

Deep Network Approximation for Smooth Functions

Jianfeng Lu ^{*} Zuowei Shen [†] Haizhao Yang [‡] Shijun Zhang [§]

Abstract

This paper establishes the (nearly) optimal approximation error characterization of deep rectified linear unit (ReLU) networks for smooth functions in terms of both width and depth simultaneously. To that end, we first prove that multivariate polynomials can be approximated by deep ReLU networks of width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ with an approximation error $\mathcal{O}(N^{-L})$. Through local Taylor expansions and their deep ReLU network approximations, we show that deep ReLU networks of width $\mathcal{O}(N \ln N)$ and depth $\mathcal{O}(L \ln L)$ can approximate $f \in C^s([0, 1]^d)$ with a nearly optimal approximation error $\mathcal{O}(\|f\|_{C^s([0,1]^d)} N^{-2s/d} L^{-2s/d})$. Our estimate is non-asymptotic in the sense that it is valid for arbitrary width and depth specified by $N \in \mathbb{N}^+$ and $L \in \mathbb{N}^+$, respectively.

Key words. Deep ReLU Network, Smooth Function, Polynomial Approximation, Function Composition, Curse of Dimensionality.

1 Introduction

Deep neural networks have made significant impacts in many fields of computer science and engineering, especially for large-scale and high-dimensional learning problems. Well-designed neural network architectures, efficient training algorithms, and high-performance computing technologies have made neural-network-based methods very successful in real applications. Especially in supervised learning; e.g., image classification and objective detection, the great advantages of neural-network-based methods over traditional learning methods have been demonstrated. Understanding the approximation capacity of deep neural networks has become a key question for revealing the power of deep learning. A large number of experiments in real applications have shown the large capacity of deep network approximation from many empirical points of view, motivating much effort in establishing the theoretical foundation of deep network approximation. One of the fundamental problems is the characterization of the optimal approximation error of deep neural networks of arbitrary depth and width.

^{*}Department of Mathematics, Department of Physics, and Department of Chemistry, Duke University (jianfeng@math.duke.edu).

[†]Department of Mathematics, National University of Singapore (matzuows@nus.edu.sg).

[‡]Department of Mathematics, Purdue University (haizhao@purdue.edu).

[§]Department of Mathematics, National University of Singapore (zhangshijun@u.nus.edu).

30 1.1 Main result

31 Previously, the quantitative characterization of the approximation power of deep
 32 feed-forward neural networks (FNNs) with rectified linear unit (ReLU) activation func-
 33 tions was provided in [41]. For ReLU FNNs with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$, the deep
 34 network approximation of $f \in C([0, 1]^d)$ admits an approximation error $\mathcal{O}(\omega_f(N^{-2/d}L^{-2/d}))$
 35 in the L^p -norm for any $p \in [1, \infty]$, where $\omega_f(\cdot)$ is the modulus of continuity of f . In par-
 36 ticular, for the class of Hölder continuous functions, the approximation error is nearly
 37 optimal.^① The next question is whether the smoothness of functions can improve the
 38 approximation error. In this paper, we investigate the deep network approximation of
 39 smaller function space, such as the smooth function space $C^s([0, 1]^d)$.

40 In Theorem 1.1 below, we prove by construction that ReLU FNNs with width
 41 $\mathcal{O}(N \ln N)$ and depth $\mathcal{O}(L \ln L)$ can approximate $f \in C^s([0, 1]^d)$ with a nearly optimal
 42 approximation error $\mathcal{O}(\|f\|_{C^s([0,1]^d)} N^{-2s/d} L^{-2s/d})$, where the norm $\|\cdot\|_{C^s([0,1]^d)}$ is defined
 43 as

$$44 \quad \|f\|_{C^s([0,1]^d)} := \max \{ \|\partial^\alpha f\|_{L^\infty([0,1]^d)} : \|\alpha\|_1 \leq s, \alpha \in \mathbb{N}^d \} \text{ for any } f \in C^s([0, 1]^d).$$

45 **Theorem 1.1.** *Given a smooth function $f \in C^s([0, 1]^d)$ with $s \in \mathbb{N}^+$, for any $N, L \in \mathbb{N}^+$,
 46 there exists a function ϕ implemented by a ReLU FNN with width $C_1(N + 2) \log_2(8N)$
 47 and depth $C_2(L + 2) \log_2(4L) + 2d$ such that*

$$48 \quad \|\phi - f\|_{L^\infty([0,1]^d)} \leq C_3 \|f\|_{C^s([0,1]^d)} N^{-2s/d} L^{-2s/d},$$

49 where $C_1 = 17s^{d+1}3^d$, $C_2 = 18s^2$, and $C_3 = 85(s + 1)^d 8^s$.

50 As we can see from Theorem 1.1, the smoothness improves the approximation error
 51 in N and L ; e.g., $s \geq d$ implies $N^{-2s/d} L^{-2s/d} \leq N^{-2} L^{-2}$. However, we would like to remark
 52 that the improved approximation error is at the price of a prefactor much larger than
 53 d^d if $s \geq d$. The proof of Theorem 1.1 will be presented in Section 2.2 and its tightness
 54 will be discussed in Section 2.3. In fact, the logarithmic terms in width and depth in
 55 Theorem 1.1 can be further reduced if the approximation error is weakened. Given any
 56 $\tilde{N}, \tilde{L} \in \mathbb{N}^+$ with

$$57 \quad \tilde{N} \geq C_1(1 + 2) \log_2(8) = 17s^{d+1}3^{d+2}d \quad \text{and} \quad \tilde{L} \geq C_2(1 + 2) \log_2(4) + 2d = 108s^2 + 2d,$$

58 there exist $N, L \in \mathbb{N}^+$ such that

$$59 \quad C_1(N + 2) \log_2(8N) \leq \tilde{N} < C_1((N + 1) + 2) \log_2(8(N + 1))$$

60 and

$$61 \quad C_2(L + 2) \log_2(4L) + 2d \leq \tilde{L} < C_2((L + 1) + 2) \log_2(4(L + 1)) + 2d.$$

62 It follows that

$$63 \quad N \geq \frac{N + 3}{4} > \frac{\tilde{N}}{4C_1 \log_2(8N + 8)} \geq \frac{\tilde{N}}{4C_1 \log_2(8\tilde{N} + 8)} = \frac{\tilde{N}}{68s^{d+1}3^d d \log_2(8\tilde{N} + 8)}$$

64 and

$$65 \quad L \geq \frac{L + 3}{4} > \frac{\tilde{L} - 2d}{4C_2 \log_2(4L + 4)} \geq \frac{\tilde{L} - 2d}{4C_2 \log_2(4\tilde{L} + 4)} = \frac{\tilde{L} - 2d}{72s^2 \log_2(4\tilde{L} + 4)}.$$

66 Thus, we have an immediate corollary.

^①“nearly optimal” up to a logarithmic factor.

67 **Corollary 1.2.** Given a function $f \in C^s([0, 1]^d)$ with $s \in \mathbb{N}^+$, for any $\tilde{N}, \tilde{L} \in \mathbb{N}^+$, there
 68 exists a function ϕ implemented by a ReLU FNN with width \tilde{N} and depth \tilde{L} such that

$$69 \quad \|\phi - f\|_{L^\infty([0,1]^d)} \leq \tilde{C}_1 \|f\|_{C^s([0,1]^d)} \left(\frac{\tilde{N}}{\tilde{C}_2 \log_2(8\tilde{N}+8)} \right)^{-2s/d} \left(\frac{\tilde{L}-2d}{\tilde{C}_3 \log_2(4\tilde{L}+4)} \right)^{-2s/d}$$

70 for any $\tilde{N} \geq 17s^{d+1}3^{d+2}d$ and $\tilde{L} \geq 108s^2 + 2d$, where $\tilde{C}_1 = 85(s+1)^d 8^s$, $\tilde{C}_2 = 68s^{d+1}3^d$, and
 71 $\tilde{C}_3 = 72s^2$.

