
Deep Learning meets Nonparametric Regression: Are Weight-Decayed DNNs Locally Adaptive?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the theory of neural network (NN) from the lens of classical nonpara-
2 metric regression problems with a focus on NN’s ability to *adaptively* estimate
3 functions with *heterogeneous smoothness* — a property of functions in Besov or
4 Bounded Variation (BV) classes. Existing work on this problem requires tuning
5 the NN architecture based on the function spaces and sample sizes. We consider a
6 “Parallel NN” variant of deep ReLU networks and show that the standard weight
7 decay is equivalent to promoting the ℓ_p -sparsity ($0 < p < 1$) of the coefficient
8 vector of an end-to-end learned function bases, i.e., a dictionary. Using this equiv-
9 alence, we further establish that by tuning only the weight decay, such Parallel
10 NN achieves an estimation error arbitrarily close to the minimax rates for both the
11 Besov and BV classes. Notably, it gets exponentially closer to minimax optimal
12 as the NN gets deeper. Our research sheds new lights on why depth matters and
13 how NNs are more powerful than kernel methods.

14 1 Introduction

15 *Why* do deep neural networks (DNNs) *work* better? They are universal function approximators [6],
16 but so are splines and kernels. They learn data-driven representations, but so are the shallower and
17 linear counterparts such as matrix factorization. There is surprisingly little theoretical understanding
18 on why DNNs are superior to these classical alternatives.

19 In this paper, we study DNNs in nonparametric regression problems — a classical branch of statis-
20 tical theory and methods with more than half a century of associated literature [25, 7, 46, 10, 23, 37,
21 33]. Nonparametric regression addresses the following fundamental problem:

- 22 • Let $y_i = f(x_i) + \text{Noise}$ for $i = 1, \dots, n$. How can we estimate a function f using data points
23 $(x_1, y_1), \dots, (x_n, y_n)$ in conjunction with the knowledge that f belongs to a function class \mathcal{F} ?

24 Function class \mathcal{F} typically imposes only weak regularity assumptions such as smoothness, which
25 makes nonparametric regression widely applicable to real-life applications under weak assumptions.

26 **Local adaptivity.** A subset of nonparametric regression techniques were shown to have the property
27 of *local adaptivity* [24] in both theory and practice. These include wavelet smoothing [10], locally
28 adaptive regression splines [24], trend filtering [40, 47] and adaptive local polynomials [2, 3]. We
29 say a nonparametric regression technique is *locally adaptive* if it can cater to local differences in
30 smoothness, hence allowing more accurate estimation of functions with varying smoothness and
31 abrupt changes.

32 In light of such a distinction, it is natural to consider the following question.

33 Are NNs *locally adaptive*, i.e., optimal in learning functions with heterogeneous
34 smoothness?

35 This is a timely question to ask, partly because the bulk of recent theory of NN leverages its asymp-
36 totic Reproducing Kernel Hilbert Space (RKHS) in the overparameterized regime [21, 5, 1]. RKHS-
37 based approaches, e.g., kernel ridge regression with any fixed kernels are *suboptimal* in estimating
38 functions with heterogeneous smoothness [9]. Therefore, existing deep learning theory based on
39 RKHS does not satisfactorily explain the advantages of neural networks over kernel methods.

40 We build upon the recent work of Suzuki [39] and Parhi and Nowak [29] who provided encouraging
41 first answers to the question above about the local adaptivity of NNs. Specifically, Parhi and Nowak
42 [29, Theorem 8] showed that a two-layer *truncated power function* activated neural network with
43 a non-standard regularization is equivalent to the *locally adaptive regression splines* (LARS) [24].
44 This connection implies that such non-standard NNs achieve the minimax rate for the (higher order)
45 bounded variation (BV) classes. We provide a detailed discussion about this work in Section B.
46 Suzuki [39] showed that multilayer ReLU DNNs can achieve minimax rate for the Besov class,
47 but requires the width, depth and an artificially imposed sparsity-level of the DNN weights to be
48 carefully calibrated according to parameters of the Besov class, thus is quite different from how
49 DNNs are typically trained in practice.

50 In this paper, we aim at addressing the same *locally adaptivity* question for a more commonly used
51 neural network with standard weight decayed training.

52 **Parallel neural networks.** We restrict our attention on a special network architecture called *parallel*
53 *neural network* [18, 15] which learns an ensemble of subnetworks — each being a multilayer ReLU
54 DNNs. Parallel NNs have been shown to be more well-behaved both theoretically [18, 51, 16, 15, 14]
55 and empirically [50, 44]. Moreover, the idea of parallel NNs was used in many successful NN
56 architectures such as SqueezeNet, ResNext and Inception (see [15] and the references therein).

57 **Weight decay.** Weight decay is a common method in deep learning to reduce overfitting. Em-
58 pirically, the regularizer is not necessarily explicit. Many tricks in deep learning, including early
59 stopping [48], quantization [20], and dropout [45] have similar effect as weight decay. In this paper,
60 we make no assumption on the training method thus there is no (implicit) regularizers apart from
61 weight decay.

62 **Summary of results.** Our main contributions are:

- 63 1. We prove that the (standard) weight decay in training an L -layer *parallel* ReLU-activated
64 neural network is equivalent to a sparse ℓ_p penalty term (where $p = 2/L$) on the linear
65 coefficients of a learned representation.
- 66 2. We show that neural networks can approximate B-spline basis functions of any order with-
67 out the need of choosing the order parameter manually. In other words, neural networks
68 can adapt to functions of different order of smoothness, and even functions with different
69 smoothness in different regions in their domain.
- 70 3. We show that the estimation error of weight decayed parallel ReLU neural network de-
71 creases polynomially with the number of samples up to a constant error for estimating
72 functions with heterogeneous smoothness in the both BV and Besov classes, and the ex-
73 ponential term in the error rate is close to the minimax rate. Notably, the method requires
74 tuning only the weight decay parameter.
- 75 4. We find that deeper models achieve closer to the optimal error rate. This result helps explain
76 why deep neural networks can achieve better performance than shallow ones empirically.

77 The above results separate NNs with any linear methods such as kernel ridge regression. To the
78 best of our knowledge, we are the first to demonstrate that standard techniques (“weight decay”
79 and ReLU activation) suffice for DNNs in achieving the optimal rates for estimating BV and Besov
80 functions.

81 2 Preliminary

82 2.1 Notation and Problem Setup.

83 We denote regular font letters as scalars, bold lower case letters as vectors and bold upper case letters
84 as matrices. $a \lesssim b$ means $a \leq Cb$ for some constant C that does not depend on a or b , and $a \approx b$
85 denotes $a \lesssim b$ and $b \lesssim a$. See Table 1 for the full list of symbols used.

Table 1: Symbols used in this paper

symbol	Meaning		
$a/\mathbf{a}/\mathbf{A}$	scalars / vectors / matrices.	$[a, b]$	$\{x \in \mathbb{R} : a \leq x \leq b\}$
$B_{p,q}^\alpha$	Besov space.	$[n]$	$\{x \in \mathbb{N} : 1 \leq x \leq n\}$.
$ \cdot _{B_{p,q}^\alpha}$	Besov quasi-norm .	$\ \cdot\ _F$	Frobenius norm.
$\ \cdot\ _{B_{p,q}^\alpha}$	Besov norm.	$\ \cdot\ _p$	ℓ_p -norm.
$M_m(\cdot)$	m^{th} order Cardinal B-spline bases.	d	Dimension of input.
$M_{m,k,\mathbf{s}}(\cdot)$	m^{th} order Cardinal B-spline basis function of resolution k at position \mathbf{s} .	M	# subnetworks in a parallel NN.
$\sigma(\cdot)$	ReLU activation function.	L	# layers in a (parallel) NN.
$\mathbf{W}_j^{(\ell)}, \mathbf{b}_j^{(\ell)}$	Weight and bias in the ℓ -th layer in the j -th subnetwork.	w	Width of a subnetwork.
		n	# samples.
		$\mathbb{R}, \mathbb{Z}, \mathbb{N}$	Set of real numbers, integers, and nonnegative integers.

86 Let f_0 be the target function to be estimated. The training dataset is $\mathcal{D}_n := \{(\mathbf{x}_i, y_i), y_i = f_0(\mathbf{x}_i) +$
87 $\epsilon_i, i \in [n]\}$, where x_i are fixed and ϵ_i are zero-mean, independent Gaussian noises with variance σ^2 .
88 In the following discussion, we assume $\mathbf{x}_i \in [0, 1]^d, f_0(\mathbf{x}_i) \in [-1, 1], \forall i$.

We will be comparing estimators under the mean square error (MSE), defined as

$$\text{MSE}(\hat{f}) := \mathbb{E}_{\mathcal{D}_n} \frac{1}{n} \sum_{i=1}^n (\hat{f}(\mathbf{x}_i) - f_0(\mathbf{x}_i))^2.$$

89 The optimal worst-case MSE is described by $R(\mathcal{F}) := \min_{\hat{f}} \max_{f_0 \in \mathcal{F}} \text{MSE}(\hat{f})$, we say that
90 \hat{f} is optimal if $\text{MSE}(\hat{f}) \approx R(\mathcal{F})$. The empirical (square error) loss is defined as $\hat{L}(\hat{f}) :=$
91 $\frac{1}{n} \sum_{i=1}^n (\hat{f}(\mathbf{x}_i) - y_i)^2$. The corresponding population loss is $L(\hat{f}) := \mathbb{E}[\frac{1}{n} \sum_{i=1}^n (\hat{f}(\mathbf{x}_i) - y'_i)^2 | \hat{f}]$
92 where y'_i are new data points. It is clear that $\mathbb{E}[L(\hat{f})] = \text{MSE}[\hat{f}] + \sigma^2$.

93 2.2 Besov Spaces and Bound Variation Space

94 **Besov space**, denoted as $B_{p,q}^\alpha$, is a flexible function class parameterized by α, p, q whose definition
95 is deferred to Section C.1. Here $\alpha \geq 0$ determines the smoothness of functions, $1 \leq p \leq \infty$
96 determines the averaging (quasi-)norm over locations, $1 \leq q \leq \infty$ determines the averaging (quasi-)
97 norm over scale which plays a relatively minor role. Smaller p is more forgiving to inhomogeneity
98 and loosely speaking, when the function domain is bounded, smaller p induces a larger function
99 space. On the other hand, it is easy to see from definition that $B_{p,q}^\alpha \subset B_{p,q'}^\alpha$, if $q < q'$. Without loss
100 of generalizability, in the following discussion we will only focus on $B_{p,\infty}^\alpha$.

101 When $p = 1$, the Besov space allows higher inhomogeneity, and it is more general than the Sobolev
102 or Hölder space.

103 **Bounded variation (BV) space** is a more interpretable class of functions with spatially hetero-
104 geneous smoothness [10]. It is defined through the total variation (TV) of a function. For
105 $(m+1)$ th differentiable function $f : [0, 1] \rightarrow \mathbb{R}$, the m th order total variation is defined as
106 $TV^{(m)}(f) := TV(f^{(m+1)}) = \int_{[0,1]} |f^{(m+1)}(x)| dx$, and the corresponding m th order Bounded
107 Variation class $BV(m) := \{f : TV(f^{(m)}) < \infty\}$. The more general definition is given in Sec-
108 tion C.2. Bounded variation class is tightly connected to Besov classes. Specifically [8]:

$$B_{1,1}^{m+1} \subset BV(m) \subset B_{1,\infty}^{m+1} \quad (1)$$

109 This allows the results derived for the Besov space to be easily applied to BV space.

110 **Minimax MSE** It is well known that minimax rate for Besov and 1D BV classes are $O(n^{-\frac{2\alpha}{2\alpha+d}})$
111 and $O(n^{-(2m+2)/(2m+3)})$ respectively . The minimax rate for *linear estimators* in 1D BV classes is
112 known to be $O(n^{-(2m+1)/(2m+2)})$ [24, 10].

113 3 Main Results: Parallel ReLU DNNs

114 Consider a parallel neural network containing M multi layer perceptrons (MLP) with ReLU activa-
115 tion functions called *subnetworks*. Each subnetwork has width w and depth L . The input is fed to

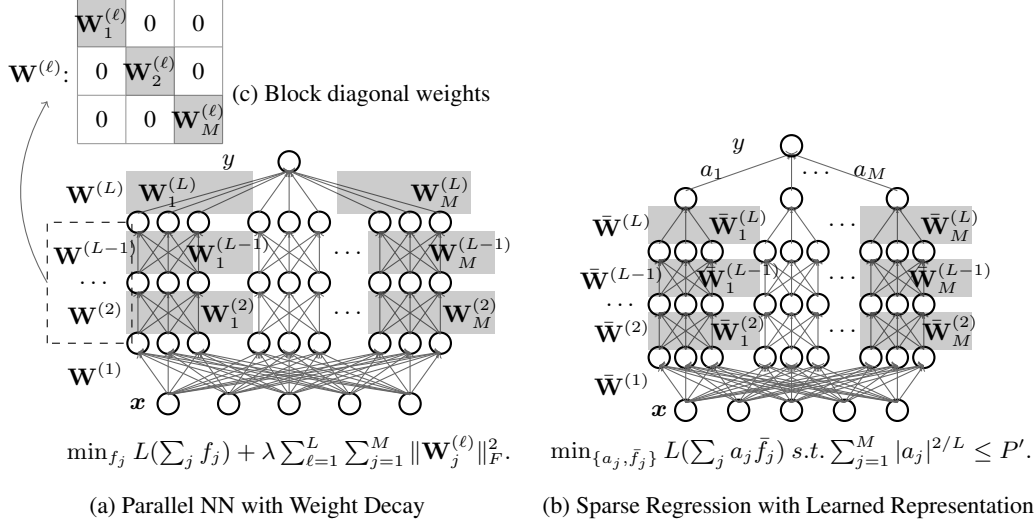


Figure 1: Parallel neural network and the equivalent sparse regression model we discovered.

116 all the subnetworks, and the output of the parallel NN is the summation of the output of each sub-
 117 network. The architecture of a parallel neural network is shown in Figure 1a. This parallel neural
 118 network is equivalent to a vanilla neural network with block diagonal weights in all but the first and
 119 the last layers (Figure 1(c)). Let $\mathbf{W}_j^{(\ell)}$ and $\mathbf{b}_j^{(\ell)}$ denote the weight and bias in the ℓ -th layer in the
 120 j -th subnetwork respectively. Training this model with weight decay returns:

$$\arg \min_{\{\mathbf{W}_j^{(\ell)}, \mathbf{b}_j^{(\ell)}\}} \hat{L}(f) + \lambda \sum_{j=1}^M \sum_{\ell=1}^L \|\mathbf{W}_j^{(\ell)}\|_F^2, \quad (2)$$

121 where $f(x) = \sum_{j=1}^M f_j(x)$ denotes the parallel neural network, $f_j(\cdot)$ denotes the j -th subnetwork,
 122 and $\lambda > 0$ is a fixed scaling factor.

123
 124 **Theorem 1.** For any fixed $\alpha - d/p > 1, q \geq 1, L \geq 3$, for any $f_0 \in B_{p,q}^\alpha$, given an L -layer parallel
 125 neural network satisfying

- 126 • The width of each subnetwork is fixed and large enough: $w \gtrsim d$. See Theorem 9 for the
 127 detail.
- 128 • The number of subnetworks is large enough: $M \gtrsim m^d n^{\frac{1-2/L}{2\alpha/d+1-2/(pL)}}$ where $m = \lceil \alpha - 1 \rceil$.

129 With proper choice of the parameter of weight decay λ , the solution \hat{f} parameterized by (2) satisfies
 130

$$\text{MSE}(\hat{f}) = \tilde{O}\left(n^{-\frac{2\alpha/d(1-2/L)}{2\alpha/d+1-2/(pL)}}\right) + \text{Const.} \quad (3)$$

131 where \tilde{O} shows the scale up to a logarithmic factor, and the trailing constant term decreases expo-
 132 nentially with L .

133 We explain the proof idea in the next section, but defer the extended form of the theorem and the full
 134 proof to Section F. Before that, we comment on a few interesting aspects of the result.

135 **Near optimal rates and the effect of depth.** The first term in the MSE bound is the estimation error
 136 and the second term is (part of) the approximation error of this NN. Recall that the minimax rate of
 137 a Besov class is $O(n^{-\frac{2\alpha}{2\alpha+d}})$ thus as the depth parameter L increases it can get arbitrarily close to
 138 the minimax rate. The constant term would be a negligible if we choose $L \gtrsim \log n$.

Corollary 2. Under the conditions of Theorem 1, for any $f_0 \in B_{p,q}^\alpha$, there is a numerical constant
 C such that when we choose $C \log n \leq L \leq 100C \log n$,

$$\text{MSE}(\hat{f}) = \tilde{O}\left(n^{-\frac{2\alpha}{2\alpha+d}(1-o(1))}\right),$$

139 where \tilde{O} hides only logarithmic factors and the $o(1)$ factor in the exponent is $O(1/\log(n))$.

140 This result says that deeper parallel neural networks achieves lower error and gets closer to the
 141 statistical limit.

142 **Overparameterization and sparsity.** We also note that the result does not depend on M as long
 143 as M is large enough. This means that the neural network can be arbitrarily overparameterized
 144 while not overfitting. The underlying reason is *sparsity*. As it will become clearer in the proof
 145 sketch, weight decayed training of a parallel L -layer ReLU NNs is equivalent to a sparse regression
 146 problem with an ℓ_p penalty assigned to the coefficient vector of a learned dictionary. Here $p = 2/L$
 147 which promotes even sparser solutions than an ℓ_1 penalty.

148 **No architecture tuning.** For any fixed L , the required architecture of the model does not depend
 149 on the dataset or the target function (n, α) except the number of subnetworks M , for which the only
 150 requirement is being large enough. As a result, one can design a model using a large guess on M ,
 151 and achieve the claimed near-optimal error rate by only tuning the weight decay parameter.

152 **Bounded variation classes.** Thanks to the Besov space embedding of the BV class (1), our theorem
 153 also implies the result for the BV class in $1D$.

154 **Corollary 3.** *If the target function is in bounded variation class $f_0 \in BV(m)$, For any fixed $L \geq 3$,
 155 for a neural network satisfying the requirements in Theorem 1 with $d = 1$ and with proper choice of
 156 the parameter of weight decay λ , the NN \hat{f} parameterized by (5) satisfies*

$$\text{MSE}(\hat{f}) = \tilde{O}(n^{-\frac{(2m+2)(1-2/L)}{2m+3-2/L}}) + \text{Const.}$$

157 where \tilde{O} shows the scale up to a logarithmic factor, and the trailing constant term decreases expo-
 158 nentially with L .

159 It is known that any linear estimators such as kernel smoothing and smoothing splines cannot have
 160 an error lower than $O(n^{-(2m+1)/(2m+2)})$ for $BV(m)$ [10]. This partly explains the advantage of
 161 DNNs over kernels.

