,q}^\beta(\Omega) = \{f \in L^p(\Omega) \mid \|f\|_{B_{p,q}^\beta} < \infty\}$.

101 Roughly speaking β represents the smoothness in each direction. If β_i is large, then a function in
 102 $B_{p,q}^\beta$ is smooth to the i th coordinate direction, otherwise, it is non-smooth to that direction. p is also
 103 an important quantity that controls the *spatial inhomogeneity* of the smoothness. If $\beta_1 = \beta_2 = \dots =$
 104 β_d , then the definition is equivalent to the usual Besov space (DeVore & Popov, 1988; DeVore et al.,
 105 1993). Suzuki (2019) analyzed curse of dimensionality of deep learning through a so-called *mixed*
 106 *smooth Besov* (m-Besov) space which imposes a stronger condition toward all directions uniformly.
 107 Particularly, it imposes stronger smoothness toward non-coordinate axis directions. Moreover, m-
 108 Besov space does *not* include the vanilla Besov space as a special case and thus cannot capture the
 109 situation that we consider in this paper.

110 Throughout this paper, for given $\beta = (\beta_1, \dots, \beta_d)^\top \in \mathbb{R}_{++}^d$, we write $\underline{\beta} := \min_i \beta_i$ (smallest
 111 smoothness) and $\bar{\beta} := \max_i \beta_i$ (largest smoothness). The approximation error of a function in

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

112 anisotropic Besov spaces is characterized by the harmonic mean of $(\beta_j)_{j=1}^d$, which corresponds to
 113 the average smoothness, and thus we define

$$\tilde{\beta} := \left(\sum_{j=1}^d 1/\beta_j \right)^{-1}. \quad (1)$$

114 The Besov space is closely related to other function spaces such as Hölder space. Let $\partial^\alpha f(x) =$
 115 $\frac{\partial^{|\alpha|} f}{\partial^{\alpha_1} x_1 \dots \partial^{\alpha_d} x_d}(x)$.

116 **Definition 3** (Hölder space $(\mathcal{C}^\beta(\Omega))$). For a smoothness parameter $\beta \in \mathbb{R}_{++}$ with $\beta \notin \mathbb{N}$, con-
 117 sider an m times differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ where $m = \lfloor \beta \rfloor$ (the largest integer
 118 less than β), and let the norm of the Hölder space $\mathcal{C}^\beta(\Omega)$ be $\|f\|_{\mathcal{C}^\beta} := \max_{|\alpha| \leq m} \|\partial^\alpha f\|_\infty +$
 119 $\max_{|\alpha|=m} \sup_{x,y \in \Omega} \frac{|\partial^\alpha f(x) - \partial^\alpha f(y)|}{\|x-y\|^{\beta-m}}$. Then, (β) -Hölder space $\mathcal{C}^\beta(\Omega)$ is defined as $\mathcal{C}^\beta(\Omega) = \{f \mid$
 120 $\|f\|_{\mathcal{C}^\beta} < \infty\}$.

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

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

- 124 1. For $\beta = (\beta_0, \dots, \beta_0)^\top \in \mathbb{R}^d$ with $\beta_0 \notin \mathbb{N}$, it holds that $\mathcal{C}^{\beta_0}(\Omega) = B_{\infty, \infty}^{\beta_0}(\Omega)$.
- 125 2. For $0 < p_1, p_2, q \leq \infty$, $p_1 \leq p_2$ and $\beta \in \mathbb{R}_{++}^d$ with $\tilde{\beta} > (1/p_1 - 1/p_2)_+^3$, it holds that⁴
 126 $B_{p_1, q}^{\beta}(\Omega) \hookrightarrow B_{p_2, q}^{\beta}(\Omega)$ for $\gamma = 1 - (1/p_1 - 1/p_2)_+ / \tilde{\beta}$.
- 127 3. For $0 < p, q_1, q_2 \leq \infty$, $q_1 < q_2$, and $\beta \in \mathbb{R}_{++}^d$, it holds that $B_{p, q_1}^{\beta} \hookrightarrow B_{p, q_2}^{\beta}$. In particular,
 128 with properties 1 and 2, if $\tilde{\beta} > 1/p$, it holds that $B_{p, q}^{\beta}(\Omega) \hookrightarrow \mathcal{C}^{\gamma \tilde{\beta}}(\Omega)$ where $\gamma = 1 - 1/(\tilde{\beta} p)$.
- 129 4. For $0 < p, q \leq \infty$ and $\beta \in \mathbb{R}_{++}^d$, if $\tilde{\beta} > 1/p$, then $B_{p, q}^{\beta}(\Omega) \hookrightarrow \mathcal{C}^0(\Omega)$.