72 Theorem 1.1 and Corollary 1.2 characterize the approximation error in terms of
 73 total number of neurons (with an arbitrary distribution in width and depth) and the
 74 smoothness of the target function to be approximated. The only result in this direction
 75 we are aware of in the literature is Theorem 4.1 of [46]. It shows that ReLU FNNs with
 76 width $2d + 10$ and depth L achieve a nearly optimal error $\mathcal{O}((\frac{L}{\ln L})^{-2s/d})$ for sufficiently
 77 large L when approximating functions in the unit ball of $C^s([0, 1]^d)$. This result is
 78 essentially a special case of Corollary 1.2 by setting $\tilde{N} = \mathcal{O}(1)$ and \tilde{L} sufficiently large.

79 1.2 Contributions and related work

80 Our key contributions can be summarized as follows.

81 (i) **Upper bound:** We provide a **quantitative** and **non-asymptotic** approximation
 82 error $\mathcal{O}(\|f\|_{C^s([0,1]^d)} N^{-2s/d} L^{-2s/d})$ when the ReLU FNN has width $\mathcal{O}(N \ln N)$ and
 83 depth $\mathcal{O}(L \ln L)$ for functions in $C^s([0, 1]^d)$ in Theorem 1.1. In real applications,
 84 the first question is to decide the network width and depth since they are two
 85 required hyper-parameters. The approximation error as a function of width and
 86 depth in this paper can directly answer this question, while the approximation
 87 results in terms of the total number of parameters in the literature cannot, because
 88 there are many architectures sharing the same number of parameters. Actually, an
 89 immediate corollary of our theorem as we shall discuss can also describe our theory
 90 in terms of the total number of parameters. Furthermore, our results contain
 91 approximation error estimates for both wide networks with fixed finite depth and
 92 deep networks with fixed finite width.

93 (ii) **Lower bound:** Through the Vapnik-Chervonenkis (VC) dimension upper bound
 94 of ReLU FNNs in [22], we prove a lower bound

$$95 \quad C(N^2 L^2 (\ln N)^3 (\ln L)^3)^{-s/d} \quad \text{for some positive constant } C$$

96 for the approximation error of the functions in the unit ball of $C^s([0, 1]^d)$ approx-
 97 imated by ReLU FNNs with width $\mathcal{O}(N \ln N)$ and depth $\mathcal{O}(L \ln L)$ in Section 2.3.
 98 Thus, the approximation error $\mathcal{O}(N^{-2s/d} L^{-2s/d})$ in Theorem 1.1 is nearly optimal
 99 for the unit ball of $C^s([0, 1]^d)$.

100 (iii) **Approximation of polynomials:** It is proved by construction in Proposition 4.1
 101 that ReLU FNNs with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ can approximate polynomials
 102 on $[0, 1]^d$ with an approximation error $\mathcal{O}(N^{-L})$. This is a non-trivial extension of
 103 the result $\mathcal{O}(2^{-L})$ for polynomial approximation by fixed-width ReLU FNNs with
 104 depth L in [44].

105 (iv) **Uniform approximation:** The approximation error in this paper is measured in
 106 the $L^\infty([0, 1]^d)$ -norm as a result of Theorem 2.1. To achieve this, given a ReLU
 107 FNN approximating the target function f uniformly well on $[0, 1]^d$ except for a
 108 small region, we develop a technique to construct a new ReLU FNN with a similar
 109 size to approximate f **uniformly** well on $[0, 1]^d$ in Theorem 2.1. This technique
 110 can be applied to improve approximation errors from the L^p -norm to the L^∞ -norm
 111 for other function spaces in general, e.g., the continuous function space in [41],
 112 which is of independent interest.

113 In particular, if we denote the best approximation error of functions in $C_u^s([0, 1]^d)$
 114 approximated by ReLU FNNs with width \tilde{N} and depth \tilde{L} as

$$115 \quad \varepsilon_{s,d}(\tilde{N}, \tilde{L}) := \sup_{f \in C_u^s([0,1]^d)} \left(\inf_{\phi \in \mathcal{NN}(\text{width} \leq \tilde{N}; \text{depth} \leq \tilde{L})} \|\phi - f\|_{L^\infty([0,1]^d)} \right) \quad \text{for any } \tilde{N}, \tilde{L} \in \mathbb{N}^+,$$

116 where $C_u^s([0, 1]^d)$ denotes the unit ball of $C^s([0, 1]^d)$ defined by

$$117 \quad C_u^s([0, 1]^d) := \{f \in C^s([0, 1]^d) : \|\partial^\alpha f\|_{L^\infty([0,1]^d)} \leq 1, \text{ for all } \alpha \in \mathbb{N}^d \text{ with } \|\alpha\|_1 \leq s\}.$$

118 By combining the upper and lower bounds stated above, we have

$$119 \quad \underbrace{C_1(s, d) \cdot \left(\tilde{N}^2 \tilde{L}^2 \ln(\tilde{N} \tilde{L}) \right)^{-s/d}}_{\text{proved in Section 2.3}} \leq \varepsilon_{s,d}(\tilde{N}, \tilde{L}) \leq \underbrace{C_2(s, d) \cdot \left(\frac{\tilde{N}^2 \tilde{L}^2}{(\ln \tilde{N} \ln \tilde{L})^2} \right)^{-s/d}}_{\text{shown in Corollary 1.2}},$$

120 where $C_1(s, d)$ and $C_2(s, d)$ are two positive constants in s and d , and $C_2(s, d)$ can be
 121 **explicitly** represented by s and d .

122 The expressiveness of deep neural networks has been studied extensively from many
 123 perspectives, e.g., in terms of combinatorics [34], topology [8], VC-dimension [7, 22, 39],
 124 fat-shattering dimension [2, 27], information theory [37], and classical approximation
 125 theory [4, 5, 9, 12, 14, 15, 20, 21, 24, 29, 32, 35, 42–45, 47]. In the early works of approximation
 126 theory for neural networks, the universal approximation theorem [15, 23, 24] without
 127 approximation errors showed that, given any $\varepsilon > 0$, there exists a sufficiently large neural
 128 network approximating a target function in a certain function space within an error ε .
 129 For one-hidden-layer neural networks and functions with integral representations, Barron
 130 [5, 6] showed an asymptotic approximation error $\mathcal{O}(\frac{1}{\sqrt{N}})$ in the L^2 -norm, leveraging
 131 an idea that is similar to Monte Carlo sampling for high-dimensional integrals. For
 132 very deep ReLU neural networks with width fixed as $\mathcal{O}(d)$ and depth $\mathcal{O}(L)$, Yarotsky
 133 [45, 46] showed that the nearly optimal approximation errors for Lipschitz continuous
 134 functions and functions in the unit ball of $C^s([0, 1]^d)$ are $\mathcal{O}(L^{-2/d})$ and $\mathcal{O}((L/\ln L)^{-2s/d})$,
 135 respectively. Note that the results are asymptotic in the sense that L is required to be
 136 sufficiently large and the prefactors of these rates are unknown. To obtain a generic result
 137 that characterizes the approximation error for arbitrary width and depth with known
 138 prefactors to guide applications, the authors of [41] demonstrated that the nearly optimal
 139 approximation error for ReLU FNNs with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ to approximate
 140 Lipschitz continuous functions on $[0, 1]^d$ is $\mathcal{O}(N^{-2/d} L^{-2/d})$. Such a nearly optimal error
 141 is further improved to an optimal one, $\mathcal{O}((N^2 L^2 \ln N)^{-1/d})$, in a more recent paper [42].
 142 In this paper, we extend this generic framework to $C^s([0, 1]^d)$ with a nearly optimal
 143 approximation error $\mathcal{O}(\|f\|_{C^s([0,1]^d)} N^{-2s/d} L^{-2s/d})$.