162 **Representation learning and adaptivity.** The results also shed a light on the role of representation
 163 learning in DNN's ability to adapt. Specifically, different from the two-layer NN in [29], which
 164 achieves the minimax rate of $BV(m)$ by choosing appropriate activation functions using each m ,
 165 each subnetwork of a parallel NN can learn to approximate the spline basis of an arbitrary order,
 166 which means that if we choose L to be sufficiently large, such Parallel NN with optimally tuned λ is
 167 simultaneously near optimal for $m = 1, 2, 3, \dots$. In fact, even if different regions of the space has
 168 different *orders* of smoothness, the parallel NN will still be able to learn appropriate basis functions
 169 in each local region. To the best of our knowledge, this is a property that none of the classical
 170 nonparametric regression methods possess.

171 **Synthesis v.s. analysis methods.** Our result could also inspire new ideas in estimator design.
 172 There are two families of methods in non-parametric estimation. One called *synthesis* framework
 173 which focuses on constructing appropriate basis functions to encode the contemplated structures
 174 and regress the data to such basis, e.g., wavelets [10]. The other is called *analysis* framework which
 175 uses analysis regularization on the data directly (see, e.g., RKHS methods [37] or trend filtering
 176 [40]). It appears to us that parallel NN is doing both simultaneously. It has a parametric family
 177 capable to synthesizing an $O(n)$ subset of an exponentially large family of basis, then *implicitly*
 178 use sparsity-inducing analysis regularization to select the relevant basis functions. In this way the
 179 estimator does not actually have to explicitly represent that exponentially large set of basis functions,
 180 thus computationally more efficient.

181 4 Proof Overview

182 We start by first proving that a parallel neural network trained with weight decay is equivalent to an
 183 ℓ_p -sparse regression problem with representation learning (Section 4.1); which helps decompose its
 184 MSE into an estimation error and approximation error. Then we bound the two terms in Section 4.2
 185 and Section 4.3 respectively.

186 4.1 Equivalence to ℓ_p Sparse Regression with a Learned Feature Representation

187 It is widely known that ReLU function is 1-homogeneous: $\sigma(ax) = a\sigma(x), \forall a \geq 0, x \in \mathbb{R}$. In any
 188 consecutive two layers in a neural network (or a subnetwork), one can multiply the weight and bias

189 in one layer with a positive constant, and divide the weight in another layer with the same constant.
 190 The neural network after such transformation is equivalent to the original one:

$$\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) = \frac{1}{c}\mathbf{W}^{(2)}\sigma(c\mathbf{W}^{(1)}\mathbf{x} + c\mathbf{b}^{(1)}), \quad \forall c > 0, \mathbf{x}. \quad (4)$$

191 This property allows us to reformulate (2) to an ℓ_p sparsity constraint problem:

192 **Proposition 4.** Fix the input dataset \mathcal{D}_n and a constant $c_1 > 0$. There exists an one-to-one mapping
 193 between $\lambda > 0$ and $P' > 0$ such that (2) is equivalent to the following problem:

$$\begin{aligned} \arg \min_{\{\bar{\mathbf{W}}_j^{(\ell)}, \bar{\mathbf{b}}_j^{(\ell)}, a_j\}} \hat{L}\left(\sum_{j=1}^M a_j \bar{f}_j\right) &= \frac{1}{n} \sum_i (y_i - \bar{f}_{1:M}(\mathbf{x}_i)^T \mathbf{a})^2 \\ \text{s.t. } \|\bar{\mathbf{W}}_j^{(1)}\|_F &\leq c_1 \sqrt{d}, \forall j \in [M], \\ \|\bar{\mathbf{W}}_j^{(\ell)}\|_F &\leq c_1 \sqrt{w}, \forall j \in [M], 2 \leq \ell \leq L, \quad \|\{a_j\}\|_{2/L} \leq P' \end{aligned} \quad (5)$$

194 where $\bar{f}_j(\cdot)$ is a subnetwork with parameters $\bar{\mathbf{W}}_j^{(\ell)}, \bar{\mathbf{b}}_j^{(\ell)}$.

195 This equivalent model is demonstrated in Figure 1b. The proof can be found in Section D.1. The
 196 constraint $\|\bar{\mathbf{W}}_j^{(1)}\|_F \lesssim \sqrt{d}, \|\bar{\mathbf{W}}_j^{(\ell)}\|_F \lesssim \sqrt{w}, \forall \ell > 1$ is typical in deep learning for better numerical
 197 stability. The equivalent model in Proposition 4 is also a parallel neural network, but it appends one
 198 layer with parameters $\{a_k\}$ at the end of the neural network and the constraint on the Frobenius
 199 norm is converted to the $2/L$ norm on the factors $\{a_k\}$. Since $L \gg 2$ in a typical application,
 200 $2/L \ll 1$ and this constraint can enforce a sparser model than that in Section B.

201 There are two useful implications of Proposition 4. First, it gives an intuitive explanation on how
 202 a weight decayed Parallel NN works. Specifically, it can be viewed as a sparse linear regression
 203 with representation learning. Second, the conversion into the constrained form allows us to adapt
 204 generic statistical learning machinery (a self-bounding argument) from Suzuki [39, Proposition 4]
 205 for studying this constrained ERM problem.

206 The adaptation is nontrivial because (1) our regression problem has a *fixed design* (so data points are
 207 not iid); (2) there is an *unconstrained* subspace with no bounded metric entropy. Specifically, our
 208 Proposition 15 shows that the MSE of the regression problem can be bounded by

$$\text{MSE}(\hat{f}) = O\left(\underbrace{\inf_{f \in \mathcal{F}} \text{MSE}(f)}_{\text{approximation error}} + \underbrace{\frac{\log \mathcal{N}(\mathcal{F}_{\parallel}, \delta, \|\cdot\|_{\infty}) + d(\mathcal{F}_{\perp})}{n}}_{\text{estimation error}} + \delta\right)$$

209 in which \mathcal{F} decomposes into $\mathcal{F}_{\parallel} \times \mathcal{F}_{\perp}$, where \mathcal{F}_{\perp} is an unconstrained subspace with finite dimension,
 210 and \mathcal{F}_{\parallel} is a compact set in the orthogonal complement with a δ -covering number of $\mathcal{N}(\mathcal{F}_{\parallel}, \delta, \|\cdot\|_{\infty})$
 211 in $\|\cdot\|_{\infty}$ -norm. This decomposes MSE into an approximation error and an estimation error. The novel
 212 analysis of these two represents the major technical contribution of this paper.

213 4.2 Estimation Error Analysis

214 The decomposition above reveals that to bound the estimation error, it suffices to compute the cov-
 215 ering number of the constraint set in the sup-norm of the function it represents.

216 Previous results that bound the covering number of neural networks [49, 39] depends on the width
 217 of the neural networks explicitly, which cannot be applied when analysing a potentially infinitely
 218 wide neural network. In this section, we leverage the ℓ_p -norm bounded coefficients to avoid the
 219 dependence in M in the covering number bound.

220 **Theorem 5.** The covering number of the model defined in (5) apart from the bias in the last layer
 221 satisfies

$$\log \mathcal{N}(\mathcal{F}, \delta) \lesssim w^{2+2/(1-2/L)} L^2 \sqrt{d} P'^{\frac{1}{1-2/L}} \delta^{-\frac{2/L}{1-2/L}} \log(wP'/\delta). \quad (6)$$

222 The proof can be found in Section D.2. It requires the following lemma:

223 **Lemma 6.** $\log \mathcal{N}(\mathcal{G}, \delta) \lesssim k \log(1/\delta)$ for some finite c_3 , and for any $g \in \mathcal{G}$, $|a| \leq 1$, we have
 224 $ag \in \mathcal{G}$. The covering number of $\mathcal{F} = \left\{ \sum_{i=1}^M a_i g_i \mid g_i \in \mathcal{G}, \|a\|_p^p \leq P, 0 < p < 1 \right\}$ for any $P > 0$
 225 satisfies

$$\log \mathcal{N}(\mathcal{F}, \epsilon) \lesssim k P^{\frac{1}{1-p}} (\delta/c_3)^{-\frac{p}{1-p}} \log(c_3 P/\delta)$$

226 up to a double logarithmic factor.

227 See Section D.3 for the proof of Lemma 6. The covering number in Theorem 5 does not depend
 228 on the number of subnetworks M . In other words, it provides a bound of estimation error for an
 229 arbitrarily wide parallel neural network as long as the total Frobenius norm is bounded.

230 4.3 Approximation Error Analysis

231 The approximation error analysis involves two steps. In Section 4.3.1, we analyse how a subnetwork
 232 can approximate a B-spline basis. Then in Section 4.3.2 we show that a sparse linear combination
 233 of B-spline bases approximates Besov functions. Both add up to the total error in approximating
 234 Besov functions with a parallel neural network (Theorem 9).

235 4.3.1 Approximation Error of B-spline Basis Function

236 As is shown in Section C.1, functions in Besov space can be alternatively represented in a sequence
 237 space via the coefficients of a cardinal B-spline basis. In this section we study the approximation
 238 ability of ReLU neural networks to B-spline basis function.

239 **Proposition 7.** *Let $M_{m,k,s}$ be the B-spline of order m with scale 2^{-k} in each dimension and position
 240 $\mathbf{s} \in \mathbb{R}^d$: $M_{m,k,s}(\mathbf{x}) := M_m(2^k(\mathbf{x} - \mathbf{s}))$, M_m is defined in (11). There exists a parallel neural
 241 network that has the structure and satisfy the constraint in Proposition 4 for d -dimensional input
 242 and one output, containing $M = O(m^d)$ subnetworks, each of which has width $w = O(d)$ and
 243 depth $L = O(\log(c(m, d)/\epsilon))$ for some constant w, c that depends only on m and d , denoted as
 244 $\tilde{M}_m(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^d$, such that*

- 245 • $|\tilde{M}_{m,k,s}(\mathbf{x}) - M_{m,k,s}(\mathbf{x})| \leq \epsilon$, if $0 \leq 2^k(x_i - s_i) \leq m + 1, \forall i \in [d]$,
- 246 • $\tilde{M}_{m,k,s}(\mathbf{x}) = 0$, otherwise.
- 247 • The weights in the last layer satisfy $\|a\|_{2/L}^{2/L} \lesssim 2^k m^d e^{2md/L}$.

248 The proof can be found in Section E.1. Note that the product of the coefficients among all the layers
 249 are proportional to 2^k , instead of 2^{km} when approximating truncated power basis functions. This is
 250 because the transformation from M_m to $M_{m,k,s}$ only scales the domain of the function by 2^k , while
 251 the codomain of the function is not changed. To apply the transformation to the neural network, one
 252 only need to scale weights in the first layer by 2^k , which is equivalent to scaling the weights in each
 253 layer by $2^{k/L}$ and adjusting the bias according.

254 4.3.2 Approximation Error in Besov Space

255 With the results given in Section 4.3.1, we can estimate the approximation error of parallel ReLU
 256 neural networks to functions in Besov space.

Proposition 8. *Let $\alpha - d/p > 1, r > 0$. For any function in Besov space $f_0 \in B_{p,q}^\alpha$ and any
 positive integer \bar{M} , there is an \bar{M} -sparse approximation using B-spline basis of order m satisfying
 $0 < \alpha < \min(m, m - 1 + 1/p)$: $\tilde{f}_{\bar{M}} = \sum_{i=1}^{\bar{M}} a_{k_i, \mathbf{s}_i} M_{m, k_i, \mathbf{s}_i}$ for any positive integer \bar{M} such that
 the approximation error is bounded as $\|\tilde{f}_{\bar{M}} - f_0\|_r \lesssim \bar{M}^{-\alpha/d} \|f_0\|_{B_{p,q}^\alpha}$, and the coefficients satisfy*

$$\|\{2^{k_i} a_{k_i, \mathbf{s}_i}\}_{k_i, \mathbf{s}_i}\|_p \lesssim \|f_0\|_{B_{p,q}^\alpha}.$$

257 The proof can be found in Section E.2.

259 **Remark 1.** *The requirement in Proposition 8: $\alpha - d/p > 1$ is stronger than the condition typically
 260 found in approximation theorem $\alpha - d/p \geq 0$ [11], so-called ‘‘Boundary of continuity’’, or the
 261 condition in Suzuki [39] $\alpha > d(1/p - 1/r)_+$. This is because although the functions in $B_{p,q}^\alpha$ when*

262 $0 \leq \alpha - d/p < 1$ can be approximated by B-spline basis, the sum of weighted coefficients may not
 263 converge. One simple example is the step function $f_{step}(x) = \mathbf{1}(x \geq 0.5)$, $f_{step} \in B_{1,\infty}^1$. Although
 264 it can be decomposed using first order B-spline basis as in (10), the summation of the coefficients is
 265 infinite. Actually one only needs a ReLU neural network with one hidden layer and two neurons to
 266 approximate this function to arbitrary precision, but the weight need to go to infinity.

267 **Theorem 9.** Under the same condition as Proposition 8, for any positive integer \bar{M} , any function
 268 in Besov space $f_0 \in B_{p,q}^\alpha$ can be approximated by a parallel neural network with no less than
 269 $O(m^d \bar{M})$ number of subnetworks satisfying:

- 270 1. Each subnetwork has width $w = O(d)$ and depth L .
- 271 2. The weights in each layer satisfy $\|\bar{\mathbf{W}}_k^{(\ell)}\|_F \leq O(\sqrt{w})$ except the first layer $\|\bar{\mathbf{W}}_k^{(1)}\|_F \leq$
 272 $O(\sqrt{d})$,
- 273 3. The scaling factors have bounded $2/L$ -norm: $\|\{a_j\}\|_{2/L}^{2/L} \lesssim m^d e^{2md/L} \bar{M}^{1-2/(pL)}$.
- 274 4. The approximation error is bounded by

$$\|\tilde{f} - f_0\|_r \leq (c_4 \bar{M}^{-\alpha/d} + c_5 e^{-c_6 L}) \|f\|_{B_{p,q}^\alpha}$$

275 where c_4, c_5, c_6 are constants that depend only on m, d and p .

276 Here \bar{M} is the number of ‘‘active’’ subnetworks, which is not to be confused with the number of
 277 subnetworks at initialization. The proof can be found in Section E.3.

278 Using the estimation error in Theorem 5 and approximation error in Theorem 9, by choosing \bar{M} to
 279 minimax the total error, we can conclude the sample complexity of parallel neural networks using
 280 weight decay, which is the main result (Theorem 1) of this paper. See Section F for the detail.

281 5 Experiment

282 We empirically compare a parallel neural network (PNN) and a vanilla ReLU neural network (NN)
 283 with smoothing spline, trend filtering (TF) [40], and wavelet denoising. Trend filtering can be
 284 viewed as a more efficient discrete spline version of locally adaptive regression spline and enjoys the
 285 same optimal rates for the BV classes. Wavelet denoising is also known to be minimax-optimal for
 286 the BV classes. The results are shown in Figure 2. We use two target functions: a Doppler function
 287 whose frequency is decreasing (Figure 2(a)-(c)), and a combination of piecewise linear function and
 288 piecewise cubic function, or ‘‘vary’’ function (Figure 2(d)-(f)). We repeat each experiment 10 times
 289 and take the average. The shallow area in Figure 2(b)(e) shows 95% confidence interval by inverting
 290 the Wald’s test. The degree of freedom (DoF) is computed based on Tibshirani [41].

291 As can be shown in the figure, both TF and wavelet denoising can adapt to the different levels of
 292 smoothness in the target function, while smoothing splines tend to be oversmoothed where the target
 293 function is less smooth (the left side in (a)(d), enlarged in (g)). The prediction of PNN is similar to
 294 TF and wavelet denoising and shows local adaptivity. Besides, the MSE of PNN almost follows the
 295 same trend as TF and wavelet denoising which is consistent with our theoretical understanding that
 296 the error rate of neural network is closer to locally adaptive methods. Notably PNN, TF and wavelet
 297 denoising achieve lower error at a much smaller degree-of-freedom than smoothing splines.

298 In a vanilla NN, weight decay is equivalent to ℓ_1 regularizer in any two successive layers, but to the
 299 best of our knowledge it does not lead to sparse representation learning unless some specific sparse
 300 structure is enforced. While our theory does not apply to vanilla neural networks, the results seem
 301 to suggest the NN behaves similar to smoothing spline and is *not* locally adaptive.

302 There are some mild drops in the best MSE one can achieve with NN vs TF in both examples. We
 303 are surprised that the drop is small because NN needs to learn the basis functions that TF essentially
 304 hard-coded. The additional price to pay for using a more adaptive and more flexible representation
 305 learning method seems not high at all.

306 In Figure 2(c)(f), we give the output [all](#) the ‘‘active’’ subnetwork, i.e. the subnetworks whose output is
 307 not a constant. Notice that the number of active subnetworks is much smaller than the initialization.
 308 This is because weight decay induces ℓ_p sparsity and the weight in most of the subnetworks reduces
 309 towards 0 after training. More details are shown in Section G.

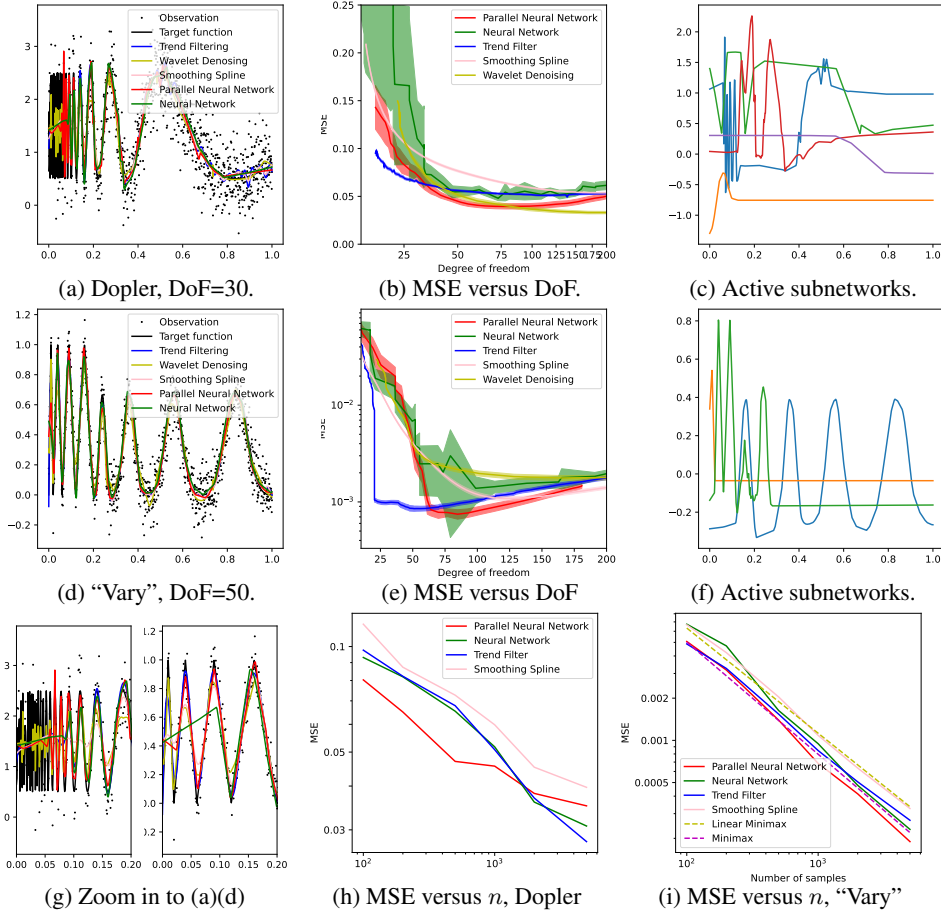


Figure 2: Numerical experiment results of the Doppler function (a-c,h), and “vary” function (d-f,g). All the “active” subnetworks are plotted in (c)(f). The horizontal axis in (b) is not linear.