130 This result is basically proven by Triebel (2011). For completeness, we provide its derivation in
 131 Appendix D. If the average smoothness $\tilde{\beta}$ is sufficiently large ($\tilde{\beta} > 1/p$), then the functions in
 132 $B_{p, q}^{\beta}$ are continuous; however, if it is small ($\tilde{\beta} < 1/p$), then they are no longer continuous. Small
 133 p indicates spatially inhomogeneous smoothness; thus, spikes and jumps appear (see Donoho &
 134 Johnstone (1998) for this perspective, from the viewpoint of wavelet analysis).

135 2.2 Model of the true function

136 As a model of the true function f° , we consider two types of models: *Affine composition model* and
 137 *deep composition model*. For a Banach space \mathcal{H} , we let $U(\mathcal{H})$ be the unit ball of \mathcal{H} .

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

$$\mathcal{H}_{\text{aff}} := \{h(Ax + b) \mid h \in U(B_{p, q}^{\beta}([0, 1]^{\tilde{d}})), A \in \mathbb{R}^{\tilde{d} \times d}, b \in \mathbb{R}^b \text{ s.t. } Ax + b \in [0, 1]^{\tilde{d}} (\forall x \in \Omega)\},$$

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

144 **(b) Deep composition model:** The *deep composition model* generalizes the affine composition
 145 model to a composition of nonlinear functions. Let $m_1 = d, m_{L+1} = 1, m_\ell$ be the dimension of the
 146 ℓ th layer, and let $\beta^{(\ell)} \in \mathbb{R}_{++}^{m_\ell}$ be the smoothness parameter in the ℓ th layer. The deep composition
 147 model is defined as

$$\mathcal{H}_{\text{deep}} := \{h_H \circ \dots \circ h_1(x) \mid h_\ell : [0, 1]^{m_\ell} \rightarrow [0, 1]^{m_{\ell+1}}, h_{\ell, k} \in U(B_{p, q}^{\beta^{(\ell)}}([0, 1]^{m_\ell})) (\forall k \in [m_{\ell+1}])\}.$$

148 Here, the interval $[0, 1]$ can be replaced by another compact interval, such as $[a_\ell, b_\ell]$, but this dif-
 149 ference can be absorbed by changing a scaling factor. The assumption $\|h_{\ell, k}\|_{B_{p, q}^{\beta^{(\ell)}}} \leq 1$ can also be
 150 relaxed, but we do not pursue that direction due to presentation simplicity. This model includes the
 151 affine composition model as a special case. However, it requires a stronger assumption to properly
 152 evaluate the estimation error on this model.

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

⁴The symbol \hookrightarrow means continuous embedding.

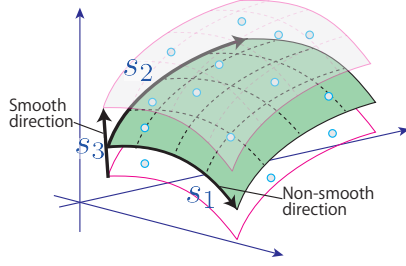


Figure 1: Near low dimensional data distribution with anisotropic smoothness of the target function. The target function has less smoothness (s_1, s_2) toward the first two coordinates on the manifold while it is almost constant toward the third coordinate (large s_3).

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

154 **(a) Linear projection** Schmidt-Hieber (2018) analyzed estimation of the following model by deep
 155 learning: $f^\circ(x) = g(w^\top x)$ where $g \in \mathcal{C}^\beta([0, 1])$ and $w \in \mathbb{R}^d$. In this example, the function f°
 156 varies along only one direction, w . Apparently, this is an example of the affine composition model.

157 **(b) Distribution on low dimensional smooth manifold** Assume that the input x is distributed
 158 on a low-dimensional smooth manifold embedded in Ω , and the smoothness of the true function
 159 f° is anisotropic along a coordinate direction on the manifold. We suppose that the low dimen-
 160 sional manifold is \tilde{d} -dimensional and $\tilde{d} \ll d$. In this situation, the true function can be written as
 161 $f^\circ(x) = h(\phi(x))$ where $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^{\tilde{d}}$ is a map that returns the coordinate of x on the manifold and
 162 h is an element in an anisotropic Besov space on $\mathbb{R}^{\tilde{d}}$. This situation appears if data is distributed
 163 on a low-dimensional sub-manifold of Ω and the target function is invariant against noise injection
 164 to some direction on the manifold at each input point x (Figure 1 illustrates this situation). One
 165 typical example of this situation is a function invariant with data augmentation (Simard et al., 2003;
 166 Krizhevsky et al., 2012). Even if the noise injection destroys low dimensionality of the data distribu-
 167 tion (i.e., $\tilde{d} = d$), an anisotropic smoothness of the target function eases the curse of dimensionality
 168 as analyzed below, which is quite different from existing works (Yang & Dunson, 2016; Bickel &
 169 Li, 2007; Nakada & Imaizumi, 2020; Schmidt-Hieber, 2019; Chen et al., 2019; Chen et al., 2019).