144 Most related works are summarized in Table 1 for the comparison of our contribu-
 145 tions in this paper and the results in the literature.

Table 1: A summary of existing approximation errors of ReLU FNNs for $\text{Lip}([0, 1]^d)$ (the Lipschitz continuous function space) and $C_u^s([0, 1]^d)$ (the unit ball of $C^s([0, 1]^d)$).

| paper | function class | width | depth | approximation error | $L^p([0, 1]^d)$ -norm | tightness | valid for |
|------------|------------------------|------------------------|------------------------|--|-----------------------|-----------------------------|-----------------------------|
| [44] | polynomial | $\mathcal{O}(1)$ | $\mathcal{O}(L)$ | $\mathcal{O}(2^{-L})$ | $p = \infty$ | | any $L \in \mathbb{N}^+$ |
| this paper | polynomial | $\mathcal{O}(N)$ | $\mathcal{O}(L)$ | $\mathcal{O}(N^{-L})$ | $p = \infty$ | | any $N, L \in \mathbb{N}^+$ |
| [40] | $\text{Lip}([0, 1]^d)$ | $\mathcal{O}(N)$ | 3 | $\mathcal{O}(N^{-2/d})$ | $p \in [1, \infty)$ | nearly tight in N | any $N \in \mathbb{N}^+$ |
| [45] | $\text{Lip}([0, 1]^d)$ | $2d + 10$ | $\mathcal{O}(L)$ | $\mathcal{O}(L^{-2/d})$ | $p = \infty$ | nearly tight in L | large $L \in \mathbb{N}^+$ |
| [41] | $\text{Lip}([0, 1]^d)$ | $\mathcal{O}(N)$ | $\mathcal{O}(L)$ | $\mathcal{O}(N^{-2/d} L^{-2/d})$ | $p \in [1, \infty]$ | nearly tight in N and L | any $N, L \in \mathbb{N}^+$ |
| [42] | $\text{Lip}([0, 1]^d)$ | $\mathcal{O}(N)$ | $\mathcal{O}(L)$ | $\mathcal{O}((N^2 L^2 \ln N)^{-1/d})$ | $p \in [1, \infty]$ | tight in N and L | any $N, L \in \mathbb{N}^+$ |
| [46] | $C_u^s([0, 1]^d)$ | $2d + 10$ | $\mathcal{O}(L)$ | $\mathcal{O}((L/\ln L)^{-2s/d})$ | $p = \infty$ | neatly tight in L | large $L \in \mathbb{N}^+$ |
| this paper | $C_u^s([0, 1]^d)$ | $\mathcal{O}(N \ln N)$ | $\mathcal{O}(L \ln L)$ | $\mathcal{O}(N^{-2s/d} L^{-2s/d})$ | $p = \infty$ | nearly tight in N and L | any $N, L \in \mathbb{N}^+$ |
| this paper | $C_u^s([0, 1]^d)$ | $\mathcal{O}(N)$ | $\mathcal{O}(L)$ | $\mathcal{O}((N/\ln N)^{-2s/d} (L/\ln L)^{-2s/d})$ | $p = \infty$ | nearly tight in N and L | any $N, L \in \mathbb{N}^+$ |

1.3 Discussion

We will discuss the comparison of our theory with existing works and the application scope in machine learning.

Approximation errors in $\mathcal{O}(N)$ and $\mathcal{O}(L)$ versus $\mathcal{O}(W)$

It is fundamental and indispensable to characterize deep network approximation in terms of width $\mathcal{O}(N)$ ^② and depth $\mathcal{O}(L)$ simultaneously in realistic applications, while the approximation in terms of the number of nonzero parameters W is probably only of interest in theory. First, networks used in practice are specified via width and depth and, therefore, Theorem 1.1 can provide an error bound for such networks. However, existing results in W cannot serve this purpose because they may be only valid for networks with other widths and depths. Theories in terms of W essentially have a single variable to control the network size in three types of structures: 1) a fixed width N and a varying depth L ; 2) a fixed depth L and a varying width N ; 3) both the width and depth are controlled by the target error ε (e.g., N is a polynomial of $\frac{1}{\varepsilon^d}$ and L is a polynomial of $\ln(\frac{1}{\varepsilon})$). Therefore, given a network with arbitrary width N and depth L , there might not be a known theory in terms of W to quantify the performance of this structure. Second, the error characterization in terms of N and L is more useful than that in terms of W , because most existing optimization and generalization analyses are based on N and L [1, 3, 10, 13, 17, 18, 25, 26], to the best of our knowledge. Approximation results in terms of N and L are more consistent with optimization and generalization analysis tools to obtain a full error analysis.

Most existing approximation theories for deep neural networks so far focus on the approximation error in the number of parameters W [4, 5, 9, 11, 12, 14, 15, 19–21, 24, 29–33, 35–38, 43–47]. Controlling two variables N and L in our theory is more challenging than controlling one variable W in the literature. The characterization of deep network approximation in terms of N and L can imply an approximation error in terms of W , while this may not be true the other way around, e.g., our theorems cannot be derived from results in [46]. Let us discuss the first type of structure mentioned in the previous paragraph, which includes the best-known result for a nearly optimal approximation error, $\mathcal{O}((W/\ln W)^{-2s/d})$, for functions in the unit ball of $C^s([0, 1]^d)$ using ReLU FNNs with W parameters [46]. As an example to show how Theorem 1.1 in terms of N and

^②For simplicity, we omit $\mathcal{O}(\cdot)$ in the following discussion.

177 L can be applied to show a similar result in terms of W . The main idea is to specify
 178 the value of N and L in Theorem 1.1 to show the desired corollary. For example, if we
 179 let $N = \mathcal{O}(1)$ in Theorem 1.1, then we have the following corollary, which is essentially
 180 equivalent to Theorem 4.1 of [46].

181 **Corollary 1.3.** *Given any function f in the unit ball of $C^s([0,1]^d)$ with $s \in \mathbb{N}^+$, there*
 182 *exists a function ϕ implemented by a ReLU FNN with W parameters such that*

$$183 \quad \|\phi - f\|_{L^\infty([0,1]^d)} \leq \mathcal{O}\left(\left(\frac{W}{\ln W}\right)^{-2s/d}\right) \quad \text{for large } W \in \mathbb{N}^+.$$

184 As we can see in this example, it is simple to derive Corollary 1.3 above and The-
 185 orem 4.1 of [46] using Theorem 1.1 in this paper. However, Theorem 1.1 cannot be
 186 derived from any existing result that characterizes approximation errors in terms of the
 187 number of parameters. Therefore, Theorem 1.1 goes beyond existing results on the
 188 approximation of deep neural networks.

189 Note that the logarithmic term in the approximation error is not significant in the
 190 case of $s > 1$ since it can be cancelled out in the sense that $\left(\frac{W}{\ln W}\right)^{-2s/d} \lesssim W^{-2\tilde{s}/d}$ for
 191 any $\tilde{s} \in (1, s)$. We remark that Theorem 3.3 of [46] provides a better approximation
 192 error by a logarithmic term: ReLU FNNs with W nonzero parameters can approximate
 193 a function f in the unit ball of $C^s([0,1]^d)$ within an error $\mathcal{O}(W^{-2s/d})$. However, the
 194 network architecture therein is relatively complex and s -dependent as stated by the
 195 authors of [46]. In fact, it contains many s -dependent blocks (sub-networks), making
 196 it difficult to implement if s is not known in applications. In contrast, our network
 197 architecture in Corollary 1.2 is simple and can be pre-specified once the width \tilde{N} and
 198 depth \tilde{L} therein are given.