310 In Figure 2(h)(i), we plot the MSE versus the number of training samples for “Doppler” and “Vary”
 311 respectively. It is clear that parallel NN works the best overall. In (i), we further compare the scaling
 312 of the MSE against the minimax rate ($n^{-4/5}$) and the minimax linear rate ($n^{-3/4}$), i.e., the best rate
 313 kernel methods could achieve. As is predicted by our theory, when n is large, the MSE of parallel
 314 neural networks and trend filtering decreases at almost the same rate as the minimax rate, while
 315 smoothing splines, as expected, is converging at the (suboptimal) minimax linear rate. Interestingly,
 316 vanilla NN seems to converge at the optimal rate too on this example. It remains an open question
 317 whether vanilla NN is merely “lucky” on this example, or it also achieves the minimax rate for all
 318 functions in $BV(m)$.

319 6 Conclusion and Discussion

320 In this paper, we show that a deep parallel neural network can be locally adaptive by tuning only
 321 the weight decay parameter. This confirms that neural networks can be nearly optimal in learning
 322 functions with heterogeneous smoothness which separates them from kernel methods. We prove
 323 that training an L layer parallel neural network with weight decay is equivalent to an $\ell_{2/L}$ -penalized
 324 regression model with representation learning. Since in typical application $L \gg 2$, weight decay
 325 promotes a sparse linear combination of the learned bases. Using this method, we proved that a
 326 parallel neural network can achieve close to the minimax rate in the Besov space and bounded
 327 variation (BV) space. Our result reveals that one do not need to specify the smoothness parameter α
 328 (or m). Neural networks can adapt to different degree of smoothness, or choose different parameters
 329 for different regions of the domain of the target function. This is a new type of adaptivity not

330 possessed by traditional adaptive nonparametric regression methods like locally adaptive regression
331 spline or trend filtering.

332 On the other hand, as the depth of neural network L increases, $2/L$ tends to 0 and the error rate
333 moves closer to the minimax rate of Besov and BV space. This indicates that when the sample size
334 is large enough, deeper models have smaller error than shallower models, and helps explain why
335 empirically deep neural networks has better performance than shallow neural networks.

336 References

- 337 [1] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang.
338 On exact computation with an infinitely wide neural net. In *Proceedings of the 33rd Interna-*
339 *tional Conference on Neural Information Processing Systems*, pages 8141–8150, 2019.
- 340 [2] Dheeraj Baby and Yu-Xiang Wang. Online forecasting of total-variation-bounded sequences.
341 In *Neural Information Processing Systems (NeurIPS)*, 2019.
- 342 [3] Dheeraj Baby and Yu-Xiang Wang. Adaptive online estimation of piecewise polynomial
343 trends. *Neural Information Processing Systems (NeurIPS)*, 2020.
- 344 [4] Andrew R Barron. Approximation and estimation bounds for artificial neural networks. *Ma-*
345 *chine learning*, 14(1):115–133, 1994.
- 346 [5] Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to
347 understand kernel learning. In *International Conference on Machine Learning*, pages 541–
348 549. PMLR, 2018.
- 349 [6] George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of*
350 *control, signals and systems*, 2(4):303–314, 1989.
- 351 [7] Carl De Boor, Carl De Boor, Etats-Unis Mathématicien, Carl De Boor, and Carl De Boor. *A*
352 *practical guide to splines*, volume 27. Springer-Verlag New York, 1978.
- 353 [8] Ronald A DeVore and George G Lorentz. *Constructive approximation*, volume 303. Springer
354 Science & Business Media, 1993.
- 355 [9] David L Donoho, Richard C Liu, and Brenda MacGibbon. Minimax risk over hyperrectangles,
356 and implications. *The Annals of Statistics*, pages 1416–1437, 1990.
- 357 [10] David L Donoho, Iain M Johnstone, et al. Minimax estimation via wavelet shrinkage. *The*
358 *annals of Statistics*, 26(3):879–921, 1998.
- 359 [11] Dinh Dũng. Optimal adaptive sampling recovery. *Advances in Computational Mathematics*,
360 34(1):1–41, 2011.
- 361 [12] Tolga Ergen and Mert Pilanci. Implicit convex regularizers of cnn architectures: Con-
362 vex optimization of two-and three-layer networks in polynomial time. *arXiv preprint*
363 *arXiv:2006.14798*, 2020.
- 364 [13] Tolga Ergen and Mert Pilanci. Convex geometry and duality of over-parameterized neural
365 networks. *Journal of machine learning research*, 2021.
- 366 [14] Tolga Ergen and Mert Pilanci. Global optimality beyond two layers: Training deep relu net-
367 works via convex programs. In *International Conference on Machine Learning*, pages 2993–
368 3003. PMLR, 2021.
- 369 [15] Tolga Ergen and Mert Pilanci. Path regularization: A convexity and sparsity inducing regular-
370 ization for parallel relu networks. *arXiv preprint arXiv:2110.09548*, 2021.
- 371 [16] Tolga Ergen and Mert Pilanci. Revealing the structure of deep neural networks via convex
372 duality. In *International Conference on Machine Learning*, pages 3004–3014. PMLR, 2021.
- 373 [17] Jerome Friedman, Trevor Hastie, and Rob Tibshirani. Regularization paths for generalized
374 linear models via coordinate descent. *Journal of statistical software*, 33(1):1, 2010.

- 375 [18] Benjamin D Haeffele and René Vidal. Global optimality in neural network training. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7331–
376 7339, 2017.
377
- 378 [19] Daniel Hsu, Sham M Kakade, and Tong Zhang. An analysis of random design linear regression.
379 *arXiv preprint arXiv:1106.2363*, 2011.
- 380 [20] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Bina-
381 rized neural networks. In *Proceedings of the 30th international conference on neural informa-
382 tion processing systems*, pages 4114–4122. Citeseer, 2016.
- 383 [21] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: convergence and
384 generalization in neural networks. In *Proceedings of the 32nd International Conference on
385 Neural Information Processing Systems*, pages 8580–8589, 2018.
- 386 [22] Seung-Jean Kim, Kwangmoo Koh, Stephen Boyd, and Dmitry Gorinevsky. ℓ_1 trend filter-
387 ing. *SIAM review*, 51(2):339–360, 2009.
- 388 [23] Stéphane Mallat. *A wavelet tour of signal processing*. Elsevier, 1999.
- 389 [24] Enno Mammen and Sara van de Geer. Locally adaptive regression splines. *The Annals of
390 Statistics*, 25(1):387–413, 1997.
- 391 [25] Elizbar A Nadaraya. On estimating regression. *Theory of Probability & Its Applications*, 9(1):
392 141–142, 1964.
- 393 [26] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever.
394 Deep double descent: Where bigger models and more data hurt. *Journal of Statistical Mechan-
395 ics: Theory and Experiment*, 2021(12):124003, 2021.
- 396 [27] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias:
397 On the role of implicit regularization in deep learning. *arXiv preprint arXiv:1412.6614*, 2014.
- 398 [28] Greg Ongie, Rebecca Willett, Daniel Soudry, and Nathan Srebro. A function space view of
399 bounded norm infinite width relu nets: The multivariate case. In *International Conference on
400 Learning Representations*, 2019.
- 401 [29] Rahul Parhi and Robert D Nowak. Banach space representer theorems for neural networks and
402 ridge splines. *J. Mach. Learn. Res.*, 22:43–1, 2021.
- 403 [30] Rahul Parhi and Robert D Nowak. What kinds of functions do deep neural networks learn?
404 insights from variational spline theory. *arXiv preprint arXiv:2105.03361*, 2021.
- 405 [31] Rahul Parhi and Robert D Nowak. Near-minimax optimal estimation with shallow relu neural
406 networks. *arXiv preprint arXiv:2109.08844*, 2021.
- 407 [32] Mert Pilanci and Tolga Ergen. Neural networks are convex regularizers: Exact polynomial-
408 time convex optimization formulations for two-layer networks. In *International Conference
409 on Machine Learning*, pages 7695–7705. PMLR, 2020.
- 410 [33] Carl Edward Rasmussen and Christopher KI Williams. *Gaussian processes for machine learn-
411 ing*. MIT Press, 2006.
- 412 [34] Veeranjaneyulu Sadhanala, Yu-Xiang Wang, Addison J Hu, and Ryan J Tibshirani. Multivari-
413 ate trend filtering for lattice data. *arXiv preprint arXiv:2112.14758*, 2021.
- 414 [35] Pedro Savarese, Itay Evron, Daniel Soudry, and Nathan Srebro. How do infinite width bounded
415 norm networks look in function space? In *Conference on Learning Theory*, pages 2667–2690.
416 PMLR, 2019.
- 417 [36] Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with relu
418 activation function. *The Annals of Statistics*, 48(4):1875–1897, 2020.
- 419 [37] Bernhard Scholkopf and Alexander J Smola. *Learning with kernels: support vector machines,
420 regularization, optimization, and beyond*. MIT press, 2001.

- 421 [38] Nathan Srebro, Jason DM Rennie, and Tommi S Jaakkola. Maximum-margin matrix factor-
422 ization. In *NIPS*, volume 17, pages 1329–1336. Citeseer, 2004.
- 423 [39] Taiji Suzuki. Adaptivity of deep relu network for learning in besov and mixed smooth besov
424 spaces: optimal rate and curse of dimensionality. *arXiv preprint arXiv:1810.08033*, 2018.
- 425 [40] Ryan J Tibshirani. Adaptive piecewise polynomial estimation via trend filtering. *The Annals*
426 *of Statistics*, 42(1):285–323, 2014.
- 427 [41] Ryan J Tibshirani. Degrees of freedom and model search. *Statistica Sinica*, pages 1265–1296,
428 2015.
- 429 [42] Ryan J Tibshirani. Equivalences between sparse models and neural networks. 2021. URL
430 <http://www.stat.cmu.edu/~ryantibs/papers/sparsitynn.pdf>.
- 431 [43] Ryan J Tibshirani. Personal communication, Jan. 24, 2022.
- 432 [44] Andreas Veit, Michael J Wilber, and Serge Belongie. Residual networks behave like ensembles
433 of relatively shallow networks. *Advances in neural information processing systems*, 29:550–
434 558, 2016.
- 435 [45] Stefan Wager, Sida Wang, and Percy Liang. Dropout training as adaptive regularization. *arXiv*
436 *preprint arXiv:1307.1493*, 2013.
- 437 [46] Grace Wahba. *Spline models for observational data*, volume 59. Siam, 1990.
- 438 [47] Yu-Xiang Wang, Alex Smola, and Ryan Tibshirani. The falling factorial basis and its statistical
439 applications. In *International Conference on Machine Learning*, pages 730–738. PMLR, 2014.
- 440 [48] Yuan Yao, Lorenzo Rosasco, and Andrea Caponnetto. On early stopping in gradient descent
441 learning. *Constructive Approximation*, 26(2):289–315, 2007.
- 442 [49] Dmitry Yarotsky. Error bounds for approximations with deep relu networks. *Neural Networks*,
443 94:103–114, 2017. ISSN 0893-6080. doi: <https://doi.org/10.1016/j.neunet.2017.07.002>. URL
444 <https://www.sciencedirect.com/science/article/pii/S0893608017301545>.
- 445 [50] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint*
446 *arXiv:1605.07146*, 2016.
- 447 [51] Hongyang Zhang, Junru Shao, and Ruslan Salakhutdinov. Deep neural networks with multi-
448 branch architectures are intrinsically less non-convex. In *The 22nd International Conference*
449 *on Artificial Intelligence and Statistics*, pages 1099–1109. PMLR, 2019.

450 Checklist

- 451 1. For all authors...
- 452 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
453 contributions and scope? [Yes]
- 454 (b) Did you describe the limitations of your work? [Yes] See Section 6.
- 455 (c) Did you discuss any potential negative societal impacts of your work? [N/A] This
456 work does not have potential negative societal impacts to the best of our knowledge.
- 457 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
458 them? [Yes]
- 459 2. If you are including theoretical results...
- 460 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 461 (b) Did you include complete proofs of all theoretical results? [Yes] All the proofs are in
462 the appendix.
- 463 3. If you ran experiments...

- 464 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
465 imental results (either in the supplemental material or as a URL)? [No] Although we
466 did not publish the code, the experiments are very simple and can be easily reproduced
467 using the information given in Section G.
- 468 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
469 were chosen)? [Yes] See Section G.
- 470 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
471 ments multiple times)? [Yes] We provided the confidence interval in Figure 2.
- 472 (d) Did you include the total amount of compute and the type of resources used (e.g.,
473 type of GPUs, internal cluster, or cloud provider)? [No] The experiments are light-
474 weighted and can be finished in an ordinary PC in a reasonable time.
- 475 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 476 (a) If your work uses existing assets, did you cite the creators? [N/A] We only used
477 synthetic dataset.
- 478 (b) Did you mention the license of the assets? [N/A]
- 479 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 480
- 481 (d) Did you discuss whether and how consent was obtained from people whose data
482 you're using/curating? [N/A]
- 483 (e) Did you discuss whether the data you are using/curating contains personally identifi-
484 able information or offensive content? [N/A]
- 485 5. If you used crowdsourcing or conducted research with human subjects...
- 486 (a) Did you include the full text of instructions given to participants and screenshots, if
487 applicable? [N/A]
- 488 (b) Did you describe any potential participant risks, with links to Institutional Review
489 Board (IRB) approvals, if applicable? [N/A]
- 490 (c) Did you include the estimated hourly wage paid to participants and the total amount
491 spent on participant compensation? [N/A]

492 **A Other related works**

493 Besides Parhi and Nowak [29] which we discussed earlier, Parhi and Nowak [30, 31] also leveraged
 494 the connections between NNs and splines. Parhi and Nowak [30] focused on characterizing the
 495 variational form of multi-layer NN. Parhi and Nowak [31] showed that two-layer ReLU activated
 496 NN achieves minimax rate for a BV class of order 1 but did not cover multilayer NNs nor BV class
 497 with order > 1 , which is our focus.

498 The connection between weight-decay regularization with sparsity-inducing penalties in two-layer
 499 NNs is folklore and used by Neyshabur et al. [27], Savarese et al. [35], Ongie et al. [28], Ergen and
 500 Pilanci [13, 16], Parhi and Nowak [29, 31], Pilanci and Ergen [32]. The key underlying technique
 501 — an application of the AM-GM inequality (which we used in this paper as well) — can be traced
 502 back to Srebro et al. [38] (see a recent exposition by Tibshirani [42]). [42] also generalized the
 503 result to multi-layered NNs, but with a simple (element-wise) connections. [14] [generalized the](#)
 504 [results to a three-layer parallel neural network, and proved its equivalence to an \$\ell_1\$ sparse model, but](#)
 505 [this requires a non-standard regularizer. Besides, \[12\] proved that training a two-layer convolution](#)
 506 [neural network \(CNN\) with weight decay induces sparsity, and points to a potential extension to](#)
 507 [these works including our work.](#)

508 The approximation-theoretic and estimation-theoretic research for neural network has a long history
 509 too [6, 4, 49, 36, 39]. Most existing work considered the Holder, Sobolev spaces and their exten-
 510 sions, which contain only homogeneously smooth functions and cannot demonstrate the advantage
 511 of NNs over kernels. The only exception is Suzuki [39] which, as we discussed earlier, requires
 512 modifications to NN architecture for each class. In contrast, we require tuning only the standard
 513 weight decay parameter.

514 **B Two-layer Neural Network with Truncated Power Activation Functions**

515 We start by recapping the result of Parhi and Nowak [29] and formalizing its implication in esti-
 516 mating BV functions. Parhi and Nowak [29] considered a two layer neural network with truncated
 517 power activation function. Let the neural network be

$$f(x) = \sum_{j=1}^M v_j \sigma^m(w_j x + b_j) + c(x), \quad (7)$$

518 where w_j, v_j denote the weight in the first and second layer respectively, b_j denote the bias in the
 519 first layer, $c(x)$ is a polynomial of order up to m , $\sigma^m(x) := \max(x, 0)^m$. Parhi and Nowak [29,
 520 Theorem 8] showed that when M is large enough, The optimization problem

$$\min_{w, v} \hat{L}(f) + \frac{\lambda}{2} \sum_{j=1}^M (|v_j|^2 + |w_j|^{2m}) \quad (8)$$

521 is equivalent to the locally adaptive regression spline:

$$\min_f \hat{L}(f) + \lambda TV(f^{(m)}(x)), \quad (9)$$

522 which optimizes over arbitrary functions that is m -times weakly differentiable. The latter was
 523 studied in Mammen and van de Geer [24], which leads to the following MSE:

524 **Theorem 10.** *Let $M \geq n - m$, and \hat{f} be the function (7) parameterized by the minimizer of (8),*
 525 *then*

$$\text{MSE}(\hat{f}) = O(n^{-(2m+2)(2m+3)}).$$

526 We show a simpler proof in the univariate case due to Tibshirani [43]:

527 *Proof.* As is shown in Parhi and Nowak [29, Theorem 8], the minimizer of (8) satisfy

$$|v_j| = |w_j|^m, \forall k$$

528 so the TV of the neural network f_{NN} is

$$\begin{aligned} TV^{(m)}(f_{NN}) &= TV^{(m)}c(x) + \sum_{j=1}^M |v_j||w_j|^m TV^{(m)}(\sigma^{(m)}(x)) \\ &= \sum_{j=1}^M |v_j||w_j|^m \\ &= \frac{1}{2} \sum_{j=1}^M (|v_j|^2 + |w_j|^{2m}) \end{aligned}$$

529 which shown that (8) is equivalent to the locally adaptive regression spline (9) as long as the number
530 of knots in (9) is no more than M . Furthermore, it is easy to check that any spline with knots no
531 more than M can be expressed as a two layer neural network (8). It suffices to prove that the solution
532 in (9) has no more than $n - m$ number of knots.

533 Mammen and van de Geer [24, Proposition 1] showed that there is a solution to (9) $\hat{f}(x)$ such that
534 $\hat{f}(x)$ is a m th order spline with a finite number of knots but did not give a bound. Let the number of
535 knots be M , we can represent \hat{f} using the truncated power basis

$$\hat{f}(x) = \sum_{j=1}^M a_j (x - t_j)_+^m + c(x) := \sum_{j=1}^M a_j \sigma_j^{(m)}(x) + c(x)$$

536 where t_j are the knots, $c(x)$ is a polynomial of order up to m , and define $\sigma_j^{(m)}(x) = (x - t_j)_+^m$.