170 **Related work** Here, we introduce some more related work and discuss their relation to our analy-
 171 sis. The statistical analysis on an anisotropic Besov space can be back to Ibragimov & Khas'minskii
 172 (1984) who considered density estimation, where the density is assumed to be included in an
 173 anisotropic Sobolev space with $p \geq 2$, and derived the minimax optimal rate $n^{-r\tilde{\beta}/(2\tilde{\beta}+1)}$ with
 174 respect to L^r -norm. Nyssbaum (1983, 1987) analyzed a nonparametric regression problem on an
 175 anisotropic Besov space. Following these results, several studies have been conducted in the liter-
 176 ature pertaining to nonparametric statistics, such as nonlinear kernel estimator Kerkyacharian et al.
 177 (2001), adaptive confidence band construction Hoffman & Lepski (2002), optimal aggregation Gaif-
 178 fas & Lecue (2011), Gaussian process estimator Bhattacharya et al. (2011, 2014), and kernel ridge
 179 regression Hang & Steinwart (2018). Basically, these studies investigated estimation problems in
 180 which the target function is in anisotropic Besov spaces, but the composition models considered in
 181 this paper have not been analyzed. Hoffman & Lepski (2002); Bhattacharya et al. (2011) consid-
 182 ered a dimension reduction model; that is, the target function is dependent on only a few variables
 183 of x , but they did not deal with more general models, such as the affine/deep composition models.
 184 The nonparametric regression problems where the input data are distributed on a low-dimensional
 185 smooth manifold has been studied as a ‘‘manifold regression’’ Yang & Dunson (2016); Bickel & Li
 186 (2007); Yang & Tokdar (2015). Such a model can be considered as a specific example of the deep
 187 composition model. In this sense, our analysis is a significant extension of these analyses.

188 3 Approximation error analysis

Here, we consider approximating the true function f° via a deep neural network and derive the approximation error. As the activation function, we consider the ReLU activation denoted by $\eta(x) = \max\{x, 0\}$ ($x \in \mathbb{R}$). For a vector x , $\eta(x)$ 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 $\Phi(L, W, S, B) := \{(\mathcal{W}^{(L)}\eta(\cdot) + b^{(L)}) \circ \dots \circ (\mathcal{W}^{(1)}x + b^{(1)}) \mid \mathcal{W}^{(L)} \in \mathbb{R}^{1 \times W}, b^{(L)} \in \mathbb{R}, \mathcal{W}^{(1)} \in \mathbb{R}^{W \times d}, b^{(1)} \in \mathbb{R}^W, \mathcal{W}^{(\ell)} \in \mathbb{R}^{W \times W}, b^{(\ell)} \in \mathbb{R}^W (1 < \ell < L), \sum_{\ell=1}^L (\|\mathcal{W}^{(\ell)}\|_0 + \|b^{(\ell)}\|_0) \leq$

$S, \max_{\ell} \|\mathcal{W}^{(\ell)}\|_{\infty} \vee \|b^{(\ell)}\|_{\infty} \leq B\}$, where $\|\cdot\|_0$ is the ℓ_0 -norm of the matrix (the number of non-zero elements of the matrix), and $\|\cdot\|_{\infty}$ is the ℓ_{∞} -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

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

189 where \mathcal{F} is the set of functions used for approximation, and \mathcal{H} is the set of target functions.

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

$$R_r(\Phi(L_1, W_1, S_1, B_1), U(B_{p,q}^{\beta}(\Omega))) \lesssim N^{-\tilde{\beta}},$$

190 by setting

$$L_1(d) := 3 + 2 \lceil \log_2 \left(\frac{3^{d \vee m}}{\epsilon c_{(d,m)}} \right) + 5 \rceil \lceil \log_2(d \vee m) \rceil, \quad W_1(d) := N W_0, \quad (2)$$

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

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

192 The proof of this proposition is provided in Appendix B. The rate $N^{-\tilde{\beta}}$ is the optimal *adaptive*
 193 approximation error rate that can be achieved by a model with N parameters (the difference be-
 194 tween adaptive and non-adaptive methods is explained in the discussion below). Note that this is
 195 an approximation error in an oracle setting and no sample complexity appears here. We notice that
 196 we can avoid the *curse of dimensionality* if the average smoothness $\tilde{\beta}$ is small. This means that
 197 if the target function is non-smooth in only a few directions and smooth in other directions, we
 198 can avoid the curse of dimensionality. In contrast, if we consider an isotropic Besov space where
 199 $\beta_1 = \dots = \beta_d (= \beta)$, then $\tilde{\beta} = d\beta$, which directly depends on the dimensionality d , and we need
 200 an exponentially large number of parameters in this situation to achieve ϵ -accuracy. Therefore, the
 201 anisotropic smoothness has a significant impact on the approximation error rate. The assumption
 202 $\tilde{\beta} > (1/p - 1/r)_+$ ensures the L_r -integrability of the target function, and the inequality (without
 203 equality) admits a near-optimal wavelet approximation of the target function in terms of L_r -norm.