199 Continuity of the weight selection

200 We would like to discuss the continuity of the weight selection as a map $\Sigma : F_{s,d} \rightarrow$
 201 \mathbb{R}^W , where $F_{s,d}$ denotes the unit ball of the d -dimensional Sobolev space with smooth-
 202 ness s . For a fixed network architecture with a fixed number of parameters W , let
 203 $g : \mathbb{R}^W \rightarrow C([0,1]^d)$ be the map of realizing a ReLU FNN from a given set of param-
 204 eters in \mathbb{R}^W to a function in $C([0,1]^d)$. Suppose that the map Σ is continuous such
 205 that $\|f - g(\Sigma(f))\|_{L^\infty([0,1]^d)} \leq \varepsilon$ for all $f \in F_{s,d}$. Then $W \geq c\varepsilon^{-d/s}$ with some constant c
 206 depending only on s . This conclusion is given in Theorem 3 of [44], which is a corollary
 207 of Theorem 4.2 of [16] in a more general form. These theorems mean that the weight
 208 selection map Σ corresponding to our constructive proof in Theorem 1.1 in this paper is
 209 not continuous, since our error is better than $\mathcal{O}(W^{-s/d})$. Theorem 4.2 of [16] is essentially
 210 a min-max criterion to evaluate weight selection maps maintaining continuity: the ap-
 211 proximation error obtained by minimizing over all continuous selections Σ and network
 212 realizations g and maximizing over all target functions is bounded below by $\mathcal{O}(W^{-s/d})$.
 213 In the worst case, a continuous weight selection cannot enjoy an approximation error
 214 beating $\mathcal{O}(W^{-s/d})$. However, Theorem 4.2 of [16] does not exclude the possibility that
 215 most functions of interest in practice may still enjoy a continuous weight selection with
 216 the approximation error in Theorem 1.1. It would be interesting in future work to in-
 217 vestigate whether continuous weight selection is possible for many functions commonly
 218 encountered in real applications.

219 **Application scope of our theory in machine learning**

220 In deep learning, given a target function f , the final goal is to train a function
 221 $\phi(\mathbf{x}; \boldsymbol{\theta})$ approximating f well, where $\phi(\mathbf{x}; \boldsymbol{\theta})$ is a function in $\mathbf{x} \in \mathcal{X}$ realized by a network
 222 architecture parameterized with $\boldsymbol{\theta} \in \mathbb{R}^W$. To get the best solution, one needs to identify
 223 the expected risk minimizer

$$224 \quad \boldsymbol{\theta}_{\mathcal{D}} := \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^W} R_{\mathcal{D}}(\boldsymbol{\theta}), \quad \text{where } R_{\mathcal{D}}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x} \sim U(\mathcal{X})} [\ell(\phi(\mathbf{x}; \boldsymbol{\theta}), f(\mathbf{x}))]$$

225 with a loss function usually taken as $\ell(y, y') = \frac{1}{2}|y - y'|^2$ and an unknown data distribution
 226 $U(\mathcal{X})$.

227 In practice, only data samples $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n$ instead of f and $U(\mathcal{X})$ are available.
 228 Thus, the empirical risk minimizer $\boldsymbol{\theta}_{\mathcal{S}}$ is used to model/approximate the expected risk
 229 minimizer $\boldsymbol{\theta}_{\mathcal{D}}$, where

$$230 \quad \boldsymbol{\theta}_{\mathcal{S}} := \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^W} R_{\mathcal{S}}(\boldsymbol{\theta}), \quad \text{where } R_{\mathcal{S}}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \ell(\phi(\mathbf{x}_i, \boldsymbol{\theta}), f(\mathbf{x}_i)). \quad (1.1)$$

231 In real applications, only a numerical solution (denoted as $\boldsymbol{\theta}_{\mathcal{N}}$) is achieved when
 232 a numerical optimization method is applied to solve (1.1). Hence, the actually learned
 233 function generated by the network is $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}})$. Since $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}})$ is the expected inference
 234 error over all possible data samples, it can quantify how good $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}})$ is. Note that

$$235 \quad R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) = \underbrace{[R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}})]}_{\text{GE}} + \underbrace{[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}})]}_{\text{OE}} + \underbrace{[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}})]}_{\leq 0 \text{ by (1.1)}} + \underbrace{[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}}) - R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})]}_{\text{GE}} + \underbrace{R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})}_{\text{AE}}$$

$$236 \quad \leq \underbrace{R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})}_{\text{Approximation error (AE)}} + \underbrace{[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}})]}_{\text{Optimization error (OE)}} + \underbrace{[R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}})] + [R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}}) - R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})]}_{\text{Generalization error (GE)}}. \quad (1.2)$$

$$237 \quad \text{Approximation error (AE)} \quad \text{Optimization error (OE)} \quad \text{Generalization error (GE)}$$

238 Constructive approximation provides an upper bound of $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$ in terms of the
 239 network size. For example, Theorem 1.1 and its corollaries provide an upper bound
 240 $\mathcal{O}(\|f\|_{C^s([0,1]^d)} N^{-2s/d} L^{-2s/d})$ of $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$ for $C^s([0,1]^d)$. The second term of (1.2) is
 241 bounded by the optimization error of the numerical algorithm applied to solve the em-
 242 pirical loss minimization problem in (1.1). The study of the bounds for the third and
 243 fourth terms is referred to as the generalization error analysis of neural networks.

244 One of the key targets in the area of deep learning is to develop algorithms to
 245 reduce $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}})$. Our analysis here provides an upper bound of the approximation error
 246 $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$ for smooth functions, which is crucial to control $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}})$. Instead of deriving
 247 an approximator to attain the error bound, deep learning algorithms aim to identify a
 248 solution $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}})$ reducing the generalization and optimization errors in (1.2). Solutions
 249 minimizing both generalization and optimization errors will lead to a good solution only
 250 if we also have a good upper bound estimate of $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$ as shown in (1.2). Independent
 251 of whether our analysis here leads to a good approximator, which is an interesting topic
 252 to pursue, the theory here does provide a key ingredient in the error analysis of deep
 253 learning algorithms.

254 We would like to emphasize that the introduction of the ReLU activation function
 255 to image classification is one of the key techniques that boost the performance of deep
 256 learning [28] with surprising generalization, which is the main reason that we focus on
 257 ReLU FNNs in this paper.

258 **Organization:** The rest of the present paper is organized as follows. In Section 2,
 259 we prove Theorem 1.1 by combining two theorems (Theorems 2.1 and 2.2) that will be
 260 proved later. We will also discuss the optimality of Theorem 1.1 in Section 2. Next,
 261 Theorem 2.1 will be proved in Section 3 while Theorem 2.2 will be shown in Section 4.
 262 Several propositions supporting Theorem 2.2 will be presented in Section 5. Finally,
 263 Section 6 concludes this paper with a short discussion.

264 2 Approximation of smooth functions

265 In this section, we will prove the quantitative approximation error in Theorem 1.1 by
 266 construction and discuss its tightness. Notation throughout the proof will be summarized
 267 in Section 2.1. The proof of Theorem 1.1 is mainly based on Theorems 2.1 and 2.2, which
 268 will be proved in Sections 3 and 4, respectively. To show the tightness of Theorem 1.1,
 269 we will introduce the VC-dimension in Section 2.3.