537 Mammen and van de Geer [24] however did not give a bound on M . Parhi and Nowak [29]'s
538 Theorem 1 implies that $M \leq n - m$. Its proof is quite technical and applies more generally to a
539 higher dimensional generalization of the BV class.

540 Tibshirani [43] communicated to us the following elegant argument to prove the same using elemen-
541 tary convex analysis and linear algebra, which we present below.

542 Define $\Pi_m(f)$ as the $L^2(P_n)$ projection of f onto polynomials of degree up to m , $\Pi_m^\perp(f) :=$
543 $f - \Pi_m(f)$. It is easy to see that

$$\Pi_m^\perp f(x) = \sum_{j=1}^M a_j \Pi_m^\perp \sigma_j^{(m)}(x)$$

544 Denote $f(x_{1:n}) := \{f(x_1), \dots, f(x_n)\} \in \mathbb{R}^n$ as a vector of all the predictions at the sample points.

$$\Pi_m^\perp \hat{f}(x_{1:n}) = \sum_{j=1}^M a_j \Pi_m^\perp \sigma_j^{(m)}(x_{1:n}) \in \Pi_m^\perp \text{conv}\{\pm \sigma_j^{(m)}(x_{1:n})\} \cdot \sum_{j=1}^M |a_j| \in \text{conv}\{\pm \Pi_m^\perp \sigma_j^{(m)}(x_{1:n})\} \cdot \sum_{j=1}^M |a_j|$$

545 where conv denotes the convex hull of a set. The convex hull $\text{conv}\{\pm \sigma_j^{(m)}(x_{1:n})\} \cdot \sum_{j=1}^M |a_j|$ is an
546 n -dimensional space, and polynomials of order up to m is an $m + 1$ dimensional space, so the set
547 defined above has dimension $n - m - 1$. By Carathéodory's theorem, there is a subset of points in
548 this space

$$\{\Pi_m^\perp \sigma_{j_k}^{(m)}(x_{1:n})\} \subseteq \{\Pi_m^\perp \sigma_j^{(m)}(x_{1:n})\}, 1 \leq k \leq n - m$$

549 such that

$$\Pi_m^\perp f(x) = \sum_{k=1}^{n-m} \tilde{a}_k \Pi_m^\perp \sigma_{j_k}^{(m)}(x), \sum_{k=1}^{n-m} |a_k| \leq 1$$

550 In other word, there exist a subset of knots $\{\tilde{t}_j, j \in [n - m]\}$ that perfectly recovers $\Pi_m^\perp \hat{f}(x)$ at all
551 the sample points, and the TV of this function is no larger than \hat{f} .

This shows that

$$\tilde{f}(x) = \sum_{j=1}^{n-m} \tilde{a}_j (x - t_j)_+^m, \text{ s.t. } \tilde{f}(x_i) = f(x_i)$$

552 for all x_i in n onbservation points.

553 The MSE of locally adaptivity regressive spline (9) was studied in Mammen and van de Geer [24,
554 Section 3], which equals the error rate given in Theorem 10. \square

555 This indicates that the neural network (7) is minimax optimal for $BV(m)$.

Let us explain a few the key observations behind this equivalence. (a) The truncated power functions (together with an m th order polynomial) spans the space of an m th order spline. (b) The neural network in (7) is equivalent to a free-knot spline with M knots (up to reparameterization). (c) A solution to (9) is a spline with at most $n - m$ knots [29, Theorem 8]. (d) Finally, by the AM-GM inequality

$$|v_j|^2 + |w_j|^{2m} \geq 2|v_j||w_j|^m = 2|c_j|$$

556 where $c_j = v_j|w_j|^m$ is the coefficient of the corresponding j th truncated power basis. The m th
557 order total variation of a spline is equal to $\sum_j |c_j|$. It is not hard to check that the loss function
558 depends only on c_j , thus the optimal solution will always take “=” in the AM-GM inequality.

559 C Introduction To Common Function Classes

560 In the following definition define Ω be the domain of the function classes, which will be omitted in
561 the definition.

562 C.1 Besov Class

563 **Definition 1.** *Modulus of smoothness:* For a function $f \in L^p(\Omega)$ for some $1 \leq p \leq \infty$, the r -th
564 modulus of smoothness is defined by

$$w_{r,p}(f, t) = \sup_{h \in \mathbb{R}^d: \|h\|_2 \leq t} \|\Delta_h^r(f)\|_p,$$

565

$$\Delta_h^r(f) := \begin{cases} \sum_{j=0}^r \binom{r}{j} (-1)^{r-j} f(x + jh), & \text{if } x \in \Omega, x + rh \in \Omega, \\ 0, & \text{otherwise.} \end{cases}$$

566 **Definition 2.** *Besov space:* For $1 \leq p, q \leq \infty, \alpha > 0, r := \lceil \alpha \rceil + 1$, define

$$|f|_{B_{p,q}^\alpha} = \begin{cases} \left(\int_{t=0}^{\infty} (t^{-\alpha} w_{r,p}(f, t))^q \frac{dt}{t} \right)^{\frac{1}{q}}, & q < \infty \\ \sup_{t>0} t^{-\alpha} w_{r,p}(f, t), & q = \infty, \end{cases}$$

567 and define the norm of Besov space as:

$$\|f\|_{B_{p,q}^\alpha} = \|f\|_p + |f|_{B_{p,q}^\alpha}.$$

568 A function f is in the Besov space $B_{p,q}^\alpha$ if $\|f\|_{B_{p,q}^\alpha}$ is finite.

569 Note that the Besov space for $0 < p, q < 1$ is also defined, but in this case it is a quasi-Banach space
570 instead of a Banach space and will not be covered in this paper.

571 Functions in Besov space can be decomposed using B-spline basis functions. Any function f in
572 Besov space $B_{p,q}^\alpha, \alpha > d/p$ can be decomposed using B-spline of order $m, m > \alpha$: let $\mathbf{x} \in \mathbb{R}^d$,

$$f(\mathbf{x}) = \sum_{k=0}^{\infty} \sum_{\mathbf{s} \in J(k)} c_{k,\mathbf{s}}(f) M_{m,k,\mathbf{s}}(\mathbf{x}) \quad (10)$$

573 where $J(k) := \{2^{-k} \mathbf{s} : \mathbf{s} \in [-m, 2^k + m]^d \subset \mathbb{Z}^d\}$, $M_{m,k,\mathbf{s}}(\mathbf{x}) := M_m(2^k(\mathbf{x} - \mathbf{s}))$, and $M_k(\mathbf{x}) =$
574 $\prod_{i=1}^d M_k(x_i)$ is the cardinal B-spline basis function which can be expressed as a polynomial:

$$\begin{aligned} M_m(x) &= \frac{1}{m!} \sum_{j=1}^{m+1} (-1)^j \binom{m+1}{j} (x-j)_+^m \\ &= ((m+1)/2)^m \frac{1}{m!} \sum_{j=1}^{m+1} (-1)^j \binom{m+1}{j} \left(\frac{x-j}{(m+1)/2} \right)_+^m, \end{aligned} \quad (11)$$

575 Furthermore, the norm of Besov space is equivalent to the sequence norm:

$$\|\{c_{k,s}\}\|_{b_{p,q}^\alpha} := \left(\sum_{k=0}^{\infty} (2^{(\alpha-d/p)k} \|\{c_{k,s}(f)\}_s\|_p)^q \right)^{1/q} \approx \|f\|_{B_{p,q}^\alpha}.$$

576 See e.g. D ng [11, Theorem 2.2] for the proof.

577 The Besov space is closely connected to other function spaces including the H lder space (\mathcal{C}^α) and
578 the Sobolev space (W_p^α). Specifically, if the domain of the functions is d -dimensional [39, 34],

- 579 • $\forall \alpha \in \mathbb{N}, B_{p,1}^\alpha \subset W_p^\alpha \subset B_{p,\infty}^\alpha$, and $B_{2,2}^\alpha = W_2^\alpha$.
- 580 • For $0 < \alpha < \infty$ and $\alpha \in \mathcal{N}$, $\mathcal{C}^\alpha = B_{\infty,\infty}^\alpha$.
- 581 • If $\alpha > d/p$, $B_{p,q}^\alpha \subset \mathcal{C}^0$.

582 C.2 Other Function Spaces

583 **Definition 3.** *H lder space:* let $m \in \mathbb{N}$, the m -th order H lder class is defined as

$$\mathcal{C}^m = \left\{ f : \max_{|a|=k} \frac{|D^a f(x) - D^a f(z)|}{\|x - z\|_2} < \infty, \forall x, z \in \Omega \right\}$$

584 where D^a denotes the weak derivative.

585 Note that fraction order of H lder space can also be defined. For simplicity, we will not cover that
586 case in this paper.

587 **Definition 4.** *Sobolev space:* let $m \in \mathcal{N}$, $1 \leq p \leq \infty$, the Sobolev norm is defined as

$$\|f\|_{W_p^m} := \left(\sum_{|a| \leq m} \|D^a f\|_p^p \right)^{1/p},$$

588 the Sobolev space is the set of functions with finite Sobolev norm:

$$W_p^m := \{f : \|f\|_{W_p^m} < \infty\}.$$

589 **Definition 5.** *Total Variation (TV):* The total variation (TV) of a function f on an interval $[a, b]$ is
590 defined as

$$TV(f) = \sup_{\mathcal{P}} \sum_{i=1}^{n_{\mathcal{P}}-1} |f(x_{i+1}) - f(x_i)|$$

591 where the \mathcal{P} is taken among all the partitions of the interval $[a, b]$.

592 In many applications, functions with stronger smoothness conditions are needed, which can be mea-
593 sured by high order total variation.

594 **Definition 6.** *High order total variation:* the m -th order total variation is the total variation of the
595 $(m - 1)$ -th order derivative

$$TV^{(m)}(f) = TV(f^{(m-1)})$$

596 **Definition 7.** *Bounded variation (BV):* The m -th order bounded variation class is the set of functions
597 whose total variation (TV) is bounded.

$$BV(m) := \{f : TV(f^{(m)}) < \infty\}.$$

598 D Proof of Estimation Error

599 D.1 Equivalence Between Parallel Neural Networks and p -norm Penalized Problems

600 **Proposition 4.** Fix the input dataset \mathcal{D}_n and a constant $c_1 > 0$. There exists an one-to-one mapping
601 between $\lambda > 0$ and $P' > 0$ such that (2) is equivalent to the following problem:

$$\begin{aligned} \arg \min_{\{\bar{\mathbf{w}}_j^{(\ell)}, \bar{b}_j^{(\ell)}, a_j\}} \hat{L} \left(\sum_{j=1}^M a_j \bar{f}_j \right) &= \frac{1}{n} \sum_i (y_i - \bar{f}_{1:M}(\mathbf{x}_i)^T \mathbf{a})^2 \\ \text{s.t. } \|\bar{\mathbf{w}}_j^{(1)}\|_F &\leq c_1 \sqrt{d}, \forall j \in [M], \\ \|\bar{\mathbf{w}}_j^{(\ell)}\|_F &\leq c_1 \sqrt{w}, \forall j \in [M], 2 \leq \ell \leq L, \quad \|\{a_j\}\|_{2/L}^2 \leq P' \end{aligned}$$

602 where $\bar{f}_j(\cdot)$ is a subnetwork with parameters $\bar{\mathbf{W}}_j^{(\ell)}, \bar{\mathbf{b}}_j^{(\ell)}$.

603 *Proof.* Using Lagrange's method, one can easily find (2) is equivalent to a constrained optimization
604 problem:

$$\arg \min_{\{\mathbf{W}_j^{(\ell)}, \mathbf{b}_j^{(\ell)}\}} \hat{L} \left(\sum_{j=1}^M f_j \right), \quad s.t. \sum_{j=1}^M \sum_{\ell=1}^L \|\mathbf{W}_j^{(\ell)}\|_F^2 \leq P \quad (12)$$

605 for some constant P that depends on λ and the dataset \mathcal{D} .

606 We make use of the property from (4) to minimize the constraint term in (12) while keeping this
607 neural network equivalent to the original one. Specifically, let $\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \dots, \mathbf{W}^{(L)}, \mathbf{b}^{(L)}$ be the
608 parameters of an L -layer neural network.

$$f(x) = \mathbf{W}^{(L)} \sigma(\mathbf{W}^{(L-1)} \sigma(\dots \sigma(\mathbf{W}^{(1)} x + \mathbf{b}^{(1)}) \dots) + \mathbf{b}^{(L-1)}) + \mathbf{b}^{(L)},$$

609 which is equivalent to

$$f(x) = \alpha_L \tilde{\mathbf{W}}^{(L)} \sigma(\alpha_{L-1} \tilde{\mathbf{W}}^{(L-1)} \sigma(\dots \sigma(\alpha_1 \tilde{\mathbf{W}}^{(1)} x + \tilde{\mathbf{b}}^{(1)}) \dots) + \tilde{\mathbf{b}}^{(L-1)}) + \tilde{\mathbf{b}}^{(L)},$$

610 as long as $\alpha_\ell > 0$, $\prod_{\ell=1}^L \alpha_\ell = \prod_{\ell=1}^L \|\mathbf{W}^{(\ell)}\|_F$, where $\tilde{\mathbf{W}}^{(\ell)} := \frac{\mathbf{W}^{(\ell)}}{\|\mathbf{W}^{(\ell)}\|_F}$. By the AM-GM inequal-
611 ity, the ℓ_2 regularizer of the latter neural network is

$$\sum_{\ell=1}^L \|\alpha_\ell \tilde{\mathbf{W}}^{(\ell)}\|_F^2 = \sum_{\ell=1}^L \alpha_\ell^2 \geq L \left(\prod_{\ell=1}^L \alpha_\ell \right)^{2/L} = L \left(\prod_{\ell=1}^L \|\mathbf{W}^{(\ell)}\|_F \right)^{2/L}$$

612 and equality is reached when $\alpha_1 = \alpha_2 = \dots = \alpha_L$. In other word, in the problem (2), it suffices to
613 consider the network that satisfies

$$\|\mathbf{W}_j^{(1)}\|_F = \|\mathbf{W}_j^{(2)}\|_F = \dots = \|\mathbf{W}_j^{(L)}\|_F, \forall j \in [M], \ell \in [L]. \quad (13)$$

614 Using (4) again, one can find that the neural network is also equivalent to

$$f(x) = \sum_{j=1}^M a_j \bar{\mathbf{W}}^{(L)} \sigma(\bar{\mathbf{W}}_j^{(L-1)} \sigma(\dots \sigma(\bar{\mathbf{W}}_j^{(1)} x + \bar{\mathbf{b}}_j^{(1)}) \dots) + \bar{\mathbf{b}}_j^{(L-1)}) + \bar{\mathbf{b}}_j^{(L)},$$

615 where

$$\|\bar{\mathbf{W}}_j^{(\ell)}\|_F \leq \beta^{(\ell)}, a_j = \frac{\prod_{\ell=1}^L \|\mathbf{W}_j^{(\ell)}\|_F}{\prod_{\ell=1}^L \beta^{(\ell)}} = \frac{\|\mathbf{W}_j^{(1)}\|_F^L}{\prod_{\ell=1}^L \beta^{(\ell)}} = \frac{(\sum_{\ell=1}^L \|\mathbf{W}_j^{(\ell)}\|_F^2 / L)^{L/2}}{\prod_{\ell=1}^L \beta^{(\ell)}}, \quad (14)$$

616 where the last two equality comes from the assumption (13). Choosing $\beta^{(\ell)} = c_1 \sqrt{w}$ expect $\ell = 1$
617 where $\beta^{(1)} = c_1 \sqrt{d}$, and scaling $\bar{\mathbf{b}}^{(\ell)}$ accordingly and taking the constraint in (12) into (14) finishes
618 the proof. \square

619 D.2 Covering Number of Parallel Neural Networks

620 **Theorem 5.** *The covering number of the model defined in (5) apart from the bias in the last layer*
621 *satisfies*

$$\log \mathcal{N}(\mathcal{F}, \delta) \lesssim w^{2+2/(1-2/L)} L^2 \sqrt{d} P^{1-\frac{1}{1-2/L}} \delta^{-\frac{2/L}{1-2/L}} \log(wP'/\delta).$$

622

The proof relies on the covering number of each subnetwork in a parallel neural network (Lemma 11), observing that $|f(x)| \leq 2^{L-1} w^{L-1} \sqrt{d}$ under the condition in Lemma 11, and then apply Lemma 6. We argue that our choice of condition on $\|\mathbf{b}^{(\ell)}\|_2$ in Lemma 11 is sufficient to analyzing the model apart from the bias in the last layer, because it guarantees that $\sqrt{w} \|\mathbf{W}^{(\ell)} \mathcal{A}_{\ell-1}(x)\|_2 \leq \|\mathbf{b}^{(\ell)}\|_2$. This leads to

$$\|\mathbf{W}^{(\ell)} \mathcal{A}_{\ell-1}(x)\|_\infty \leq \|\mathbf{W}^{(\ell)} \mathcal{A}_{\ell-1}(x)\|_2 \leq \sqrt{w} \|\mathbf{b}^{(\ell)}\|_2 \leq \|\mathbf{b}^{(\ell)}\|_\infty$$

623 If this condition is not met, $\mathbf{W}^{(\ell)}\mathcal{A}_{\ell-1}(\mathbf{x}) + b^{(\ell)}$ is either always positive or always negative
624 for all feasible \mathbf{x} along at least one dimension. If $(\mathbf{W}^{(\ell)}\mathcal{A}_{\ell-1}(\mathbf{x}) + b^{(\ell)})_i$ is always negative,
625 one can replace $b^{(\ell)}_i$ with $-\max_{\mathbf{x}} \|\mathbf{W}^{(\ell)}\mathcal{A}_{\ell-1}(\mathbf{x})\|_{\infty}$ without changing the output of this model
626 for any feasible \mathbf{x} . If $(\mathbf{W}^{(\ell)}\mathcal{A}_{\ell-1}(\mathbf{x}) + b^{(\ell)})_i$ is always positive, one can replace $b^{(\ell)}_i$ with
627 $\max_{\mathbf{x}} \|\mathbf{W}^{(\ell)}\mathcal{A}_{\ell-1}(\mathbf{x})\|_{\infty}$, and adjust the bias in the next layer such that the output of this model
628 is not changed for any feasible \mathbf{x} . In either cases, one can replace the bias $b^{(\ell)}$ with another one with
629 smaller norm while keeping the model equivalent except the bias in the last layer.