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

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

$$R_{\infty}(\Phi(L, W, S, B), \mathcal{H}_{\text{deep}}) \lesssim \max_{\ell \in [H]} N^{-\tilde{\beta}^{*(\ell)}}, \quad (4)$$

209 where $\tilde{\beta}^{*(\ell)} = \tilde{\beta}^{(\ell)} \prod_{k=\ell+1}^H [(\beta^{(k)} - 1/p) \wedge 1]$, and $L = \sum_{\ell=1}^H (L_1(m_{\ell}) + 1)$, $W = \max_{\ell} (W_1(m_{\ell}) \vee$
 210 $m_{\ell+1})$, $S = \sum_{\ell=1}^H (S_1(m_{\ell}) + 3m_{\ell+1})$, $B = \max_{\ell} B_1(m_{\ell})$.

211 The proof can be found in Appendix B.1. Since the model is more general than the vanilla
 212 anisotropic Besov space, we require a stronger assumption $\tilde{\beta}^{(\ell)} > 1/p$ on $\tilde{\beta}^{(\ell)}$ than the condi-
 213 tion in Proposition 2. This is because we need to bound the Hölder smoothness of the remaining
 214 layers to bound the influence of the approximation error in the internal layers to the entire function.
 215 Hölder smoothness is ensured according to the embedding property under this condition (Proposi-
 216 tion 1). This Hölder smoothness assumption affects the approximation error rate. The convergence
 217 rate $\tilde{\beta}^{*(\ell)}$ in Eq. (4) is different from that in Eq. (8). This is because the approximation error
 218 in the internal layers are propagated through the remaining layers with Hölder smoothness and its
 219 amplitude is controlled by the Hölder smoothness.

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

227 4 Estimation error analysis

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

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

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

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

242 **Theorem 2.** *Assume the same condition as in Theorem 6; in particular, suppose $0 < p, q \leq \infty$
 243 and $\tilde{\beta} > (1/p - 1/2)_+$ for $\tilde{\beta} \in \mathbb{R}_{++}^{\tilde{d}}$. Moreover, we assume that the distribution P_X has a density
 244 p_X such that $\|p_X\|_\infty \leq R$ for a constant $R > 0$. If $f^\circ \in \mathcal{H}_{\text{aff}} \cap L^\infty(\Omega)$, and $\|f^\circ\|_\infty \leq F$ for
 245 $F \geq 1$; then, letting $(W, L, S, B) = (L_1(\tilde{d}), W_1(\tilde{d}), S_1(\tilde{d}), (\tilde{d}C + 1)B_1(\tilde{d}))$ as in Theorem 6 with
 246 $N \asymp n^{\frac{1}{2\tilde{\beta}+1}}$, we obtain*

$$\mathbb{E}_{D_n} [\|f^\circ - \hat{f}\|_{L^2(P_X)}^2] \lesssim n^{-\frac{2\tilde{\beta}}{2\tilde{\beta}+1}} \log(n)^3,$$

247 where $\mathbb{E}_{D_n}[\cdot]$ indicates the expectation with respect to the training data D_n .

248 The proof is provided in Appendix C. We will show that the convergence rate $n^{-2\tilde{\beta}/(2\tilde{\beta}+1)}$ is min-
 249 imax optimal in Section 5 (see also Kerkycharian & Picard (1992); Donoho et al. (1996); Donoho
 250 & Johnstone (1998); Giné & Nickl (2015) for ordinary Besov spaces). The L^∞ -norm constraint
 251 $\|f^\circ\|_\infty \leq F$ is used to derive a uniform bound on the discrepancy between the population and the
 252 empirical L^2 -norm. Without this condition, the convergence rate could be slower.

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

255 **Theorem 3.** *Suppose that $0 < p, q \leq \infty$ and $\tilde{\beta}^{(\ell)} > 1/p$ for all $\ell \in [H]$. If $f^\circ \in \mathcal{H}_{\text{deep}} \cap L^\infty(\Omega)$,
 256 and $\|f\|_\infty \leq F$ for $F \geq 1$, then we obtain*

$$\mathbb{E}_{D_n} [\|f^\circ - \hat{f}\|_{L^2(P_X)}^2] \lesssim \max_{\ell \in [H]} n^{-2\tilde{\beta}^{(\ell)}/(2\tilde{\beta}^{(\ell)}+1)} \log(n)^3,$$

257 where $\tilde{\beta}^{(\ell)}$ is defined in Theorem 1, and (L, W, S, B) is as given in Theorem 1 with $N \asymp$
 258 $\max_{\ell \in [L]} n^{\frac{1}{2\tilde{\beta}^{(\ell)}+1}}$.