270 2.1 Notation

271 Now let us summarize the main notation of this paper as follows.

- 272 • Let \mathbb{R} , \mathbb{Q} , and \mathbb{Z} denote the set of real numbers, rational numbers, and integers,
 273 respectively.
- 274 • Let \mathbb{N} and \mathbb{N}^+ denote the set of natural numbers and positive natural numbers,
 275 respectively. That is, $\mathbb{N}^+ = \{1, 2, 3, \dots\}$ and $\mathbb{N} = \mathbb{N}^+ \cup \{0\}$.
- 276 • Vectors and matrices are denoted in a bold font. Standard vectorization is adopted
 277 in matrix and vector computation. For example, a scalar plus a vector means
 278 adding the scalar to each entry of the vector. Additionally, “[” and “]” are used
 279 to partition matrices (vectors) into blocks, e.g., $\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}$ and $\mathbf{v} = \begin{bmatrix} v_1 \\ \vdots \\ v_d \end{bmatrix} =$
 280 $[v_1, \dots, v_d]^T \in \mathbb{R}^d$.
- 281 • Let $\mathbf{1}_S$ be the characteristic (indicator) function on a set S ; i.e., $\mathbf{1}_S$ is equal to 1
 282 on S and 0 outside S .
- 283 • Let $\mathcal{B}(\mathbf{x}, r) \subseteq \mathbb{R}^d$ be the closed ball with a center $\mathbf{x} \subseteq \mathbb{R}^d$ and a radius $r \geq 0$.
- 284 • Similar to “min” and “max”, let $\text{mid}(x_1, x_2, x_3)$ be the middle value of three inputs
 285 x_1 , x_2 , and x_3 ^③. For example, $\text{mid}(2, 1, 3) = 2$ and $\text{mid}(3, 2, 3) = 3$.
- 286 • The set difference of two sets A and B is denoted by $A \setminus B := \{x : x \in A, x \notin B\}$.
- 287 • For a real number $p \in [1, \infty)$, the p -norm of $\mathbf{x} = [x_1, x_2, \dots, x_d]^T \in \mathbb{R}^d$ is defined by

$$288 \quad \|\mathbf{x}\|_p := \left(|x_1|^p + |x_2|^p + \dots + |x_d|^p\right)^{1/p}.$$

^③“mid” can be defined via $\text{mid}(x_1, x_2, x_3) = x_1 + x_2 + x_3 - \max(x_1, x_2, x_3) - \min(x_1, x_2, x_3)$, which can be implemented by a ReLU FNN.

- 289 • For any $x \in \mathbb{R}$, let $\lfloor x \rfloor := \max\{n : n \leq x, n \in \mathbb{Z}\}$ and $\lceil x \rceil := \min\{n : n \geq x, n \in \mathbb{Z}\}$.
- 290 • Assume $\mathbf{n} \in \mathbb{N}^d$; then $f(\mathbf{n}) = \mathcal{O}(g(\mathbf{n}))$ means that there exists positive C independent of \mathbf{n} , f , and g such that $f(\mathbf{n}) \leq Cg(\mathbf{n})$ when all entries of \mathbf{n} go to $+\infty$.
- 291
- 292 • The modulus of continuity of a continuous function $f \in C([0, 1]^d)$ is defined as

293
$$\omega_f(r) := \sup \{|f(\mathbf{x}) - f(\mathbf{y})| : \|\mathbf{x} - \mathbf{y}\|_2 \leq r, \mathbf{x}, \mathbf{y} \in [0, 1]^d\} \quad \text{for any } r \geq 0.$$

- 294 • A d -dimensional multi-index is a d -tuple $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_d]^T \in \mathbb{N}^d$. Several related notations are listed below.
- 295

- 296 – $\|\boldsymbol{\alpha}\|_1 = |\alpha_1| + |\alpha_2| + \dots + |\alpha_d|$;
- 297 – $\mathbf{x}^\boldsymbol{\alpha} = x_1^{\alpha_1} x_2^{\alpha_2} \dots x_d^{\alpha_d}$, where $\mathbf{x} = [x_1, x_2, \dots, x_d]^T$;
- 298 – $\boldsymbol{\alpha}! = \alpha_1! \alpha_2! \dots \alpha_d!$;
- 299 – $\partial^\boldsymbol{\alpha} = \frac{\partial^{\alpha_1}}{\partial x_1^{\alpha_1}} \frac{\partial^{\alpha_2}}{\partial x_2^{\alpha_2}} \dots \frac{\partial^{\alpha_d}}{\partial x_d^{\alpha_d}}$.

- 300 • For any closed cube $Q \subseteq \mathbb{R}^d$ and a real number $r > 0$, let rQ denote the closed cube which shares the same center of Q and whose sidelength is the product of r and the sidelength of Q .
- 301
- 302

- 303 • Given any $K \in \mathbb{N}^+$ and $\delta \in (0, \frac{1}{K})$, define a trifling region $\Omega([0, 1]^d, K, \delta)$ of $[0, 1]^d$ as
- 304

305
$$\Omega([0, 1]^d, K, \delta) := \bigcup_{i=1}^d \left\{ \mathbf{x} = [x_1, x_2, \dots, x_d]^T \in [0, 1]^d : x_i \in \bigcup_{k=1}^{K-1} \left(\frac{k}{K} - \delta, \frac{k}{K} \right) \right\}. \quad (2.1)$$

306 In particular, $\Omega([0, 1]^d, K, \delta) = \emptyset$ if $K = 1$. See Figure 1 for two examples of the trifling region.

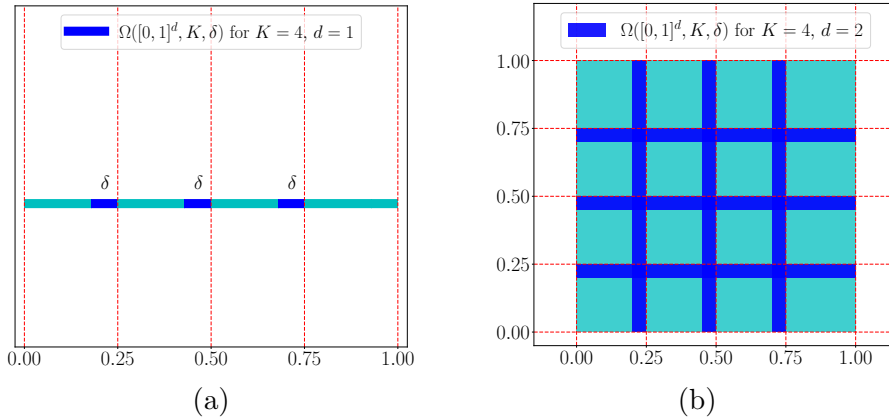


Figure 1: Two examples of the trifling region. (a) $K = 4, d = 1$. (b) $K = 4, d = 2$.

307

- 308 • Given $E \subseteq \mathbb{R}^d$, let $C^s(E)$ denote the set containing all functions, all k -th order partial derivatives of which exist and are continuous on E for any $k \in \mathbb{N}$ with $0 \leq k \leq s$. In particular, $C^0(E)$, also denoted by $C(E)$, is the set of continuous
- 309
- 310

311

functions on E . For the case $s = \infty$, $C^\infty(E) = \bigcap_{s=0}^\infty C^s(E)$. The C^s -norm is defined by

312

313

$$\|f\|_{C^s(E)} := \max \{ \|\partial^\alpha f\|_{L^\infty(E)} : \alpha \in \mathbb{N}^d \text{ with } \|\alpha\|_1 \leq s \}.$$

314

Generally, E is assigned as $[0, 1]^d$ in this paper. In particular, the closed unit ball of $C^s([0, 1]^d)$ is denoted by

315

316

$$C_u^s([0, 1]^d) := \{f \in C^s([0, 1]^d) : \|f\|_{C^s([0, 1]^d)} \leq 1\}.$$

317

- We use “ \mathcal{NN} ” to mean “functions implemented by ReLU FNNs” for short and use Python-type notation to specify a class of functions implemented by ReLU FNNs with several conditions. To be precise, we use $\mathcal{NN}(c_1; c_2; \dots; c_m)$ to denote the function set containing all functions implemented by ReLU FNN architectures satisfying m conditions given by $\{c_i\}_{1 \leq i \leq m}$, each of which may specify the number of inputs (#input), the number of outputs (#output), the total number of nodes in all hidden layers (#neuron), the number of hidden layers (depth), the number of total parameters (#parameter), and the width in each hidden layer (widthvec), the maximum width of all hidden layers (width), etc. For example, if $\phi \in \mathcal{NN}(\text{\#input} = 2; \text{widthvec} = [100, 100]; \text{\#output} = 1)$, then ϕ is a function satisfying the following conditions.