630 **Lemma 11.** *Let $\mathcal{F} \subseteq \{f : \mathbb{R}^d \rightarrow \mathbb{R}\}$ denote the set of L -layer neural network (or a subnetwork in
631 a parallel neural network) with width w in each hidden layer. It has the form*

$$\begin{aligned} f(\mathbf{x}) &= \mathbf{W}^{(L)}\sigma(\mathbf{W}^{(L-1)}\sigma(\dots\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})\dots) + \mathbf{b}^{(L-1)}) + \mathbf{b}^{(L)}, \\ \mathbf{W}^{(1)} &\in \mathbb{R}^{w \times d}, \|\mathbf{W}^{(1)}\|_F \leq \sqrt{d}, \mathbf{b}^{(1)} \in \mathbb{R}^w, \|\mathbf{b}^{(1)}\|_2 \leq \sqrt{dw}, \\ \mathbf{W}^{(\ell)} &\in \mathbb{R}^{w \times w}, \|\mathbf{W}^{(\ell)}\|_F \leq \sqrt{w}, \mathbf{b}^{(\ell)} \in \mathbb{R}^w, \|\mathbf{b}^{(\ell)}\|_2 \leq 2^{\ell-1}w^{\ell-1}\sqrt{dw}, \quad \forall \ell = 2, \dots, L-1, \\ \mathbf{W}^{(L)} &\in \mathbb{R}^{1 \times w}, \|\mathbf{W}^{(L)}\|_F \leq \sqrt{w}, b^{(L)} = 0 \end{aligned} \tag{15}$$

632 and $\sigma(\cdot)$ is the ReLU activation function, the input satisfy $\|\mathbf{x}\|_2 \leq 1$, then the supremum norm
633 δ -covering number of \mathcal{F} obeys

$$\log \mathcal{N}(\mathcal{F}, \delta) \leq c_7 L w^2 \log(1/\delta) + c_8$$

634 where c_7 is a constant depending only on d , and c_8 is a constant that depend on d, w and L .

635 *Proof.* First study two neural networks which differ by only one layer. Let g_{ℓ}, g'_{ℓ} be two neural net-
636 works satisfying (15) with parameters $\mathbf{W}_1, \mathbf{b}_1, \dots, \mathbf{W}_L, \mathbf{b}_L$ and $\mathbf{W}'_1, \mathbf{b}'_1, \dots, \mathbf{W}'_L, \mathbf{b}'_L$ respectively.
637 Furthermore, the parameters in these two models are the same except the ℓ -th layer, which satisfy

$$\|\mathbf{W}_{\ell} - \mathbf{W}'_{\ell}\|_F \leq \epsilon, \|\mathbf{b}_{\ell} - \mathbf{b}'_{\ell}\|_2 \leq \tilde{\epsilon}.$$

638 Denote the model as

$$g_{\ell}(\mathbf{x}) = \mathcal{B}_{\ell}(\mathbf{W}_{\ell}\mathcal{A}_{\ell}(\mathbf{x}) + \mathbf{b}_{\ell}), g'_{\ell}(\mathbf{x}) = \mathcal{B}_{\ell}(\mathbf{W}'_{\ell}\mathcal{A}_{\ell}(\mathbf{x}) + \mathbf{b}'_{\ell})$$

639 where $\mathcal{A}_{\ell}(\mathbf{x}) = \sigma(\mathbf{W}_{\ell-1}\sigma(\dots\sigma(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1)\dots) + \mathbf{b}_{\ell-1})$ denotes the first $\ell - 1$ layers in the neural
640 network, and $\mathcal{A}_{\ell}(\mathbf{x}) = \mathbf{W}_L\sigma(\dots\sigma(\mathbf{W}_{\ell+1}\sigma(\mathbf{x}) + \mathbf{b}_{\ell+1})\dots) + \mathbf{b}_L$ denotes the last $L - \ell - 1$ layers,
641 with definition $\mathcal{A}_1(\mathbf{x}) = \mathbf{x}, \mathcal{B}_L(\mathbf{x}) = \mathbf{x}$.

642 Now focus on bounding $\|\mathcal{A}(\mathbf{x})\|$. Let $\mathbf{W} \in \mathbb{R}^{m \times m'}$, $\|\mathbf{W}\|_F \leq \sqrt{m'}$, $\mathbf{x} \in \mathbb{R}^{m'}$, $\mathbf{b} \in \mathbb{R}^m$, $\|\mathbf{b}\|_2 \leq$
643 \sqrt{m}

$$\begin{aligned} \|\sigma(\mathbf{W}\mathbf{x} + \mathbf{b})\|_2 &\leq \|\mathbf{W}\mathbf{x} + \mathbf{b}\|_2 \\ &\leq \|\mathbf{W}\|_2\|\mathbf{x}\|_2 + \|\mathbf{b}\|_2 \\ &\leq \|\mathbf{W}\|_F\|\mathbf{x}\|_2 + \|\mathbf{b}\|_2 \\ &\leq \sqrt{m'}\|\mathbf{x}\|_2 + \sqrt{m} \end{aligned}$$

644 where we make use of $\|\cdot\|_2 \leq \|\cdot\|_F$. Because of that,

$$\begin{aligned} \|\mathcal{A}_2(\mathbf{x})\|_2 &\leq \sqrt{d} + \sqrt{dw} \leq 2\sqrt{dw}, \\ \|\mathcal{A}_3(\mathbf{x})\|_2 &\leq \sqrt{w}\|\mathcal{A}_2(\mathbf{x})\|_2 + 2w\sqrt{dw} \leq 4w\sqrt{dw}, \\ &\dots \\ \|\mathcal{A}_{\ell}(\mathbf{x})\|_2 &\leq \sqrt{w}\|\mathcal{A}_{\ell-1}(\mathbf{x})\|_2 \leq 2\sqrt{dw}(2w)^{\ell-2}. \end{aligned} \tag{16}$$

645 Then focus on $\mathcal{B}(\mathbf{x})$. Let $\mathbf{W} \in \mathbb{R}^{m \times m'}$, $\|\mathbf{W}\|_F \leq \sqrt{m'}$, $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{m'}$, $\mathbf{b} \in \mathbb{R}^m$, $\|\mathbf{b}\|_2 \leq \sqrt{m}$.
646 Furthermore, $\|\mathbf{x} - \mathbf{x}'\|_2 \leq \epsilon$, then

$$\|\sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) - \sigma(\mathbf{W}\mathbf{x}' + \mathbf{b})\|_2 \leq \|\mathbf{W}(\mathbf{x} - \mathbf{x}')\|_2 \leq \|\mathbf{W}\|_F\|\mathbf{x} - \mathbf{x}'\|_2$$

647 which indicates that $\|\mathcal{B}(\mathbf{x}) - \mathcal{B}(\mathbf{x}')\|_2 \leq (\sqrt{w})^{L-\ell}\|\mathbf{x} - \mathbf{x}'\|_2$

648 Finally, for any $\mathbf{W}, \mathbf{W}' \in \mathbb{R}^{m \times m'}$, $\mathbf{x} \in \mathbb{R}^{m'}$, $\mathbf{b}, \mathbf{b}' \in \mathbb{R}^m$, one have

$$\begin{aligned} \|(\mathbf{W}\mathbf{x} + \mathbf{b}) - (\mathbf{W}'\mathbf{x} + \mathbf{b}')\|_2 &= \|(\mathbf{W} - \mathbf{W}')\mathbf{x} + (\mathbf{b} - \mathbf{b}')\|_2 \\ &\leq \|\mathbf{W} - \mathbf{W}'\|_2 \|\mathbf{x}\|_2 + \|\mathbf{b} - \mathbf{b}'\|_2. \\ &\leq \|\mathbf{W} - \mathbf{W}'\|_F \|\mathbf{x}\|_2 + \sqrt{m} \|\mathbf{b} - \mathbf{b}'\|_\infty. \end{aligned}$$

649 In summary,

$$\begin{aligned} |g_\ell(\mathbf{x}) - g'_\ell(\mathbf{x})| &= |\mathcal{B}_\ell(\mathbf{W}_\ell \mathcal{A}_\ell(\mathbf{x}) + \mathbf{b}_\ell) - \mathcal{B}_\ell(\mathbf{W}'_\ell \mathcal{A}_\ell(\mathbf{x}) + \mathbf{b}'_\ell)| \\ &\leq (\sqrt{w})^{L-\ell} \|(\mathbf{W}_\ell \mathcal{A}_\ell(\mathbf{x}) + \mathbf{b}_\ell) - (\mathbf{W}'_\ell \mathcal{A}_\ell(\mathbf{x}) + \mathbf{b}'_\ell)\|_2 \\ &\leq (\sqrt{w})^{L-\ell} (\|\mathbf{W}_\ell - \mathbf{W}'_\ell\|_F \|\mathcal{A}_\ell(\mathbf{x})\|_2 + \|\mathbf{b}_\ell - \mathbf{b}'_\ell\|_2) \\ &\leq 2^{(\ell-1)} w^{(L+\ell-3)/2} d^{1/2} \epsilon + w^{(L-\ell)/2} \bar{\epsilon} \end{aligned}$$

650 Let $f(x), f'(x)$ be two neural networks satisfying (15) with parameters $W_1, b_1, \dots, W_L, b_L$ and
651 $W'_1, b'_1, \dots, W'_L, b'_L$ respectively, and $\|W_\ell - W'_\ell\|_F \leq \epsilon_\ell$, $\|b_\ell - b'_\ell\|_F \leq \tilde{\epsilon}_\ell$. Further define f_ℓ be the
652 neural network with parameters $W_1, b_1, \dots, W_\ell, b_\ell, W'_{\ell+1}, b'_{\ell+1}, \dots, W'_L, b'_L$, then

$$\begin{aligned} |f(x) - f'(x)| &\leq |f(x) - f_1(x)| + |f_1(x) - f_2(x)| + \dots + |f_{L-1}(x) - f'(x)| \\ &\leq \sum_{\ell=1}^L 2^{(\ell-2)} d^{1/2} w^{(L+\ell-3)/2} \epsilon + w^{(L-\ell)/2} \bar{\epsilon} \end{aligned}$$

For any $\delta > 0$, one can choose

$$\epsilon_\ell = \frac{\delta}{2^\ell w^{(L+\ell-3)/2} d^{1/2}}, \tilde{\epsilon}_\ell = \frac{\delta}{2 w^{(L-\ell)/2}}$$

653 such that $|f(x) - f'(x)| \leq \delta$.

654 On the other hand, the ϵ -covering number of $\{\mathbf{W} \in \mathbb{R}^{m \times m'} : \|\mathbf{W}\|_F \leq \sqrt{m'}\}$ on Frobenius norm
655 is no larger than $(2\sqrt{m'}/\epsilon + 1)^{m \times m'}$, and the $\bar{\epsilon}$ -covering number of $\{\mathbf{b} \in \mathbb{R}^m : \|\mathbf{b}\|_2 \leq 1\}$ on
656 infinity norm is no larger than $(2/\bar{\epsilon} + 1)^m$. The entropy of this neural network can be bounded by

$$\log \mathcal{N}(f; \delta) \leq w^2 L \log(2^{L+1} w^{L-1} / \delta + 1) + wL \log(2^{L-1} w^{(L-1)/2} d^{1/2} / \delta + 1)$$

657

□

658 D.3 Covering Number of p -Norm Constrained Linear Combination

659 **Lemma 6.** $\log \mathcal{N}(\mathcal{G}, \delta) \lesssim k \log(1/\delta)$ for some finite c_3 , and for any $g \in \mathcal{G}, |a| \leq 1$, we have
660 $ag \in \mathcal{G}$. The covering number of $\mathcal{F} = \left\{ \sum_{i=1}^M a_i g_i \mid g_i \in \mathcal{G}, \|a\|_p^p \leq P, 0 < p < 1 \right\}$ for any $P > 0$
661 satisfies

$$\log \mathcal{N}(\mathcal{F}, \epsilon) \lesssim k P^{\frac{1}{1-p}} (\delta/c_3)^{-\frac{p}{1-p}} \log(c_3 P/\delta)$$

662 up to a double logarithmic factor.

663 *Proof.* Let ϵ be a positive constant. Without the loss of generality, we can sort the coefficients in
664 descending order in terms of their absolute values. There exists a positive integer \mathcal{M} (as a function
665 of ϵ), such that $|a_i| \geq \epsilon$ for $i \leq \mathcal{M}$, and $|a_i| < \epsilon$ for $i > \mathcal{M}$.

666 By definition, $\mathcal{M}\epsilon^p \leq \sum_{i=1}^{\mathcal{M}} |a_i|^p \leq P$ so $\mathcal{M} \leq P/\epsilon^p$, and $|a_i|^p \leq P, |a_i| \leq P^{1/p}$ for all i .
667 Furthermore,

$$\sum_{i>\mathcal{M}} |a_i| = \sum_{i>\mathcal{M}} |a_i|^p |a_i|^{1-p} < \sum_{i>\mathcal{M}} |a_i|^p \epsilon^{1-p} \leq P \epsilon^{1-p}$$

668 Let $\tilde{g}_i = \arg \min_{g \in \tilde{\mathcal{G}}} \|g - \frac{a_i}{P^{1/p}} g_i\|_\infty$ where $\tilde{\mathcal{G}}$ is the δ' -covering set of \mathcal{G} . By definition of the
669 covering set,

$$\begin{aligned} \left\| \sum_{i=1}^M a_i g_i(x) - \sum_{i=1}^{\mathcal{M}} P^{1/p} \tilde{g}_i(x) \right\|_\infty &\leq \left\| \sum_{i=1}^{\mathcal{M}} (a_i g_i(x) - P^{1/p} \tilde{g}_i(x)) \right\|_\infty + \left\| \sum_{i=\mathcal{M}+1}^M a_i g_i(x) \right\|_\infty \\ &\leq \mathcal{M} P^{1/p} \delta' + c_3 P \epsilon^{1-p}. \end{aligned}$$

(17)

670 Choosing

$$\epsilon = (\delta/2c_3P)^{\frac{1}{1-p}}, \delta' \approx P^{-\frac{1}{p(1-p)}} (\delta/2c_3)^{\frac{1}{1-p}}/2, \quad (18)$$

671 we have $\mathcal{M} \leq P^{\frac{1}{1-p}} (\delta/2c_3)^{-\frac{p}{1-p}}, \mathcal{M}P^{1/p}\delta' \leq \delta/2, c_3P\epsilon^{1-p} \leq \delta/2$, so (17) $\leq \delta$. One can
672 compute the covering number of \mathcal{F} by

$$\log \mathcal{N}(\mathcal{F}, \delta) \leq \mathcal{M} \log \mathcal{N}(\mathcal{G}, \delta') \lesssim k\mathcal{M} \log(1/\delta') \quad (19)$$

673 Taking (18) into (19) finishes the proof. \square

674 E Proof of Approximation Error

675 E.1 Approximation of Neural Networks to B-spline Basis Functions

676 **Proposition 7.** *Let $M_{m,k,s}$ be the B-spline of order m with scale 2^{-k} in each dimension and
677 position $s \in \mathbb{R}^d$: $M_{m,k,s}(\mathbf{x}) := M_m(2^k(\mathbf{x} - s))$, M_m is defined in (11). There exists a parallel
678 neural network that has the structure and satisfy the constraint in Proposition 4 for d -dimensional
679 input and one output, containing $M = O(m^d)$ subnetworks, each of which has width $w = O(d)$
680 and depth $L = O(\log(c(m, d)/\epsilon))$ for some constant w, c that depends only on m and d , denoted
681 as $\tilde{M}_m(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^d$, such that*

- 682 • $|\tilde{M}_{m,k,s}(\mathbf{x}) - M_{m,k,s}(\mathbf{x})| \leq \epsilon$, if $0 \leq 2^k(x_i - s_i) \leq m + 1, \forall i \in [d]$,
- 683 • $\tilde{M}_{m,k,s}(\mathbf{x}) = 0$, otherwise.
- 684 • The weights in the last layer satisfy $\|a\|_{2/L}^{2/L} \lesssim 2^k m^d e^{2md/L}$.

685 We follow the method developed in Yarotsky [49], Suzuki [39], while putting our attention on bound-
686 ing the Frobenius norm of the weights.

687 **Lemma 12** (Yarotsky [49, Proposition 3]). *: There exists a neural network with two-dimensional
688 input and one output $f_{\times}(x, y)$, with constant width and depth $O(\log(1/\delta))$, and the weight in each
689 layer is bounded by a global constant c_1 , such that*

- 690 • $|f_{\times}(x, y) - xy| \leq \delta, \forall 0 \leq x, y \leq 1$,
- 691 • $f_{\times}(x, y) = 0, \forall x = 0$ or $y = 0$.

692 We first prove a special case of Proposition 7 on the unscaled, unshifted B-spline basis function by
693 fixing $k = 0, s = 0$:

694 **Proposition 13.** *There exists a parallel neural network that has the structure and satisfy the con-
695 straint in Proposition 4 for d -dimensional input and one output, containing $M = \lceil (m + 1)/2 \rceil^d =$
696 $O(m^d)$ subnetworks, each of which has width $w = O(d)$ and depth $L = O(\log(c(m, d)/\epsilon))$ for
697 some constant w, c that depends only on m and d , denoted as $\tilde{M}_m(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^d$, such that*

- 698 • $|\tilde{M}_m(\mathbf{x}) - M_m(\mathbf{x})| \leq \epsilon$, if $0 \leq x_i \leq m + 1, \forall i \in [d]$, while $M_m(\cdot)$ denote m -th order
699 B-spline basis function, and c only depends on m and d .
- 700 • $\tilde{M}_m(\mathbf{x}) = 0$, if $x_i \leq 0$ or $x_i \geq m + 1$ for any $i \in [d]$.
- 701 • The weights in the last layer satisfy $\|a\|_{2/L}^{2/L} \lesssim m^d e^{2md/L}$.

702 *Proof.* We first show that one can use a neural network with constant width w_0 , depth $L \approx$
703 $\log(m/\epsilon_1)$ and bounded norm $\|W^{(1)}\|_F \leq O(\sqrt{d}), \|W^{(\ell)}\|_F \leq O(\sqrt{w}), \forall \ell = 2, \dots, L$ to
704 approximate truncated power basis function up to accuracy ϵ_1 in the range $[0, 1]$. Let $m =$
705 $\sum_{i=0}^{\lceil \log_2 m \rceil} m_i 2^i, m_i \in \{0, 1\}$ be the binary digits of m , and define $\bar{m}_j = \sum_{i=0}^j m_i, \gamma = \lceil \log_2 m \rceil$,

706 then for any x

$$\begin{aligned}
x_+^m &= x_+^{\tilde{m}_\gamma} \times (x_+^{2^\gamma})^{m_\gamma} \\
[x_+^{\tilde{m}_\gamma}, x_+^{2^\gamma}] &= [x_+^{\tilde{m}_{\gamma-1}} \times (x_+^{2^{\gamma-1}})^{m_{\gamma-1}}, x_+^{2^{\gamma-1}} \times x_+^{2^{\gamma-1}}] \\
&\dots \\
[x_+^{\tilde{m}_2}, x_+^4] &= [x_+^{\tilde{m}_1} \times (x_+^2)^{m_1}, x_+^2 \times x_+^2] \\
[x_+^{\tilde{m}_1}, x_+^2] &= [x_+^{\tilde{m}_0} \times x_+^{m_0}, x_+ \times x_+]
\end{aligned} \tag{20}$$

707 Notice that each line of equation only depends on the line immediately below. Replacing the
708 multiply operator \times with the neural network approximation shown in Lemma 12 demonstrates the
709 architecture of such neural network approximation. For any $x, y \in [0, 1]$, let $|f_\times(x, y) - xy| \leq$
710 δ , $|x - \tilde{x}| \leq \delta_1$, $|y - \tilde{y}| \leq \delta_2$, then $|f_\times(\tilde{x}, \tilde{y}) - xy| \leq \delta_1 + \delta_2 + \delta$. Taking this into (20) shows that
711 $\epsilon_1 \approx 2^\gamma \delta \approx m\delta$, where ϵ_1 is the upper bound on the approximate error to truncated power basis of
712 order m and δ is the approximation error to a single multiply operator as in Lemma 12.