259 The proof is provided in Appendix C. We will show that this is also minimax optimal in Theorem 4.
 260 Because of the Hölder continuity, the convergence rate becomes slower than the affine composition
 261 model (that is, $\tilde{\beta}^{(\ell)} \leq \tilde{\beta}^{(\ell)}$). However, this slower rate is unavoidable in terms of the minimax
 262 optimal rate. Schmidt-Hieber (2018) analyzed the same situation for the Hölder class which corre-
 263 sponds to $\beta_1^{(\ell)} = \dots = \beta_d^{(\ell)}$ ($\forall \ell$) and $p = q = \infty$. Our analysis far extends their analysis to the
 264 setting of anisotropic Besov spaces in which the parameters $\beta^{(\ell)}, p, q$ have much more freedom.

265 From these two bounds (Theorems 2 and 3), we can see that as the smoothness $\tilde{\beta}$ becomes large,
 266 the convergence rates faster. If the target function is included in the isotropic Besov space with
 267 smoothness $\beta_1 = \dots = \beta_d (= \underline{\beta})$, then the estimation error becomes

$$\text{(Isotropic Besov)} \quad n^{-2\underline{\beta}/(2\underline{\beta}+d)}.$$

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

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

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

276 5 Minimax optimal rate

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

$$R_*(\mathcal{F}^\circ) := \inf_{\hat{f}} \sup_{f^\circ \in \mathcal{F}^\circ} \mathbb{E}_{D_n} [\|\hat{f} - f^\circ\|_{L^2(P_X)}^2],$$

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

282 **Theorem 4. (a) Affine composition model:** For $0 < p, q \leq \infty$ and $\beta \in \mathbb{R}_{++}^d$, assume that
 283 $\tilde{\beta} > \max\{1/p - 1/2, 1/p - 1, 0\}$. Then, the minimax optimal risk of the affine composition model
 284 is lower bounded as $R_*(\mathcal{H}_{\text{aff}}) \gtrsim n^{-\frac{2\tilde{\beta}}{2\tilde{\beta}+1}}$. **(b) Deep composition model:** For $0 < p, q \leq \infty$ and
 285 $\beta^{(\ell)} \in \mathbb{R}_{++}^d$ ($\ell = 1, \dots, H$), assume that $\tilde{\beta}^{(\ell)} > 1/p$. Let $\epsilon > 0$ be arbitrarily small for $q < \infty$, and
 286 let $\epsilon = 0$ for $q = 0$. Let $\tilde{\beta}^{*(\ell)} = \tilde{\beta}^{(\ell)} \prod_{k=\ell+1}^H [(\underline{\beta}^{(k)} - 1/p + \epsilon) \wedge 1]$, and $\tilde{\beta}^{**} := \min_{\ell} \tilde{\beta}^{*(\ell)}$. Then, the
 287 minimax optimal risk of the deep composition model is lower bounded as $R_*(\mathcal{H}_{\text{deep}}) \gtrsim n^{-\frac{2\tilde{\beta}^{**}}{2\tilde{\beta}^{**}+1}}$.

288 The proof is provided in Appendix E (see also Ibragimov & Khas'minskii (1984); Nyssbaum
 289 (1987)). From this theorem, we can see that the estimation error of deep learning shown in The-
 290 orems 2 and 3 indeed achieve the minimax optimal rate up to a poly-log(n) factor.

291 6 Suboptimality of linear estimators

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

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

where $X^n = (x_1, \dots, x_n)$ and $\varphi_i(x; X^n)$ ($i = 1, \dots, n$) 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 $\hat{f}(x) = k_{x, X^n} (k_{X^n, X^n} + \lambda \mathbf{I})^{-1} Y^n$, 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:

$$R_*^{(\text{lin})}(\mathcal{F}^\circ) := \inf_{\hat{f}: \text{linear}} \sup_{f^\circ \in \mathcal{F}^\circ} \mathbb{E}_{D_n} [\|f^\circ - \hat{f}\|_{L^2(P_X)}^2],$$

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

297 **Theorem 5.** (i) Suppose that the input distribution P_X is uniform distribution on $\Omega = [0, 1]^d$ and
 298 $\tilde{\beta} > 1/p$, $1 \leq p, q \leq \infty$. Then, the minimax optimal rate of the linear estimators is lower bounded
 299 as

$$R_*^{(\text{lin})}(U(B_{p,q}^\beta)) \gtrsim n^{-\frac{2\tilde{\beta}-v}{2\tilde{\beta}+1-v}}, \quad (6)$$

300 where $v = 2(1/p - 1/2)_+$. (ii) In addition to the above conditions, we assume that $\tilde{d} \leq d$, $\underline{\beta} =$
 301 $\beta_1 = \dots = \beta_{\tilde{d}}$ and $0 < p \leq 2$. Let $a_d = 1$ when $\tilde{d} < d/2$ and $a_d = 0$ when $\tilde{d} \geq d/2$. Then, the
 302 minimax rate of the linear estimators on the affine composition model is lower bounded by

$$R_*^{(\text{lin})}(\mathcal{H}_{\text{aff}}) \gtrsim n^{-\frac{2(\underline{\beta}-\tilde{d}/p+d/2+a_d)}{2(\underline{\beta}-\tilde{d}/p+d/2+a_d)+d}}. \quad (7)$$