328

– ϕ maps from \mathbb{R}^2 to \mathbb{R} .

329

– ϕ is implemented by a ReLU FNN with two hidden layers and the number of nodes in each hidden layer being 100.

330

331

- Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ denote the rectified linear unit (ReLU), i.e. $\sigma(x) = \max\{0, x\}$. With

332

the abuse of notation, we define $\sigma : \mathbb{R}^d \rightarrow \mathbb{R}^d$ as $\sigma(\mathbf{x}) = \begin{bmatrix} \max\{0, x_1\} \\ \vdots \\ \max\{0, x_d\} \end{bmatrix}$ for any

333

$$\mathbf{x} = [x_1, \dots, x_d]^T \in \mathbb{R}^d.$$

334

- For a function $\phi \in \mathcal{NN}(\text{\#input} = d; \text{widthvec} = [N_1, N_2, \dots, N_L]; \text{\#output} = 1)$, if we set $N_0 = d$ and $N_{L+1} = 1$, then the architecture of the network implementing ϕ can be briefly described as follows:

335

336

$$\mathbf{x} = \tilde{\mathbf{h}}_0 \xrightarrow{\mathbf{W}_0, \mathbf{b}_0} \mathbf{h}_1 \xrightarrow{\sigma} \tilde{\mathbf{h}}_1 \quad \dots \quad \xrightarrow{\mathbf{W}_{L-1}, \mathbf{b}_{L-1}} \mathbf{h}_L \xrightarrow{\sigma} \tilde{\mathbf{h}}_L \xrightarrow{\mathbf{W}_L, \mathbf{b}_L} \mathbf{h}_{L+1} = \phi(\mathbf{x}),$$

337

where $\mathbf{W}_i \in \mathbb{R}^{N_{i+1} \times N_i}$ and $\mathbf{b}_i \in \mathbb{R}^{N_{i+1}}$ are the weight matrix and the bias vector in the i -th affine linear transform \mathcal{L}_i in ϕ , respectively, i.e.,

338

339

$$\mathbf{h}_{i+1} = \mathbf{W}_i \cdot \tilde{\mathbf{h}}_i + \mathbf{b}_i =: \mathcal{L}_i(\tilde{\mathbf{h}}_i) \quad \text{for } i = 0, 1, \dots, L$$

340

and

341

342

$$\tilde{\mathbf{h}}_i = \sigma(\mathbf{h}_i) \quad \text{for } i = 1, 2, \dots, L.$$

343

In particular, ϕ can be represented in a form of function compositions as follows

344

$$\phi = \mathcal{L}_L \circ \sigma \circ \mathcal{L}_{L-1} \circ \sigma \circ \dots \circ \sigma \circ \mathcal{L}_1 \circ \sigma \circ \mathcal{L}_0,$$

345

which has been illustrated in Figure 2.

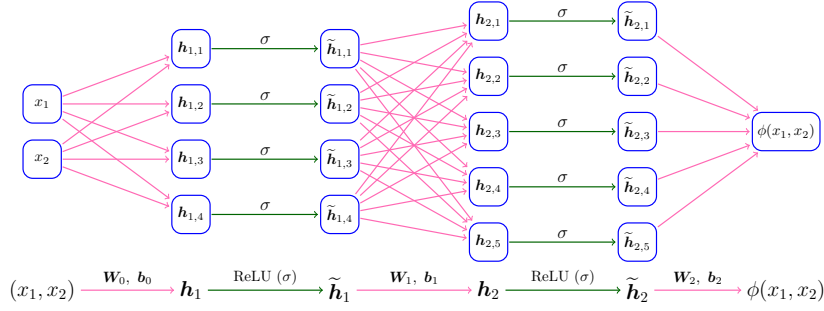


Figure 2: An example of a ReLU FNN with width 5 and depth 2.

- 346 • The expression “a network (architecture) with (of) width N and depth L ” means
- 347 – The maximum width of this network (architecture) for all **hidden** layers is
- 348 no more than N .
- 349 – The number of **hidden** layers of this network (architecture) is no more than
- 350 L .
- 351 • For any $\theta \in [0, 1)$, suppose its binary representation is $\theta = \sum_{\ell=1}^{\infty} \theta_{\ell} 2^{-\ell}$ with $\theta_{\ell} \in$
- 352 $\{0, 1\}$. We introduce a special notation $\text{bin}0.\theta_1\theta_2\cdots\theta_L$ to denote the L -term binary
- 353 representation of θ , i.e., $\text{bin}0.\theta_1\theta_2\cdots\theta_L := \sum_{\ell=1}^L \theta_{\ell} 2^{-\ell} \approx \theta$.

354 2.2 Proof of Theorem 1.1

355 The introduction of the trifling region $\Omega([0, 1]^d, K, \delta)$ is due to the fact that ReLU

356 FNNs cannot approximate a step function uniformly well (as the ReLU activation func-

357 tion is continuous), which is also the reason for the main difficulty in obtaining approxi-

358 mation errors in the $L^{\infty}([0, 1]^d)$ -norm in our previous papers [40, 41]. The trifling region

359 is a key technique to simplify the proofs of theories in [40, 41] as well as the proof of

360 Theorem 1.1.

361 First, we present Theorem 2.1 to show that, as long as good uniform approximation

362 by a ReLU FNN can be obtained outside the trifling region, the uniform approximation

363 error can also be well controlled inside the trifling region when the network size is slightly

364 increased. Second, as a simplified version of Theorem 1.1 ignoring the approximation

365 error in the trifling region $\Omega([0, 1]^d, K, \delta)$, Theorem 2.2 shows the existence of a ReLU

366 FNN approximating a target smooth function uniformly well outside the trifling region.

367 Finally, Theorems 2.1 and 2.2 immediately lead to Theorem 1.1. Theorem 2.1 can

368 be applied to improve the theories in [40, 41] to obtain approximation errors in the

369 $L^{\infty}([0, 1]^d)$ -norm.

370 **Theorem 2.1.** *Given any $\varepsilon > 0$, $N, L, K \in \mathbb{N}^+$, and $\delta \in (0, \frac{1}{3K}]$, assume $f \in C([0, 1]^d)$*

371 *and $\tilde{\phi}$ is a function implemented by a ReLU FNN with width N and depth L . If*

$$372 \quad |\tilde{\phi}(\mathbf{x}) - f(\mathbf{x})| \leq \varepsilon \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta),$$

373 *then there exists a new function ϕ implemented by a ReLU FNN with width $3^d(N + 4)$*

374 *and depth $L + 2d$ such that*

$$375 \quad |\phi(\mathbf{x}) - f(\mathbf{x})| \leq \varepsilon + d \cdot \omega_f(\delta) \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

376 **Theorem 2.2.** Assume that $f \in C^s([0, 1]^d)$ satisfies $\|\partial^\alpha f\|_{L^\infty([0, 1]^d)} \leq 1$ for any $\alpha \in \mathbb{N}^d$
377 with $\|\alpha\|_1 \leq s$. For any $N, L \in \mathbb{N}^+$, there exists a function ϕ implemented by a ReLU
378 FNN with width $16s^{d+1}d(N+2)\log_2(8N)$ and depth $18s^2(L+2)\log_2(4L)$ such that

$$379 \quad |\phi(\mathbf{x}) - f(\mathbf{x})| \leq 84(s+1)^d 8^s N^{-2s/d} L^{-2s/d} \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta),$$

380 where $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor$ and δ is an arbitrary number in $(0, \frac{1}{3K}]$.