713 A univariate B-spline basis can be expressed using truncated power basis, and observing that it is
714 symmetric around $(m+1)/2$:

$$\begin{aligned}
M_m(x) &= \frac{1}{m!} \sum_{j=1}^{m+1} (-1)^j \binom{m+1}{j} (x-j)_+^m \\
&= \frac{1}{m!} \sum_{j=1}^{\lceil (m+1)/2 \rceil} (-1)^j \binom{m+1}{j} (\min(x, m+1-x) - j)_+^m \\
&= \frac{((m+1)/2)^m}{m!} \sum_{j=1}^{\lceil (m+1)/2 \rceil} (-1)^j \binom{m+1}{j} \left(\frac{\min(x, m+1-x) - j}{(m+1)/2} \right)_+^m,
\end{aligned}$$

715 A multivariate (d -dimensional) B-spline basis function can be expressed as the product of truncated
716 power basis functions and thus can be decomposed as

$$\begin{aligned}
M_m(\mathbf{x}) &= \prod_{i=1}^d M_m(x_i) \\
&= \frac{((m+1)/2)^{md}}{(m!)^d} \prod_{i=1}^d \left(\sum_{j=1}^{\lceil (m+1)/2 \rceil} (-1)^j \binom{m+1}{j} \left(\frac{\min(x_i, m+1-x) - j}{(m+1)/2} \right)_+^m \right) \tag{21} \\
&= \frac{((m+1)/2)^{md}}{(m!)^d} \sum_{j_1, \dots, j_d=1}^{\lceil (m+1)/2 \rceil} \prod_{i=1}^d (-1)^{j_i} \binom{m+1}{j_i} \left(\frac{\min(x, m+1-x) - j_i}{(m+1)/2} \right)_+^m
\end{aligned}$$

717 Using Lemma 12, one can construct a parallel neural network containing $M = \lceil (m+1)/2 \rceil^d =$
718 $O(m^d)$ subnetworks, and each subnetwork corresponds to one polynomial term in (21). Using the
719 results above, the approximation of this constructed neural network can be bounded by

$$\frac{((m+1)/2)^{md}}{(m!)^d} \sum_{j_1, \dots, j_d=1}^{\lceil (m+1)/2 \rceil} \prod_{i=1}^d (-1)^{j_i} \binom{m+1}{j_i} \epsilon_1 \lesssim e^{md} \epsilon_1$$

720 where we applied Stirling's approximation and δ and ϵ_1 has the same definition as above. Choosing
721 $\delta = \frac{\epsilon}{d(e^{2m}\sqrt{m+1})}$, and recall $\epsilon_1 \approx m\delta$ proves the approximation error.

722 To bound the norm of the factors $\|a\|_{2/L}^{2/L}$, first observe that

$$\begin{aligned}
|a_{j_1, \dots, j_d}| &= \frac{((m+1)/2)^{md}}{(m!)^d} \frac{1}{(m+1)/2} \prod_{i=1}^d \binom{m+1}{j_i} \\
&\leq \frac{((m+1)/2)^{md}}{(m!)^d} \frac{2^{md}}{(m+1)/2} = O(e^{md})
\end{aligned}$$

723 where the first inequality is from $\binom{m+1}{j_i} \leq 2^{m+1}$, the last equality is from Stirling's approximation.
 724 Finally,

$$\|a\|_{2/L}^{2/L} \leq m^d \max_j |a_j|^{2/L} \lesssim m^d e^{2md/L}$$

725 which finishes the proof. \square

726 The proof of the Proposition 7 for general k, s follows by appending one more layer in the front, as
 727 we show below.

728 *Proof of Proposition 7.* Using the neural network proposed in Proposition 13, one can construct a
 729 neural network for approximating $M_{m,k,s}$ by adding one layer before the first layer:

$$\sigma(2^k \mathbf{I}_d \mathbf{x} - 2^k \mathbf{s})$$

730 The unused neurons in the first hidden layer is zero padded. The Frobenius norm of the weight is
 731 $2^k \|\mathbf{I}_d\|_F = 2^k \sqrt{d}$. Following the proof of Proposition 4, rescaling the weight in this layer by 2^{-k} ,
 732 and the weight matrix in the last layer by 2^k , and scaling the bias properly, one can verify that this
 733 neural network satisfy the statement. \square

734 E.2 Sparse approximation of Besov functions using B-spline wavelets

Proposition 8. *Let $\alpha - d/p > 1, r > 0$. For any function in Besov space $f_0 \in B_{p,q}^\alpha$ and any
 positive integer \bar{M} , there is an \bar{M} -sparse approximation using B-spline basis of order m satisfying
 $0 < \alpha < \min(m, m - 1 + 1/p)$: $\check{f}_{\bar{M}} = \sum_{i=1}^{\bar{M}} a_{k_i, s_i} M_{m, k_i, s_i}$ for any positive integer \bar{M} such that
 the approximation error is bounded as $\|\check{f}_{\bar{M}} - f_0\|_r \lesssim \bar{M}^{-\alpha/d} \|f_0\|_{B_{p,q}^\alpha}$, and the coefficients satisfy*

$$\|\{2^{k_i} a_{k_i, s_i}\}_{k_i, s_i}\|_p \lesssim \|f_0\|_{B_{p,q}^\alpha}.$$

735

736 The proof is divided into three steps:

- 737 1. Bound the 0-norm and the 1-norm of the coefficients of B-spline basis in order to approxi-
 738 mate an arbitrary function in Besov space up to any $\epsilon > 0$.
- 739 2. Bound p -norm of the coefficients of B-spline basis functions where $0 < p < 1$ using the
 740 results above .
- 741 3. Add the approximation to neural network to B-spline basis computed in Section 4.3.1 into
 742 Step 2.

743 *Proof.* Dũng [11, Theorem 3.1] Suzuki [39, Lemma 2] proposed an adaptive sampling recovery
 744 method that approximates a function in Besov space. The method is divided into two cases: when
 745 $p \geq r$, and when $p < r$.

746 When $p \geq r$, there exists a sequence of scalars $\lambda_j, \mathbf{j} \in P^d(\mu), P_d(\mu) := \{\mathbf{j} \in \mathbb{Z}^d : |j_i| \leq \mu, \forall i \in$
 747 $d\}$ for some positive μ , for arbitrary positive integer \bar{k} , the linear operator

$$Q_{\bar{k}}(f, \mathbf{x}) = \sum_{\mathbf{s} \in J(\bar{k}, m, d)} a_{\bar{k}, \mathbf{s}}(f) M_{\bar{k}, \mathbf{s}}(\mathbf{x}), \quad a_{\bar{k}, \mathbf{s}}(f) = \sum_{\mathbf{j} \in \mathbb{Z}^d, P^d(\mu)} \lambda_j \bar{f}(\mathbf{s} + 2^{-\bar{k}} \mathbf{j})$$

748 has bounded approximation error

$$\|f - Q_{\bar{k}}(f, x)\|_r \leq C 2^{-\alpha \bar{k}} \|f\|_{B_{p,q}^\alpha},$$

749 where \bar{f} is the extrapolation of f , $J(\bar{k}, m, d) := \{\mathbf{s} : 2^{\bar{k}} \mathbf{s} \in \mathbb{Z}^d, -m/2 \leq 2^{\bar{k}} s_i \leq 2^{\bar{k}} + m/2, \forall i \in$
 750 $[d]\}$. See Dũng [11, 2.6-2.7] for the detail of the extrapolation as well as references for options of
 751 sequence λ_j .

752 Furthermore, $Q_{\bar{k}}(f) \in B_{p,q}^\alpha$ so it can be decomposed in the form (10) with $M = \sum_{k=0}^{\bar{k}} (2^k + m -$
 753 $1)^d \lesssim 2^{\bar{k}d}$ components and $\|\{\tilde{c}_{k,s}\}_{k,s}\| \lesssim \|Q_{\bar{k}}(f)\|_{B_{p,q}^\alpha} \lesssim \|f\|_{B_{p,q}^\alpha}$ where $\tilde{c}_{k,s}$ is the coefficients of
 754 the decomposition of $Q_{\bar{k}}(f)$. Choosing $\bar{k} \approx \log_2 M/d$ leads to the desired approximation error.

755 On the other hand, when $p < r$, there exists a greedy algorithm that constructs

$$G(f) = Q_{\bar{k}}(f) + \sum_{k=\bar{k}+1}^{k^*} \sum_{j=1}^{n_k} c_{k,s_j}(f) M_{k,s_j}$$

756 where $\bar{k} \approx \log_2(M)$, $k^* = \lceil \epsilon^{-1} \log(\lambda M) \rceil + \bar{k} + 1$, $n_k = \lfloor \lambda M 2^{-\epsilon(k-\bar{k})} \rfloor$ for some $0 < \epsilon <$
 757 $\alpha/\delta - 1$, $\delta = d(1/p - 1/r)$, $\lambda > 0$, such that

$$\|f - G(f)\|_r \leq \bar{M}^{-\alpha/d} \|f\|_{B_{p,q}^\alpha}$$

758 and

$$\sum_{k=0}^{\bar{k}} (2^k + m - 1)^d + \sum_{k=\bar{k}+1}^{k^*} n_k \leq \bar{M}.$$

759 See Dũng [11, Theorem 3.1] for the detail.

760 Finally, since $\alpha - d/p > 1$,

$$\begin{aligned} \|\{2^{k_i} c_{k_i, s_i}\}_{k_i, s_i}\|_p &\leq \sum_{k=0}^{\bar{k}} 2^k \|\{c_{k_i, s_i}\}_{s_i}\|_p \\ &= \sum_{k=0}^{\bar{k}} 2^{(1-(\alpha-d/p))k} (2^{(\alpha-d/p)k} \|\{c_{k_i, s_i}\}_{s_i}\|_p) \\ &\lesssim \sum_{k=0}^{\bar{k}} 2^{(1-(\alpha-d/p))k} \|f\|_{B_{p,q}^\alpha} \\ &\approx \|f\|_{B_{p,q}^\alpha} \end{aligned} \tag{22}$$

761 where the first line is because for arbitrary vectors $\mathbf{a}_i, i \in [n]$, $\|\sum_{i=1}^n \mathbf{a}_i\|_p \leq \sum_{i=1}^n \|\mathbf{a}_i\|_p$, the
 762 third line is because the sequence norm of B-spline decomposition is equivalent to the norm in
 763 Besov space (see Section C.1). \square

764 Note that when $\alpha - d/p = 1$, the sequence norm (22) is bounded (up to a factor of constant) by
 765 $k^* \|f\|_{B_{p,q}^\alpha}$, which can be proven by following (22) except the last line. This adds a logarithmic term
 766 with respect to \bar{M} compared with the result in Proposition 8. This will add a logarithmic factor to
 767 the MSE. We will not focus on this case in this paper of simplicity.

768 E.3 Sparse approximation of Besov functions using Parallel Neural Networks

769 **Theorem 9.** *Under the same condition as Proposition 8, for any positive integer \bar{M} , any function*
 770 *in Besov space $f_0 \in B_{p,q}^\alpha$ can be approximated by a parallel neural network with no less than*
 771 *$O(m^d \bar{M})$ number of subnetworks satisfying:*

- 772 1. Each subnetwork has width $w = O(d)$ and depth L .
- 773 2. The weights in each layer satisfy $\|\bar{\mathbf{W}}_k^{(\ell)}\|_F \leq O(\sqrt{w})$ except the first layer $\|\bar{\mathbf{W}}_k^{(1)}\|_F \leq$
 774 $O(\sqrt{d})$,
- 775 3. The scaling factors have bounded $2/L$ -norm: $\|\{a_j\}\|_{2/L}^{2/L} \lesssim m^d e^{2md/L} \bar{M}^{1-2/(pL)}$.
- 776 4. The approximation error is bounded by

$$\|\tilde{f} - f_0\|_r \leq (c_4 \bar{M}^{-\alpha/d} + c_5 e^{-c_6 L}) \|f\|_{B_{p,q}^\alpha}$$

777 where c_4, c_5, c_6 are constants that depend only on m, d and p .

778 We first prove the following lemma.

Lemma 14. *For any $a \in \mathbb{R}^{\bar{M}}$, $0 < p' < p$, it holds that:*

$$\|a\|_{p'}^{p'} \leq \bar{M}^{1-p'/p} \|a\|_p^{p'}.$$

Proof.

$$\sum_i |a_i|^{p'} = \langle \mathbf{1}, |\mathbf{a}|^{p'} \rangle \leq \left(\sum_i 1 \right)^{1 - \frac{p'}{p}} \left(\sum_i (|a_i|^{p'})^{\frac{p}{p'}} \right)^{\frac{p'}{p}} = \bar{M}^{1 - \frac{p'}{p}} \|\mathbf{a}\|_p^{p'}$$

779 The first inequality uses a Holder's inequality with conjugate pair $\frac{p}{p'}$ and $1/(1 - \frac{p'}{p})$. \square

780 *Proof of Theorem 9.* Using Proposition 8, one can construct \bar{M} number of PNN each $O(m^d)$ sub-
781 networks according to Proposition 7, and in each PNN, such that each PNN represents one B-spline
782 basis function. The weights in the last layer of each PNN is scaled to match the coefficients in Propo-
783 sition 8. Taking p' in Lemma 14 as $2/L$ and combining with Proposition 7 finishes the proof. \square

784 F Proof of the Main Theorem

785 **Theorem 1 extended form.** For any fixed $\alpha - d/p > 1, q \geq 1, L \geq 3$, for any $f_0 \in B_{p,q}^\alpha$, given an
786 L -layer parallel neural network satisfying

- 787 • The width of each subnetwork is fixed and large enough: $w \gtrsim d$. See Theorem 9 for the
788 detail.
- 789 • The number of subnetworks is large enough: $M \gtrsim m^d n^{\frac{1-2/L}{2\alpha/d+1-2/(pL)}}$ where $m = \lceil \alpha - 1 \rceil$.

790 With proper choice of the parameter of weight decay λ , the solution \hat{f} parameterized by (2) satisfies

$$\text{MSE}(\hat{f}) = \tilde{O} \left(\left(\frac{w^{4-4/L} L^{2-4/L}}{n^{1-2/L}} \right)^{\frac{2\alpha/d}{2\alpha/d+1-2/(pL)}} + e^{-c_6 L} \right)$$

791 where \tilde{O} shows the scale up to a logarithmic factor, and c_6 is the constant defined in Theorem 9.

792 *Proof.* First recall the relationship between covering number (entropy) and estimation error:

793 **Proposition 15.** Let $\mathcal{F} \subseteq \{\mathbb{R}^d \rightarrow [-F, F]\}$ be a set of functions. Assume that \mathcal{F} can be decomposed
794 into two orthogonal spaces $\mathcal{F} = \mathcal{F}_\parallel \times \mathcal{F}_\perp$ where \mathcal{F}_\perp is an affine space with dimension of N . Let
795 $f_0 \in \{\mathbb{R}^d \rightarrow [-F, F]\}$ be the target function and \hat{f} be the least squares estimator in \mathcal{F} :

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n (y_i - f(x_i))^2, y_i = f_0(x_i) + \epsilon_i, \epsilon_i \sim \mathcal{N}(0, \sigma^2) \text{ i.i.d.},$$

796 then it holds that

$$\text{MSE}(\hat{f}) \leq \tilde{O} \left(\arg \min_{f \in \mathcal{F}} \text{MSE}(f) + \frac{N + \log \mathcal{N}(\mathcal{F}_\parallel, \delta) + 2}{n} + (F + \sigma)\delta \right).$$

797 The proof of Proposition 15 is deferred to the section below. We choose \mathcal{F} as the set of functions
798 that can be represented by a parallel neural network as stated, the (null) space $\mathcal{F}_\perp = \{f : f(\mathbf{x}) =$
799 $\text{constant}\}$ be the set of functions with constant output, which has dimension 1. This space captures
800 the bias in the last layer, while the other parameters contributes to the projection in \mathcal{F}_\parallel . See Sec-
801 tion D.2 for how we handle the bias in the other layers. One can find that \mathcal{F}_\parallel is the set of functions
802 that can be represented by a parallel neural network as stated, and further satisfy $\sum_{i=1}^n f(\mathbf{x}_i) = 0$.
803 Because $\mathcal{F}_\parallel \subseteq \mathcal{F}$, $\mathcal{N}(\mathcal{F}_\parallel, \delta) \leq \mathcal{N}(\mathcal{F}, \delta)$ for all $\delta > 0$, and the latter is studied in Theorem 5.