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 $\tilde{d} = 1$, $p = 1$ and $d \gg \tilde{d}$, then the obtained minimax rate of linear estimators and the estimation error rate of deep learning can be summarized as

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

303 by Theorem 2 when $\beta > 1$ (which can be checked by noticing $\tilde{d} = \underline{\beta}/\tilde{\beta} = 1$ in this situation). We
 304 can see that the dependence on the dimensionality of linear estimators is significantly worse than that
 305 of deep learning. This indicates poor adaptivity of linear estimators to the intrinsic dimensionality
 306 of data. Actually, as d becomes large, the rate for the linear estimator approaches to $1/\sqrt{n}$ but
 307 that for the deep learning is not affected by d and still faster than $1/\sqrt{n}$. To show the theorem, we
 308 used the ‘‘convex-hull argument’’ developed by Hayakawa & Suzuki (2019); Donoho & Johnstone
 309 (1998). We combined this technique with the so-called Irie-Miyake’s integral representation (Irie
 310 & Miyake, 1988; Hornik et al., 1990). Note that this difference appears because there is an affine
 311 transformation in the first layer of the affine composition model. Deep learning is flexible against
 312 such a coordinate transform so that it can find directions to which the target function is smooth. In
 313 contrast, kernel methods do not have such adaptivity because there is no feature extraction layer.
 314 (ii) The lower bound (6) states that when $p < 2$ (that is, $v > 0$), the minimax rate of the linear
 315 estimators is outperformed by that of deep learning (Theorem 2). This is due to the ‘‘adaptivity’’ of
 316 deep learning. When p is small, the smoothness of the target function is less homogeneous, and it
 317 requires an adaptive approximation scheme to achieve the best estimation error. Linear estimators
 318 do not have adaptivity and thus fail to achieve the minimax optimal rate. Our bound (6) extends
 319 the result by Zhang et al. (2002) to a multivariate anisotropic Besov space while Zhang et al. (2002)
 320 investigated the univariate space ($d = 1$).

321 7 Conclusion

322 We investigated the approximation error and estimation error of deep learning in the anisotropic
 323 Besov spaces. It was proved that the convergence rate is determined by the average of the anisotropic
 324 smoothness, which results in milder dependence on the input dimensionality. If the smoothness is
 325 highly anisotropic, deep learning can avoid overfitting. We also compared the error rate of deep
 326 learning with that of linear estimators and showed that deep learning has better dependence on
 327 the input dimensionality. Moreover, it was shown that deep learning can achieve the adaptive rate
 328 and outperform non-adaptive approximation methods and linear estimators if the homogeneity p of
 329 smoothness is small. These analyses strongly support the practical success of deep learning from a
 330 theoretical perspective.

331 **Limitations of this work** Our work does not cover the optimization aspect of deep learning. It
 332 is assumed that the regularized least squares (5) can be executed. It would be nice to combine
 333 our study with recent developments of non-convex optimization techniques (Vempala & Wibisono,
 334 2019; Suzuki & Akiyama, 2021).