381 We first prove Theorem 1.1 by assuming Theorems 2.1 and 2.2 are true. The proofs
382 of Theorems 2.1 and 2.2 can be found in Sections 3 and 4, respectively.

383 *Proof of Theorem 1.1.* We may assume $\|f\|_{C^s([0, 1]^d)} > 0$ since $\|f\|_{C^s([0, 1]^d)} = 0$ is a trivial
384 case. Define $\tilde{f} := \frac{f}{\|f\|_{C^s([0, 1]^d)}} \in C_u^s([0, 1]^d)$. Set $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor$ and choose a small
385 $\delta \in (0, \frac{1}{3K}]$ such that

$$386 \quad d \cdot \omega_{\tilde{f}}(\delta) \leq N^{-2s/d} L^{-2s/d}.$$

387 Clearly, $\|\partial^\alpha \tilde{f}\|_{L^\infty([0, 1]^d)} \leq 1$ for any $\alpha \in \mathbb{N}^d$ with $\|\alpha\|_1 \leq s$. By Theorem 2.2, there
388 exists a function $\tilde{\phi}$ implemented by a ReLU FNN with width $16s^{d+1}d(N+2)\log_2(8N)$
389 and depth $18s^2(L+2)\log_2(4L)$ such that

$$390 \quad |\tilde{\phi}(\mathbf{x}) - \tilde{f}(\mathbf{x})| \leq 84(s+1)^d 8^s N^{-2s/d} L^{-2s/d} =: \varepsilon \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta).$$

391 By Theorem 2.1, there exists a new function $\tilde{\phi}$ implemented by a ReLU FNN with width

$$392 \quad 3^d(16s^{d+1}d(N+2)\log_2(8N) + 4) \leq 17s^{d+1}3^d d(N+2)\log_2(8N)$$

393 and depth $18s^2(L+2)\log_2(4L) + 2d$ such that

$$394 \quad \begin{aligned} \|\tilde{\phi} - \tilde{f}\|_{L^\infty([0, 1]^d)} &\leq \varepsilon + d \cdot \omega_{\tilde{f}}(\delta) = 84(s+1)^d 8^s N^{-2s/d} L^{-2s/d} + d \cdot \omega_{\tilde{f}}(\delta) \\ &\leq 85(s+1)^d 8^s N^{-2s/d} L^{-2s/d}. \end{aligned}$$

395 Finally, set $\phi = \|f\|_{C^s([0, 1]^d)} \cdot \tilde{\phi}$; then

$$396 \quad \begin{aligned} \|\phi - f\|_{L^\infty([0, 1]^d)} &= \|f\|_{C^s([0, 1]^d)} \cdot \|\tilde{\phi} - \tilde{f}\|_{L^\infty([0, 1]^d)} \\ &\leq 85(s+1)^d 8^s \|f\|_{C^s([0, 1]^d)} N^{-2s/d} L^{-2s/d}, \end{aligned}$$

397 and ϕ can also be implemented by a ReLU FNN with width $17s^{d+1}3^d d(N+2)\log_2(8N)$
398 and depth $18s^2(L+2)\log_2(4L) + 2d$. So we finish the proof. \square

399 2.3 Optimality of Theorem 1.1

400 In this section, we will show that the approximation error in Theorem 1.1 is nearly
401 tight in terms of VC-dimension. The key is the VC-dimension upper bound of ReLU
402 FNNs in [22] will lead to a contradiction if our approximation is not optimal. This
403 idea was used in [44] to prove its tightness for ReLU FNNs of width $\mathcal{O}(d)$ and depth
404 sufficiently large to approximate smooth functions.

405 Let us first present the definitions of VC-dimension and related concepts. Let H be
 406 a class of functions mapping from a general domain \mathcal{X} to $\{0, 1\}$. We say H shatters the
 407 set $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\} \subseteq \mathcal{X}$ if

$$408 \quad \left| \left\{ \left[h(\mathbf{x}_1), h(\mathbf{x}_2), \dots, h(\mathbf{x}_m) \right]^T \in \{0, 1\}^m : h \in H \right\} \right| = 2^m,$$

409 where $|\cdot|$ means the size of a set. This equation means, given any $\theta_i \in \{0, 1\}$ for $i =$
 410 $1, 2, \dots, m$, there exists $h \in H$ such that $h(\mathbf{x}_i) = \theta_i$ for all i . For a general function set \mathcal{F}
 411 mapping from \mathcal{X} to \mathbb{R} , we say \mathcal{F} shatters $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\} \subseteq \mathcal{X}$ if $\mathcal{T} \circ \mathcal{F}$ does, where

$$412 \quad \mathcal{T}(t) := \begin{cases} 1, & t \geq 0, \\ 0, & t < 0 \end{cases} \quad \text{and} \quad \mathcal{T} \circ \mathcal{F} := \{\mathcal{T} \circ f : f \in \mathcal{F}\}.$$

413 For any $m \in \mathbb{N}^+$, we define the growth function of H as

$$414 \quad \Pi_H(m) := \max_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m \in \mathcal{X}} \left| \left\{ \left[h(\mathbf{x}_1), h(\mathbf{x}_2), \dots, h(\mathbf{x}_m) \right]^T \in \{0, 1\}^m : h \in H \right\} \right|.$$

415 **Definition 2.3** (VC-dimension). Let H be a class of functions from \mathcal{X} to $\{0, 1\}$. The
 416 VC-dimension of H , denoted by $\text{VCDim}(H)$, is the size of the largest shattered set,
 417 namely,

$$418 \quad \text{VCDim}(H) := \sup \left(\{0\} \cup \{m \in \mathbb{N}^+ : \Pi_H(m) = 2^m\} \right).$$

419 Let \mathcal{F} be a class of functions from \mathcal{X} to \mathbb{R} . The VC-dimension of \mathcal{F} , denoted by
 420 $\text{VCDim}(\mathcal{F})$, is defined by $\text{VCDim}(\mathcal{F}) := \text{VCDim}(\mathcal{T} \circ \mathcal{F})$, where

$$421 \quad \mathcal{T}(t) := \begin{cases} 1, & t \geq 0, \\ 0, & t < 0 \end{cases} \quad \text{and} \quad \mathcal{T} \circ \mathcal{F} := \{\mathcal{T} \circ f : f \in \mathcal{F}\}.$$

422 In particular, the expression ‘‘VC-dimension of a network (architecture)’’ means the VC-
 423 dimension of the function set that consists of all functions implemented by this network
 424 (architecture).

425 Recall that $C_u^s([0, 1]^d)$ denotes the unit ball of $C^s([0, 1]^d)$. Theorem 2.4 below shows
 426 that the best possible approximation error of functions in $C_u^s([0, 1]^d)$ approximated by
 427 functions in \mathcal{F} is bounded by a formula characterized by $\text{VCDim}(\mathcal{F})$.

428 **Theorem 2.4.** *Given any $s, d \in \mathbb{N}^+$, there exists a (small) positive constant $C_{s,d}$ deter-*
 429 *mined by s and d such that: For any $\varepsilon > 0$ and a function set \mathcal{F} with all elements defined*
 430 *on $[0, 1]^d$, if $\text{VCDim}(\mathcal{F}) \geq 1$ and*

$$431 \quad \inf_{\phi \in \mathcal{F}} \|\phi - f\|_{L^\infty([0,1]^d)} \leq \varepsilon \quad \text{for any } f \in C_u^s([0, 1]^d), \quad (2.2)$$

432 then $\text{VCDim}(\mathcal{F}) \geq C_{s,d} \varepsilon^{-d/s}$. ^④

^④In fact, $C_{s,d}$ can be expressed by s and d with a **explicitly** formula as we remark in the proof of this theorem. However, the formula may be very complicated.