804 In Theorem 1, the width of each subnetwork is no less than what is required in Theorem 9, while the
805 depth and norm constraint are the same, so the approximation error is no more that that in Theorem 9.
806 Choosing $r = 2, p = 2/L$, and taking Theorem 5 and Theorem 9 into this Proposition 15, one gets

$$\text{MSE}(\hat{f}) \lesssim \bar{M}^{-2\alpha/d} + \frac{w^{2+2/(1-2/L)} L^2}{n} \bar{M}^{-\frac{1-2/(pL)}{1-2/L}} \delta^{-\frac{2/L}{1-2/L}} (\log(\bar{M}/\delta) + 3) + \delta,$$

where $\|f\|_{B_{p,q}^\alpha}$, m and d taken as constants. The stated MSE is obtained by choosing

$$\delta \approx \frac{w^{4-4/L} L^{2-4/L} \bar{M}^{1-2/(pL)}}{n^{1-2/L}}, \bar{M} \approx \left(\frac{n^{1-2/L}}{w^{4-4/L} L^{2-4/L}} \right)^{\frac{1}{2\alpha/d+1-2/(pL)}}$$

807 Note that there exists a weight decay parameter λ' such that the $(2/L)$ -norm of the coefficients
 808 of the parallel neural network satisfy that $\|\{a_j\}\|_{2/L}^{2/L} = m^d e^{2md/L} \|\{\tilde{a}_{j,\bar{M}}\}\|_{2/L}^{2/L}$ where $\{\tilde{a}_{j,\bar{M}}\}$
 809 is the coefficient of the particular \bar{M} -sparse approximation, although $\{a_j\}$ is not necessarily \bar{M}
 810 sparse. Empirically, one only need to guarantee that during initialization, the number of subnetworks
 811 $M \geq \bar{M}$ such that the \bar{M} -sparse approximation is feasible, thus the approximation error bound
 812 from Theorem 9 can be applied. Theorem 9 also says that $\|\{a_j\}\|_{2/L}^{2/L} = m^d e^{2md/L} \|\{\tilde{a}_{j,\bar{M}}\}\|_{2/L}^{2/L} \lesssim$
 813 $\bar{M}^{1-2/pL}$, thus we can apply the covering number bound from Theorem 5 with $P' = \bar{M}^{1-2/pL}$.
 814 Finally, if λ is optimally chosen, then it achieves a smaller MSE than this particular λ' , which has
 815 been proven to be no more than $O(\bar{M}^{-\alpha/d})$ and completes the proof. \square

816

817 *Proof of Proposition 15.* For any function $f \in \mathcal{F}$, define $f_\perp = \arg \min_{h \in \mathcal{F}_\perp} \sum_{i=1}^n (f(\mathbf{x}_i) -$
 818 $h(\mathbf{x}_i))^2$ be the projection of f to \mathcal{F}_\perp , and define $f_\parallel = f - f_\perp$ be the projection to the orthogo-
 819 nal complement. Note that f_\parallel is not necessarily in \mathcal{F}_\parallel . However, if $f \in \mathcal{F}$, then $f_\parallel \in \mathcal{F}_\parallel$. $y_{i\perp}$ and
 820 $y_{i\parallel}$ are defined by creating a function f_y such that $f_y(\mathbf{x}_i) = y_i, \forall i$, e.g. via interpolation. Because
 821 \mathcal{F}_\parallel and \mathcal{F}_\perp are orthonormal, the empirical loss and population loss can be decomposed in the same
 822 way:

$$\begin{aligned} L_\parallel(f) &= \frac{1}{n} \sum_{i=1}^n (f_\parallel(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 + \frac{n-N}{n} \sigma^2, & L_\perp(f) &= \frac{1}{n} \sum_{i=1}^n (f_\perp(\mathbf{x}_i) - f_{0\perp}(\mathbf{x}_i))^2 + \frac{N}{n} \sigma^2, \\ \hat{L}_\parallel(f) &= \frac{1}{n} \sum_{i=1}^n (f_\parallel(\mathbf{x}_i) - y_{i\parallel})^2, & \hat{L}_\perp(f) &= \frac{1}{n} \sum_{i=1}^n (f_\perp(\mathbf{x}_i) - y_{i\perp})^2, \\ MSE_\parallel(f) &= \mathbb{E}_{\mathcal{D}} \left[\frac{1}{n} \sum_{i=1}^n (f_\parallel(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 \right], & MSE_\perp(f) &= \mathbb{E}_{\mathcal{D}} \left[\frac{1}{n} \sum_{i=1}^n (f_\perp(\mathbf{x}_i) - f_{0\perp}(\mathbf{x}_i))^2 \right], \end{aligned}$$

823 such that $L(f) = L_\parallel(f) + L_\perp(f), \hat{L}(f) = \hat{L}_\parallel(f) + \hat{L}_\perp(f)$. This can be verified by de-
 824 composing \hat{f}, f_0 and y into two orthogonal components as shown above, and observing that
 825 $\sum_{i=1}^n f_{1\perp}(\mathbf{x}_i) f_{2\parallel}(\mathbf{x}_i) = 0, \forall f_1, f_2$.

826 **First prove the following claim**

827 **Claim 16.** Assume that $\hat{f} = \arg \min_{f \in \mathcal{F}} \hat{L}(f)$ is the empirical risk minimizer. Then $\hat{f}_\perp =$
 828 $\arg \min_{f \in \mathcal{F}_\perp} \hat{L}_\perp(f), \hat{f}_\parallel = \arg \min_{f \in \mathcal{F}_\parallel} \hat{L}_\parallel(f)$, where \hat{f}_\perp is the projections of \hat{f} in \mathcal{F}_\perp , and
 829 $\hat{f}_\parallel = \hat{f} - \hat{f}_\perp$ respectively.

830 *Proof.* Since $\hat{f} \in \mathcal{F}$, by definition $\hat{f}_\parallel \in \mathcal{F}_\parallel$. Assume that there exist $\hat{f}'_\perp, \hat{f}'_\parallel$, and either $\hat{L}_\perp(\hat{f}'_\perp) <$
 831 $\hat{L}_\perp(\hat{f}_\perp)$, or $\hat{L}_\parallel(\hat{f}'_\parallel) < \hat{L}_\parallel(\hat{f}_\parallel)$. Then

$$\begin{aligned} \hat{L}(\hat{f}') &= \hat{L}(\hat{f}'_\perp + \hat{f}'_\parallel) = \hat{L}_\perp(\hat{f}'_\perp + \hat{f}'_\parallel) + \hat{L}_\parallel(\hat{f}'_\perp + \hat{f}'_\parallel) = \hat{L}_\perp(\hat{f}'_\perp) + \hat{L}_\parallel(\hat{f}'_\parallel) \\ &< \hat{L}_\perp(\hat{f}_\perp) + \hat{L}_\parallel(\hat{f}_\parallel) = \hat{L}_\perp(\hat{f}_\perp + \hat{f}_\parallel) + \hat{L}_\parallel(\hat{f}_\perp + \hat{f}_\parallel) = \hat{L}(\hat{f}) \end{aligned}$$

832 which shows that \hat{f} is not the minimizer of $\hat{L}(f)$ and violates the assumption. \square

833

834 **Then we bound** $MSE_{\perp}(f)$. We convert this part into a finite dimension least square problem:

$$\begin{aligned}
\hat{f}_{\perp} &= \arg \min_{f \in \mathcal{F}_{\perp}} \hat{L}_{\perp}(f) \\
&= \arg \min_{f \in \mathcal{F}_{\perp}} \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - f_{0\perp}(\mathbf{x}_i) - \epsilon_{i\perp})^2 \\
&= \arg \min_{f \in \mathcal{F}_{\perp}} \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - f_{0\perp}(\mathbf{x}_i) - \epsilon_{i\perp})^2 + \epsilon_{i\parallel}^2 \\
&= \arg \min_{f \in \mathcal{F}_{\perp}} \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - f_{0\perp}(\mathbf{x}_i) - \epsilon_{i\perp} - \epsilon_{i\parallel})^2 \\
&= \arg \min_{f \in \mathcal{F}_{\perp}} \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - f_{0\perp}(\mathbf{x}_i) - \epsilon_i)^2
\end{aligned}$$

835 The forth line comes from our assumption that \mathcal{F}_{\perp} is orthogonal to \mathcal{F}_{\parallel} , so $\forall f \in \mathcal{F}_{\perp}, f + f_{0\perp} + \epsilon_{\perp}$
836 is orthogonal to ϵ_{\parallel} .

837 Let the basis function of \mathcal{F}_{\perp} be h_1, h_2, \dots, h_N , the above problem can be reparameterized as

$$\arg \min_{\boldsymbol{\theta} \in \mathbb{R}^N} \frac{1}{n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|^2$$

838 where $\mathbf{X} \in \mathbb{R}^{n \times N} : X_i = h_j(\mathbf{x}_i), \mathbf{y} = \mathbf{y}_{0\perp} + \boldsymbol{\epsilon}, \mathbf{y}_{0\perp} = [f_{0\perp}(x_1), \dots, f_{0\perp}(x_n)], \boldsymbol{\epsilon} = [\epsilon_1, \dots, \epsilon_n]$.
839 This problem has a closed-form solution

$$\boldsymbol{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

840 Observe that $f_{0\perp} \in \mathcal{F}_{\perp}$, let $\mathbf{y}_{0\perp} = \mathbf{X}\boldsymbol{\theta}^*$, The MSE of this problem can be computed by

$$\begin{aligned}
L(\hat{f}_{\perp}) &= \frac{1}{n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}_{0\perp}\|^2 = \frac{1}{n} \|\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\boldsymbol{\theta}^* + \boldsymbol{\epsilon}) - \mathbf{X}\boldsymbol{\theta}^*\|^2 \\
&= \frac{1}{n} \|\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}\|^2
\end{aligned}$$

841 Observing that $\Pi := \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is an idempotent and independent projection whose rank is
842 N , and that $\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T] = \sigma^2 \mathbf{I}$, we get

$$MSE_{\perp}(\hat{f}_{\perp}) = \mathbb{E}[L(\hat{f}_{\perp})] = \frac{1}{n} \|\Pi\boldsymbol{\epsilon}\|^2 = \frac{1}{n} \text{tr}(\Pi\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T) = \frac{\sigma^2}{n} \text{tr}(\Pi)$$

843 which concludes that

$$MSE_{\perp}(\hat{f}) = O\left(\frac{N}{n} \sigma^2\right). \quad (23)$$

844 See also [19, Proposition 1].

845 **Next we study** $MSE_{\parallel}(\hat{f})$. Denote $\tilde{\sigma}_{\parallel}^2 = \frac{1}{n} \sum_{i=1}^n \epsilon_{i\parallel}^2, E = \max_i |\epsilon_i|$. Using Jensen's inequality and
846 union bound, we have

$$\exp(t\mathbb{E}[E]) \leq \mathbb{E}[\exp(tE)] = \mathbb{E}[\max \exp(t|\epsilon_i|)] \leq \sum_{i=1}^n \mathbb{E}[\exp(t|\epsilon_i|)] \leq 2n \exp(t^2 \sigma^2 / 2)$$

847 Taking expectation over both sides, we get

$$\mathbb{E}[E] \leq \frac{\log 2n}{t} + \frac{t\sigma^2}{2}$$

848 maximizing the right hand side over t yields

$$\mathbb{E}[E] \leq \sigma \sqrt{2 \log 2n}.$$

849 Let $\tilde{\mathcal{F}}_{\parallel}$ be the covering set of $\mathcal{F}_{\parallel} = \{f_{\parallel} : f \in \mathcal{F}\}$. For any $\tilde{f}_{\parallel} \in \tilde{\mathcal{F}}_{\parallel}$,

$$\begin{aligned}
L_{\parallel}(f_j) - \hat{L}_{\parallel}(f_j) &= \frac{1}{n} \sum_{i=1}^n (f_{j\parallel}(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 - \frac{1}{n} \sum_{i=1}^n (\tilde{f}_{\parallel}(\mathbf{x}_i) - y_{i\parallel})^2 + \frac{n-N}{n} \sigma^2 \\
&= \frac{1}{n} \sum_{i=1}^n \epsilon_i (2\tilde{f}_{\parallel}(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i) - y_{i\parallel}) + \frac{n-N}{n} \sigma^2 \\
&= \frac{1}{n} \sum_{i=1}^n \epsilon_i (2\tilde{f}_{\parallel}(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i) - y_{i\parallel}) + \frac{n-N}{n} \sigma^2 \\
&= \frac{1}{n} \sum_{i=1}^n \epsilon_i (2\tilde{f}_{\parallel}(\mathbf{x}_i) - 2f_{0\parallel}(\mathbf{x}_i)) + \frac{n-N}{n} \sigma^2 - \tilde{\sigma}_{\parallel}^2
\end{aligned}$$

850 The first term can be bounded using Bernstein's inequality: let $h_i = \epsilon_i(f_{j\parallel}(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))$, by
851 definition $|h_i| \leq 2EF$,

$$\begin{aligned}
\text{Var}[h_i] &= \mathbb{E}[\epsilon_i^2 (\tilde{f}_{\parallel}(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2] \\
&= (\tilde{f}_{\parallel}(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 \mathbb{E}[\epsilon_i^2] \\
&= (\tilde{f}_{\parallel}(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 \sigma^2
\end{aligned}$$

852 using Bernstein's inequality, for any $\tilde{f}_{\parallel} \in \tilde{\mathcal{F}}_{\parallel}$, with probably at least $1 - \delta_p$,

$$\begin{aligned}
\frac{1}{n} \sum_{i=1}^n \epsilon_i (2\tilde{f}_{\parallel}(\mathbf{x}_i) - 2f_{0\parallel}(\mathbf{x}_i)) &= \frac{2}{n} \sum_{i=1}^n h_i \\
&\leq \frac{2}{n} \sqrt{2 \sum_{i=1}^n (\tilde{f}_{\parallel}(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 \sigma^2 \log(1/\delta_p) + \frac{8EF \log(1/\delta_p)}{3n}} \\
&= 2 \sqrt{\left(L_{\parallel}(\tilde{f}_{\parallel}) - \frac{n-N}{n} \sigma^2\right) \frac{2\sigma^2 \log(1/\delta_p)}{n} + \frac{8EF \log(1/\delta_p)}{3n}} \\
&\leq \epsilon \left(L_{\parallel}(\tilde{f}_{\parallel}) - \frac{n-N}{n} \sigma^2\right) + \frac{8\sigma^2 \log(1/\delta_p)}{n\epsilon} + \frac{8EF \log(1/\delta_p)}{3n}
\end{aligned}$$

853 the last inequality holds true for all $\epsilon > 0$. The union bound shows that with probably at least $1 - \delta$,
854 for all $\tilde{f}_{\parallel} \in \tilde{\mathcal{F}}_{\parallel}$,

$$\begin{aligned}
L_{\parallel}(\tilde{f}_{\parallel}) - \hat{L}_{\parallel}(\tilde{f}_{\parallel}) &\leq \epsilon \left(L_{\parallel}(\tilde{f}_{\parallel}) - \frac{n-N}{n} \sigma^2\right) + \frac{8\sigma^2 \log(\mathcal{N}(\mathcal{F}_{\parallel}, \delta)/\delta_p)}{n\epsilon} + \frac{8EF \log(\mathcal{N}(\mathcal{F}_{\parallel}, \delta)/\delta_p)}{3n} \\
&\quad + \frac{n-N}{n} \sigma^2 - \tilde{\sigma}_{\parallel}^2.
\end{aligned}$$

855 By rearranging the terms and using the definition of $L(\tilde{f}_{\parallel})$, we get

$$(1 - \epsilon) \left(L_{\parallel}(\tilde{f}_{\parallel}) - \frac{n-N}{n} \sigma^2\right) \leq \hat{L}_{\parallel}(\tilde{f}_{\parallel}) + \frac{8\sigma^2 \log(\mathcal{N}(\mathcal{F}_{\parallel}, \delta)/\delta_p)}{n\epsilon} + \frac{8EF \log(\mathcal{N}(\mathcal{F}_{\parallel}, \delta)/\delta_p)}{3n} - \tilde{\sigma}_{\parallel}^2.$$

856 Taking the expectation (over \mathcal{D}) on both sides, and notice that $\mathbb{E}[\tilde{\sigma}_{\parallel}^2] = \frac{n-N}{n} \sigma^2$. Furthermore, for
857 any random variable X , $\mathbb{E}[X] = \int_{-\infty}^{\infty} x dP(X \leq x)$, we get

$$\begin{aligned}
&\max_{\tilde{f}_{\parallel} \in \tilde{\mathcal{F}}_{\parallel}} \left((1 - \epsilon) \text{MSE}_{\parallel}(\tilde{f}_{\parallel}) - \mathbb{E}[\hat{L}_{\parallel}(\tilde{f}_{\parallel})] \right) \\
&\leq \left(\frac{8\sigma^2}{n\epsilon} + \frac{8F\sigma\sqrt{2\log 2n}}{3n} \right) \left(\log \mathcal{N}(\mathcal{F}_{\parallel}, \delta) - \int_{\delta=0}^1 \log(\delta_p) d\delta_p \right) - \frac{n-N}{n} \sigma^2 \quad (24) \\
&= \left(\frac{8\sigma^2}{n\epsilon} + \frac{8F\sigma\sqrt{2\log 2n}}{3n} \right) (\log \mathcal{N}(\mathcal{F}_{\parallel}, \delta) + 1) - \frac{n-N}{n} \sigma^2.
\end{aligned}$$

858 where the integration can be computed by replacing δ with e^x . Though it is not integrable under
859 Riemann integral, it is integrable under Lebesgue integration.

860 Similarly, let $\check{f}_{\parallel} = \arg \min_{f \in \mathcal{F}_{\parallel}} L_{\parallel}(f)$,

$$L_{\parallel}(\check{f}_{\parallel}) - \hat{L}_{\parallel}(\check{f}_{\parallel}) = \frac{1}{n} \sum_{i=1}^n \epsilon_i (2\check{f}_{\parallel}(\mathbf{x}_i) - 2f_{0\parallel}(\mathbf{x}_i)) + \frac{n-N}{n} \sigma^2 - \tilde{\sigma}_{\parallel}^2$$

861 with probably at least $1 - \delta_q$, for any $\epsilon > 0$,

$$\begin{aligned} -\frac{1}{n} \sum_{i=1}^n \epsilon_i (2\check{f}_{\parallel}(\mathbf{x}_i) - 2f_{0\parallel}(\mathbf{x}_i)) &\leq \epsilon \left(L_{\parallel}(\check{f}_{\parallel}) - \frac{n-N}{n} \sigma^2 \right) + \frac{8\sigma^2 \log(1/\delta_p)}{n\epsilon} + \frac{8EF \log(1/\delta_p)}{3n}, \\ \hat{L}_{\parallel}(\check{f}_{\parallel}) &\leq (1 + \epsilon) \left(L_{\parallel}(\check{f}_{\parallel}) - \frac{n-N}{n} \sigma^2 \right) + \frac{8\sigma^2 \log(1/\delta_p)}{n\epsilon} + \frac{8EF \log(1/\delta_q)}{3n} + \tilde{\sigma}_{\parallel}^2. \end{aligned}$$

862 Taking the expectation on both sides,

$$\mathbb{E}[\hat{L}_{\parallel}(\check{f}_{\parallel})] \leq (1 + \epsilon) \text{MSE}_{\parallel}(\check{f}_{\parallel}) + \frac{8\sigma^2}{n\epsilon} + \frac{8F\sigma\sqrt{2\log 2n}}{3n} + \frac{n-N}{n} \sigma^2. \quad (25)$$

863 Finally, let $\hat{f}_* := \arg \min_{f \in \tilde{\mathcal{F}}_{\parallel}} \sum_{i=1}^n (\hat{f}_{\parallel}(\mathbf{x}_i) - f(\mathbf{x}_i))^2$ be the projection of \hat{f}_{\parallel} in its δ -covering
864 space,

$$\begin{aligned} \text{MSE}_{\parallel}(\hat{f}_{\parallel}) &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n (\hat{f}_{\parallel}(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 \right] \\ &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n (\hat{f}_*(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 + \frac{1}{n} \sum_{i=1}^n (\hat{f}_{\parallel}(\mathbf{x}_i) - \hat{f}_*(\mathbf{x}_i))(\hat{f}_{\parallel}(\mathbf{x}_i) + \hat{f}_*(\mathbf{x}_i) - 2f_{0\parallel}(\mathbf{x}_i)) \right] \\ &\leq \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n (\hat{f}_*(\mathbf{x}_i) - f_{0\parallel}(\mathbf{x}_i))^2 \right] + 4F\delta \\ &= \text{MSE}_{\parallel}(\hat{f}_*(\mathbf{x}_i)) + 4F\delta, \end{aligned}$$

865 and similarly

$$\hat{L}_{\parallel}(\hat{f}_*) \leq \hat{L}_{\parallel}(\hat{f}_{\parallel}) + (4F + 2E)\delta. \quad (26)$$