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

338 **References**

- 339 M. Belkin and P. Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation.
340 *Neural computation*, 15(6):1373–1396, 2003.
- 341 A. Bhattacharya, D. Pati, and D. B. Dunson. Adaptive dimension reduction with a gaussian process
342 prior. *arXiv preprint arXiv:1111.1044*, 1445, 2011.
- 343 A. Bhattacharya, D. Pati, and D. Dunson. Anisotropic function estimation using multi-bandwidth
344 gaussian processes. *Annals of statistics*, 42(1):352, 2014.
- 345 P. J. Bickel and B. Li. Local polynomial regression on unknown manifolds. In *Complex datasets*
346 *and inverse problems*, pp. 177–186. Institute of Mathematical Statistics, 2007.
- 347 M. Chen, H. Jiang, W. Liao, and T. Zhao. Nonparametric Regression on Low-Dimensional Mani-
348 folds using Deep ReLU Networks. *arXiv e-prints*, art. arXiv:1908.01842, Aug 2019.
- 349 M. Chen, H. Jiang, W. Liao, and T. Zhao. Efficient approximation of deep relu networks for functions
350 on low dimensional manifolds. In *Advances in Neural Information Processing Systems*, pp. 8172–
351 8182, 2019.
- 352 G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control,*
353 *Signals, and Systems (MCCS)*, 2(4):303–314, 1989.
- 354 J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of Deep Bidirectional
355 Transformers for Language Understanding. *arXiv e-prints*, art. arXiv:1810.04805, Oct 2018.
- 356 R. A. DeVore. Nonlinear approximation. *Acta Numerica*, 7:51–150, 1998.
- 357 R. A. DeVore and V. A. Popov. Interpolation of Besov spaces. *Transactions of the American*
358 *Mathematical Society*, 305(1):397–414, 1988.
- 359 R. A. DeVore, G. Kyriazis, D. Leviatan, and V. M. Tikhomirov. Wavelet compression and
360 nonlinear-widths. *Advances in Computational Mathematics*, 1(2):197–214, 1993.
- 361 D. L. Donoho and I. M. Johnstone. Minimax estimation via wavelet shrinkage. *The Annals of*
362 *Statistics*, 26(3):879–921, 1998.
- 363 D. L. Donoho, I. M. Johnstone, G. Kerkyacharian, and D. Picard. Density estimation by wavelet
364 thresholding. *The Annals of Statistics*, 24(2):508–539, 1996.
- 365 D. Dũng. Optimal adaptive sampling recovery. *Advances in Computational Mathematics*, 34(1):
366 1–41, 2011.
- 367 S. Gaiffas and G. Lecue. Hyper-sparse optimal aggregation. *Journal of Machine Learning Research*,
368 12(Jun):1813–1833, 2011.
- 369 E. Giné and R. Nickl. *Mathematical Foundations of Infinite-Dimensional Statistical Models*. Cam-
370 bridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2015.
- 371 X. Glorot, A. Bordes, and Y. Bengio. Deep sparse rectifier neural networks. In *Proceedings of the*
372 *14th International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings*
373 *of Machine Learning Research*, pp. 315–323, 2011.
- 374 H. Hang and I. Steinwart. Optimal learning with anisotropic gaussian svms. *arXiv preprint*
375 *arXiv:1810.02321*, 2018.
- 376 S. Hayakawa and T. Suzuki. On the minimax optimality and superiority of deep neural network
377 learning over sparse parameter spaces. *arXiv preprint arXiv:1905.09195*, 2019.
- 378 M. Hoffman and O. Lepski. Random rates in anisotropic regression (with a discussion and a rejoinder
379 by the authors). *The Annals of Statistics*, 30(2):325–396, 04 2002.
- 380 K. Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2):
381 251–257, 1991.

- 382 K. Hornik, M. Stinchcombe, and H. White. Universal approximation of an unknown mapping and
383 its derivatives using multilayer feedforward networks. *Neural Networks*, 3(5):551–560, 1990.
- 384 I. Ibragimov and R. Khas'minskii. More on the estimation of distribution densities. *Journal of Soviet*
385 *Mathematics*, 25(3):1155–1165, 1984.
- 386 B. Irie and S. Miyake. Capabilities of three-layered perceptrons. In *IEEE 1988 International Con-*
387 *ference on Neural Networks*, pp. 641–648, 1988.
- 388 G. Kerkyacharian and D. Picard. Density estimation in Besov spaces. *Statistics & Probability*
389 *Letters*, 13:15–24, 1992.
- 390 G. Kerkyacharian, O. Lepski, and D. Picard. Nonlinear estimation in anisotropic multi-index de-
391 noising. *Probability Theory and Related Fields*, 121(2):137–170, Oct 2001.
- 392 A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional
393 neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- 394 V. Myronyuk. Trigonometric approximations and kolmogorov widths of anisotropic besov classes
395 of periodic functions of several variables. *Ukrainian Mathematical Journal*, 66(8), 2015.
- 396 V. V. Myronyuk. Kolmogorov widths of the anisotropic besov classes of periodic functions of many
397 variables. *Ukrainian Mathematical Journal*, 68(5):718–727, Oct 2016.
- 398 V. V. Myronyuk. Widths of the anisotropic besov classes of periodic functions of several variables.
399 *Ukrainian Mathematical Journal*, 68(8):1238–1251, Jan 2017. ISSN 1573-9376. doi: 10.1007/
400 s11253-017-1290-1. URL <https://doi.org/10.1007/s11253-017-1290-1>.
- 401 V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceed-*
402 *ings of the 27th International Conference on Machine Learning*, pp. 807–814, 2010.
- 403 R. Nakada and M. Imaizumi. Adaptive approximation and generalization of deep neural network
404 with intrinsic dimensionality. *Journal of Machine Learning Research*, 21(174):1–38, 2020. URL
405 <http://jmlr.org/papers/v21/20-002.html>.
- 406 S. M. Nikol'skii. *Approximation of functions of several variables and imbedding theorems*, volume
407 205. Springer-Verlag Berlin Heidelberg, 1975.
- 408 M. Nyssbaum. Optimal filtration of a function of many variables in white gaussian noise. *Problems*
409 *of Information Transmission*, 19:23–29, 1983.
- 410 M. Nyssbaum. Nonparametric estimation of a regression function that is smooth in a domain in \mathbb{R}^k .
411 *Theory of Probability & Its Applications*, 31(1):108–115, 1987.
- 412 A. Radford, L. Metz, and S. Chintala. Unsupervised Representation Learning with Deep Convoluti-
413 onal Generative Adversarial Networks. *arXiv e-prints*, art. arXiv:1511.06434, Nov 2015.
- 414 J. Schmidt-Hieber. Nonparametric regression using deep neural networks with ReLU activation
415 function. *ArXiv preprint arXiv:1708.06633(v3)*, 2018.
- 416 J. Schmidt-Hieber. Deep ReLU network approximation of functions on a manifold. *arXiv preprint*
417 *arXiv:1908.00695*, 2019.
- 418 P. Y. Simard, D. Steinkraus, and J. C. Platt. Best practices for convolutional neural networks applied
419 to visual document analysis. In *Proceedings of the Seventh International Conference on Document*
420 *Analysis and Recognition-Volume 2*, pp. 958. IEEE Computer Society, 2003.
- 421 S. Sonoda and N. Murata. Neural network with unbounded activation functions is universal approx-
422 imator. *Applied and Computational Harmonic Analysis*, 43(2):233–268, 2017.
- 423 T. Suzuki. Adaptivity of deep ReLU network for learning in Besov and mixed smooth Besov spaces:
424 optimal rate and curse of dimensionality. In *International Conference on Learning Representa-*
425 *tions (ICLR2019)*, 2019. URL <https://openreview.net/forum?id=H1ebTsActm>.