866 We can conclude that

$$\begin{aligned} \text{MSE}_{\parallel}(\hat{f}_{\parallel}) &\leq \frac{1}{1-\epsilon} \left(\mathbb{E}[\hat{L}_{\parallel}(\hat{f}_*)] + \left(\frac{8\sigma^2}{n\epsilon} + \frac{8F\sigma\sqrt{2\log 2n}}{3n} \right) (\log \mathcal{N}(\mathcal{F}_{\parallel}, \delta) + 1) - \frac{n-N}{n} \sigma^2 \right) \\ &\quad + 4F\delta \\ &\leq \frac{1}{1-\epsilon} \left(\mathbb{E}[\hat{L}_{\parallel}(\hat{f}_{\parallel})] + (4F + \sigma\sqrt{8\log 2n})\delta \right) \\ &\quad + \left(\frac{8\sigma^2}{n\epsilon} + \frac{8F\sigma\sqrt{2\log 2n}}{3n} \right) (\log \mathcal{N}(\mathcal{F}_{\parallel}, \delta) + 1) - \frac{n-N}{n} \sigma^2 + 4F\delta \\ &\leq \frac{1}{1-\epsilon} \left(\mathbb{E}[\hat{L}_{\parallel}(\check{f}_{\parallel})] + (4F + \sigma\sqrt{8\log 2n})\delta \right) \\ &\quad + \left(\frac{8\sigma^2}{n\epsilon} + \frac{8F\sigma\sqrt{2\log 2n}}{3n} \right) (\log \mathcal{N}(\mathcal{F}_{\parallel}, \delta) + 1) - \frac{n-N}{n} \sigma^2 + 4F\delta \\ &\leq \frac{1+\epsilon}{1-\epsilon} \text{MSE}_{\parallel}(\check{f}_{\parallel}) + \frac{1}{n} \left(\frac{8\sigma^2}{\epsilon} + \frac{8F\sigma\sqrt{2\log 2n}}{3} \right) \left(\frac{\log \mathcal{N}(\mathcal{F}_{\parallel}, \delta) + 2}{1-\epsilon} \right) \\ &\quad + \left(4F + \frac{4F + \sigma\sqrt{8\log 2n}}{1-\epsilon} \right) \delta, \end{aligned}$$

867 where the first line comes from (24), and second comes from (26), the third line is because
868 $\hat{f}_{\parallel} = \arg \min_{f \in \mathcal{F}_{\parallel}} \hat{L}_{\parallel}(f)$, and the last line comes from (25). We also use that fact that $\hat{L}_{\parallel}(\hat{f}) \leq$
869 $\hat{L}_{\parallel}(f), \forall f$. Noticing that $\text{MSE}(\hat{f}) = \text{MSE}_{\parallel}(\hat{f}) + \text{MSE}_{\perp}(\hat{f})$, combining this with (23) finishes the
870 proof. \square

871 **G Detailed experimental setup**

872 **G.1 Target Functions**

873 The doppler function used in Figure 2(d)-(f) is

$$f(x) = \sin(4/(x + 0.01)) + 1.5.$$

874 The “vary” function used in Figure 2(g)-(i) is

$$\begin{aligned} f(x) = & M_1(x/0.01) + M_1((x - 0.02)/0.02) + M_1((x - 0.06)/0.03) \\ & + M_1((x - 0.12)/0.04) + M_3((x - 0.2)/0.02) + M_3((x - 0.28)/0.04) \\ & + M_3((x - 0.44)/0.06) + M_3((x - 0.68)/0.08), \end{aligned}$$

875 where M_1, M_3 are first and third order Cardinal B-spline bases functions respectively. We uni-
 876 formly take 256 samples from 0 to 1 in the piecewise cubic function experiment, and uniformly
 877 1000 samples from 0 to 1 in the doppler function and “vary” function experiment. We add zero
 878 mean independent (white) Gaussian noise to the observations. The standard derivation of noise is
 879 0.4 in the doppler function experiment and 0.1 in the “vary” function experiment.

880 **G.2 Training/Fitting Method**

881 In the piecewise polynomial function (“vary”) experiment, the depth of the PNN $L = 10$, the width
 882 of each subnetwork $w = 10$, and the model contains $M = 500$ subnetworks. The depth of NN is also
 883 10, and the width is 200 such that the NN and PNN have almost the same number of parameters. In
 884 the doppler function experiment, the depth of the PNN $L = 12$, the width of each subnetwork $w =$
 885 10, and the model contains $M = 2000$ subnetworks, because this problem requires a more complex
 886 model to fit. The depth of NN is 12, and the width is 400. We used Adam optimizer with learning rate
 887 of 10^{-3} . We first train the neural network layer by layer without weight decay. Specifically, we start
 888 with a two-layer neural network with the same number of subnetworks and the same width in each
 889 subnetwork, then train a three layer neural network by initializing the first layer using the trained
 890 two layer one, until the desired depth is reached. After that, we turn the weight decay parameter and
 891 train it until convergence. In both trend filtering and smoothing spline experiment, the order is 3,
 892 and in wavelet denoising experiment, we use sym4 wavelet with soft thresholding. We implement
 893 the trend filtering problem according to Tibshirani [40] using CVXPY, and use MOSEK to solve
 894 the convex optimization problem. We directly call R function *smooth.spline* to solve smoothing
 895 spline.

896 **G.3 Post Processing**

897 The degree of freedom of smoothing spline is returned by the solver in R, which is rounded to the
 898 nearest integer when plotting. To estimate the degree of freedom of trend filtering, for each choice
 899 of λ , we repeated the experiment for 10 times and compute the average number of nonzero knots as
 900 estimated degree of freedom. For neural networks, we use the definition [41]:

$$2\sigma^2 df = \mathbb{E}\|\mathbf{y}' - \hat{\mathbf{y}}\|_2^2 - \mathbb{E}\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 \quad (27)$$

901 where df denotes the degree of freedom, σ^2 is the variance of the noise, \mathbf{y} are the labels, $\hat{\mathbf{y}}$ are
 902 the predictions and \mathbf{y}' are independent copy of \mathbf{y} . We find that estimating (27) directly by sampling
 903 leads to large error when the degree of freedom is small. Instead, we compute

$$2\sigma^2 \hat{df} = \hat{\mathbb{E}}\|\mathbf{y}_0 - \hat{\mathbf{y}}\|_2^2 - \hat{\mathbb{E}}\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 + \hat{\mathbb{E}}\|\mathbf{y} - \bar{y}_0\|_2^2 - \|\mathbf{y}_0 - \bar{y}_0\|_2^2 \quad (28)$$

904 where \hat{df} is the estimated degree of freedom, $\hat{\mathbb{E}}$ denotes the empirical average (sample mean), \mathbf{y}_0 is
 905 the target function and \bar{y}_0 is the mean of the target function in its domain.

906 **Proposition 17.** *The expectation of (28) over the dataset \mathcal{D} equals (27).*

Proof.

$$\begin{aligned}
2\sigma^2 \hat{d}f &= \mathbb{E}_{\mathcal{D}} [\hat{\mathbb{E}} \|\mathbf{y}_0 - \hat{\mathbf{y}}\|_2^2 - \hat{\mathbb{E}} \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 + \hat{\mathbb{E}} \|\mathbf{y} - \bar{y}_0\|_2^2 - \|\mathbf{y}_0 - \bar{y}_0\|_2^2] \\
&= \mathbb{E} \|\mathbf{y}_0 - \hat{\mathbf{y}}\|_2^2 - \mathbb{E} \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 + \mathbb{E}_{\mathcal{D}} [\hat{\mathbb{E}} [(\mathbf{y} - \mathbf{y}_0)(\mathbf{y} + \mathbf{y}_0 - 2\bar{y}_0)]] \\
&= \mathbb{E} \|\mathbf{y}_0 - \hat{\mathbf{y}}\|_2^2 - \mathbb{E} \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 + \mathbb{E} \left[\sum_{i=1}^n \epsilon_i (2y_i + \epsilon_i - 2\bar{y}_0) \right] \\
&= \mathbb{E} \|\mathbf{y}_0 - \hat{\mathbf{y}}\|_2^2 - \mathbb{E} \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 + n\sigma^2 \\
&= \mathbb{E} \|\mathbf{y}' - \hat{\mathbf{y}}\|_2^2 - \mathbb{E} \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2
\end{aligned}$$

907 where \mathcal{D} denotes the dataset. In the third line, we make use of the fact that $\mathbb{E}[\epsilon_i] = 0$, $\mathbb{E}[\epsilon_i^2] = \sigma^2$,
908 and in the last line, we make use of $\mathbb{E}[\epsilon_i'] = 0$, $\mathbb{E}[\epsilon_i'^2] = \sigma^2$, and ϵ_i' are independent of y_i and $y_{0,i}$. \square

909 One can easily check that a “zero predictor” (a predictor that always predict \bar{y}_0 , and it always predicts
910 0 if the target function has zero mean) always has an estimated degree of freedom of 0.

911 In Figure 2(h)(i), we take the minimum MSE over different choices of λ , and plot the average over
912 10 runs. Due to optimization issue, sometimes the neural networks are stuck at bad local minima
913 and the empirical loss is larger than the global minimum by orders of magnitude. To deal with this
914 problem, in Figure 2(h)(i), we manually detect these results by removing the experiments where
915 the MSE is larger than 1.5 times the average MSE under the same setting, and remove them before
916 computing the average.

917 G.4 More experimental results

918 G.4.1 Regularization weight vs degree-of-freedom

919 As we explained in the previous section, the degree of freedom is the exact information-theoretic
920 measure of the generalization gap. A Larger degree-of-freedom implies more overfitting.

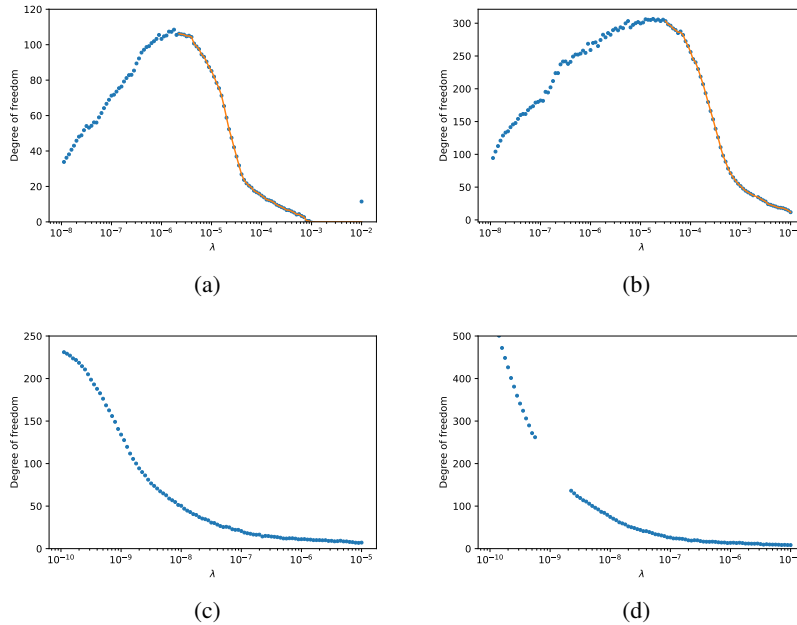


Figure 3: The relationship between degree of freedom and the scaling factor of the regularizer λ . The solid line shows the result after denoising. (a)(b) in a NN. (c)(d) In trend filtering. (a)(c): the piecewise cubic function. (b)(d) the doppler function.

921 In figure Figure 3, we show the relationship between the estimated degree of freedom and the scaling
922 factor of the regularizer λ in a parallel neural network and in trend filtering. As is shown in the

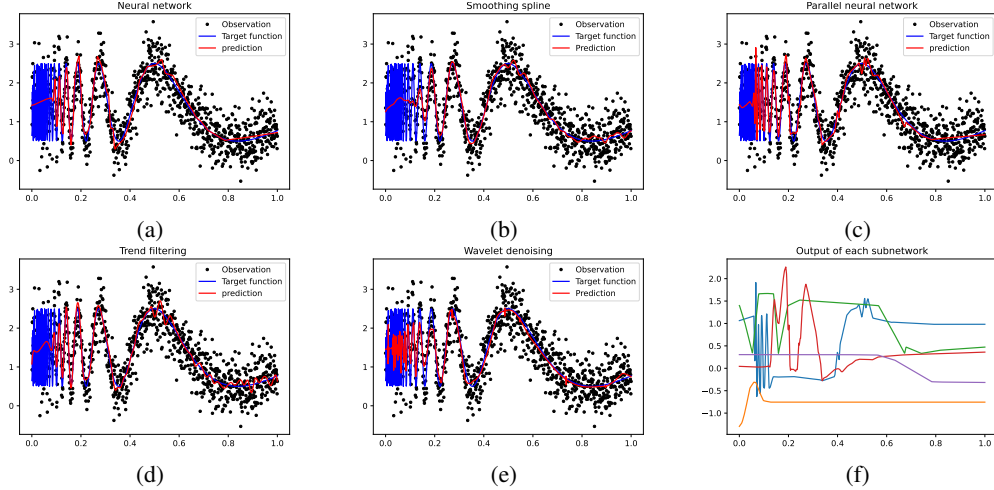


Figure 4: More experiments results of Doppler function.

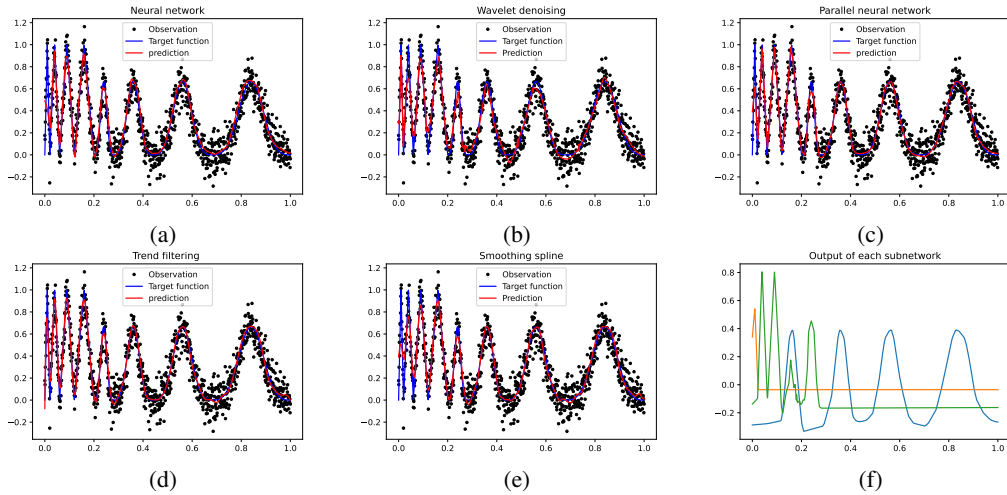


Figure 5: More experiments results of the “vary” function.

923 figure, generally speaking as λ decreases towards 0, the degree of freedom should increase too.
 924 However, for parallel neural networks, if λ is very close to 0, the estimated degree of freedom will
 925 not increase although the degree of freedom is much smaller than the number of parameters —
 926 actually even smaller than the number of subnetworks. Instead, it actually decreases a little. This
 927 effect has not been observed in other nonparametric regression methods, e.g. trend filtering, which
 928 overfits every noisy datapoint perfectly when $\lambda \rightarrow 0$. But for the neural networks, even if we do
 929 not regularize at all, the amount of overfitting is still relatively mild 30/256 vs 80/1000. In our
 930 experiments using neural networks, when λ is small, we denoise the estimated degree of freedom
 931 using isotonic regression.

932 We do not know the exact reason of this curious observation. Our hypothesis is that it might be
 933 related to issues with optimization, i.e., the optimizer ends up at a local minimum that generalizes
 934 better than a global minimum; or it could be connected to the “double descent” behavior of DNN
 935 [26] under over-parameterization.

936 G.4.2 Detailed numerical results

937 In order to allow the readers to view our result in detail, we plot the numerical experiment results of
 938 each method separately in Figure 4 and Figure 5.

939 **G.4.3 Practical equivalence between the weight-decayed two-layer NN and L1-Trend**
 940 **Filtering**

941 In this section we investigate the equivalence of two-layer NN and the locally adaptive regression
 942 splines from Section B. In the special case when $m = 1$ the special regularization reduces to weight
 943 decay and the non-standard truncated power activation becomes ReLU. We compare L1 trend fil-
 944 tering [22] (shown to be equivalent to locally adaptive regression splines by Tibshirani [40]) and
 945 an overparameterized version of the neural network for all regularization parameter $\lambda > 0$, i.e.,
 946 a regularization path. The results are shown in Figure 6. It is clear that as the weight decay in-
 947 creases, it induces sparsity in the number of knots it selects similarly to L1-Trend Filtering, and the
 948 regularization path matches up nearly perfectly even though NNs are also learning knots locations.

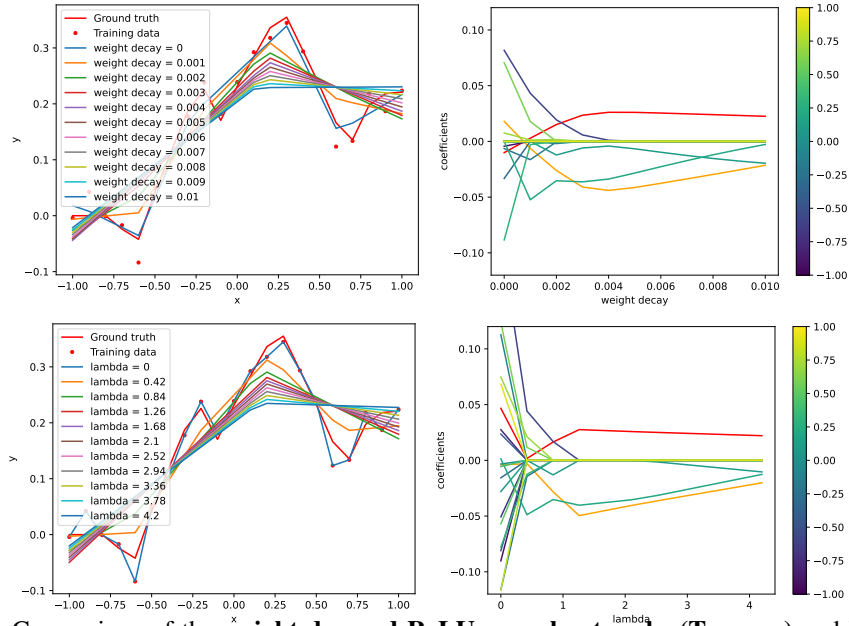


Figure 6: Comparison of the **weight decayed ReLU neural networks (Top row)** and **L1 Trend Filtering (Bottom row)** with different regularization parameters. The left column shows the fitted functions and the right column shows the *regularization path* (in the flavor of [17]) of the coefficients of the truncated power basis at individual data points (the free-knots learned by NN are snapped to the nearest input x to be comparable).