426 T. Suzuki and S. Akiyama. Benefit of deep learning with non-convex noisy gradient descent: Prov-
427 able excess risk bound and superiority to kernel methods. In *International Conference on Learn-*
428 *ing Representations*, 2021. URL <https://openreview.net/forum?id=2m0g1wEafh>.

429 V. Temlyakov. *Approximation of Periodic Functions*. Nova Science Publishers, 1993.

430 J. B. Tenenbaum, V. De Silva, and J. C. Langford. A global geometric framework for nonlinear
431 dimensionality reduction. *science*, 290(5500):2319–2323, 2000.

432 H. Triebel. Entropy numbers in function spaces with mixed integrability. *Revista matemática com-*
433 *plutense*, 24(1):169–188, 2011.

434 S. Vempala and A. Wibisono. Rapid convergence of the unadjusted langevin algorithm: Isoperimetry
435 suffices. In *Advances in Neural Information Processing Systems*, pp. 8094–8106, 2019.

436 J. Vybiral. Function spaces with dominating mixed smoothness. *Dissertationes Math. (Rozprawy*
437 *Mat.)*, 436:3–73, 2006.

438 Y. Yang and D. B. Dunson. Bayesian manifold regression. *The Annals of Statistics*, 44(2):876–905,
439 2016.

440 Y. Yang and S. T. Tokdar. Minimax-optimal nonparametric regression in high dimensions. *The*
441 *Annals of Statistics*, 43(2):652–674, 2015.

442 D. Yarotsky. Error bounds for approximations with deep relu networks. *Neural Networks*, 94:
443 103–114, 2017.

444 S. Zhang, M.-Y. Wong, and Z. Zheng. Wavelet threshold estimation of a regression function with
445 random design. *Journal of Multivariate Analysis*, 80(2):256–284, 2002.

446 Checklist

- 447 1. For all authors...
- 448 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
449 contributions and scope? [Yes]
- 450 (b) Did you describe the limitations of your work? [Yes]
- 451 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 452 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
453 them? [Yes]
- 454 2. If you are including theoretical results...
- 455 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Each theorem
456 explicitly describes assumptions on which the theorem is established.
- 457 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix in the
458 supplementary material. All proofs are included there.
- 459 3. If you ran experiments...
- 460 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
461 mental results (either in the supplemental material or as a URL)? [N/A]
- 462 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
463 were chosen)? [N/A]
- 464 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
465 ments multiple times)? [N/A]
- 466 (d) Did you include the total amount of compute and the type of resources used (e.g., type
467 of GPUs, internal cluster, or cloud provider)? [N/A]
- 468 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 469 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 470 (b) Did you mention the license of the assets? [N/A]

- 471 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
472
- 473 (d) Did you discuss whether and how consent was obtained from people whose data
474 you're using/curating? [N/A]
- 475 (e) Did you discuss whether the data you are using/curating contains personally identifi-
476 able information or offensive content? [N/A]
- 477 5. If you used crowdsourcing or conducted research with human subjects...
- 478 (a) Did you include the full text of instructions given to participants and screenshots, if
479 applicable? [N/A]
- 480 (b) Did you describe any potential participant risks, with links to Institutional Review
481 Board (IRB) approvals, if applicable? [N/A]
- 482 (c) Did you include the estimated hourly wage paid to participants and the total amount
483 spent on participant compensation? [N/A]