433 This theorem demonstrates the connection between the VC-dimension of \mathcal{F} and
 434 the approximation error using elements of \mathcal{F} to approximate functions in $C_u^s([0, 1]^d)$.
 435 To be precise, the best possible approximation error is controlled by $\text{VCDim}(\mathcal{F})^{-s/d}$ up
 436 to a constant. It is shown in [22] that the VC-dimension of ReLU FNNs with a fixed
 437 architecture with W parameters and L layers has an upper bound $\mathcal{O}(WL \ln W)$. It
 438 follows that the VC-dimension of ReLU FNNs with width N and depth L is bounded
 439 by $\mathcal{O}(N^2 L \cdot L \cdot \ln(N^2 L)) \leq \mathcal{O}(N^2 L^2 \ln(NL))$. That is, $\text{VCDim}(\mathcal{F}) \leq \mathcal{O}(N^2 L^2 \ln(NL))$,
 440 where

$$441 \quad \mathcal{F} = \mathcal{NN}(\#\text{input} = d; \text{width} \leq N; \text{depth} \leq L; \#\text{output} = 1).$$

442 Hence, the approximation error of functions in $C_u^s([0, 1]^d)$, approximated by ReLU FNNs
 443 with width N and depth L , has a lower bound

$$444 \quad C(s, d) \cdot (N^2 L^2 \ln(NL))^{-s/d}$$

445 for some positive constant $C(s, d)$ determined by s and d . When the width and depth
 446 become $\mathcal{O}(N \ln N)$ and $\mathcal{O}(L \ln L)$, respectively, the lower bound of the approximation
 447 error becomes

$$448 \quad C(s, d) \cdot (N^2 L^2 (\ln N)^3 (\ln L)^3)^{-s/d}$$

449 for some positive constant $C(s, d)$ determined by s and d . These two lower bounds mean
 450 that our approximation errors in Theorem 1.1 and Corollary 1.2 are nearly optimal.

451 Now let us present the detailed proof of Theorem 2.4.

452 *Proof of Theorem 2.4.* To find a subset of \mathcal{F} shattering $\mathcal{O}(\varepsilon^{-d/s})$ points in $[0, 1]^d$, we
 453 divided the proof into two steps.

- 454 • Construct $\{f_\chi : \chi \in \mathcal{X}\} \subseteq C_u^s([0, 1]^d)$ that scatters $\mathcal{O}(\varepsilon^{-d/s})$ points, where \mathcal{X} is a
 455 function set defined later.
- 456 • Design $\phi_\chi \in \mathcal{F}$, for each $\chi \in \mathcal{X}$, based on f_χ and Equation (2.2) such that $\{\phi_\chi : \chi \in \mathcal{X}\} \subseteq \mathcal{F}$ also shatters $\mathcal{O}(\varepsilon^{-d/s})$ points.
 457

458 The details of these two steps can be found below.

459 **Step 1:** Construct $\{f_\chi : \chi \in \mathcal{X}\} \subseteq C_u^s([0, 1]^d)$ that scatters $\mathcal{O}(\varepsilon^{-d/s})$ points.

460 Let $K = \mathcal{O}(\varepsilon^{-1/s})$ be an integer determined later and divide $[0, 1]^d$ into K^d non-
 461 overlapping sub-cubes $\{Q_\beta\}_\beta$ as follows:

$$462 \quad Q_\beta := \{\mathbf{x} = [x_1, x_2, \dots, x_d]^T \in [0, 1]^d : x_i \in [\frac{\beta_i}{K}, \frac{\beta_i+1}{K}] \text{ for } i = 1, 2, \dots, d\}$$

463 for any index vector $\beta = [\beta_1, \beta_2, \dots, \beta_d]^T \in \{0, 1, \dots, K-1\}^d$.

464 There exists $\tilde{g} \in C^\infty(\mathbb{R}^d)$ such that $\tilde{g}(\mathbf{0}) = 1$ and $\tilde{g}(\mathbf{x}) = 0$ for $\|\mathbf{x}\|_2 \geq 1/3$.^⑤ Then,
 465 $g := \tilde{g}/\tilde{C}_{s,d} \in C_u^s([0, 1]^d)$ by setting $\tilde{C}_{s,d} := \|\tilde{g}\|_{C^s([0, 1]^d)} > 0$.

466 Define

$$467 \quad \mathcal{X} := \{\chi : \chi \text{ is a map from } \{0, 1, \dots, K-1\}^d \text{ to } \{-1, 1\}\}$$

^⑤In fact, such a function \tilde{g} is called “bump function”. An example can be attained by setting $\tilde{g}(\mathbf{x}) = C \exp(\frac{1}{\|\mathbf{x}\|_2^2 - 1})$ if $\|\mathbf{x}\|_2 < 1/3$ and $\tilde{g}(\mathbf{x}) = 0$ if $\|\mathbf{x}\|_2 \geq 1/3$, where C is a proper constant such that $\tilde{g}(\mathbf{0}) = 1$.

468 and

$$469 \quad g_\beta := K^{-s}g(K(\mathbf{x} - \mathbf{x}_{Q_\beta})) \quad \text{for each } \beta \in \{0, 1, \dots, K-1\}^d,$$

470 where \mathbf{x}_{Q_β} is the center of Q_β .

471 Next, for each $\chi \in \mathcal{X}$, we can define f_χ via

$$472 \quad f_\chi(\mathbf{x}) := \sum_{\beta \in \{0, 1, \dots, K-1\}^d} \chi(\beta)g_\beta(\mathbf{x}).$$

473 Then $f_\chi \in C_u^s([0, 1]^d)$ for each $\chi \in \mathcal{X}$, since it satisfies the following two conditions.

474 • By the definition of g_β and χ , we have

$$475 \quad \{\mathbf{x} : \chi(\beta)g_\beta(\mathbf{x}) \neq 0\} \subseteq \mathcal{B}(\mathbf{x}_{Q_\beta}, \frac{1}{3K}) \subseteq \frac{2}{3}Q_\beta \quad \text{for each } \beta \in \{0, 1, \dots, K-1\}^d,$$

476 which implies that $f_\chi \in C^\infty([0, 1]^d)$.

477 • For any $\mathbf{x} \in Q_\beta$, $\beta \in \{0, 1, \dots, K-1\}^d$, and $\boldsymbol{\alpha} \in \mathbb{N}^d$ with $\|\boldsymbol{\alpha}\|_1 \leq s$,

$$478 \quad \partial^\alpha f_\chi(\mathbf{x}) = \chi(\beta)\partial^\alpha g_\beta(\mathbf{x}) = K^{-s}\chi(\beta)K^{\|\boldsymbol{\alpha}\|_1}\partial^\alpha g(K(\mathbf{x} - \mathbf{x}_\beta)),$$

479 from which we deduce $|\partial^\alpha f_\chi(\mathbf{x})| = |K^{-(s-\|\boldsymbol{\alpha}\|_1)}\partial^\alpha g(K(\mathbf{x} - \mathbf{x}_\beta))| \leq 1$.

480 It is easy to check that $\{f_\chi : \chi \in \mathcal{X}\} \subseteq C_u^s([0, 1]^d)$ can shatter $K^d = \mathcal{O}(\varepsilon^{-d/s})$ points in
481 $[0, 1]^d$.

482 **Step 2:** Construct $\{\phi_\chi : \chi \in \mathcal{X}\}$ that also scatters $\mathcal{O}(\varepsilon^{-d/s})$ points.

483 By Equation (2.2), for each $\chi \in \mathcal{X}$, there exists $\phi_\chi \in \mathcal{F}$ such that

$$484 \quad \|\phi_\chi - f_\chi\|_{L^\infty([0, 1]^d)} \leq \varepsilon + \varepsilon/2.$$

485 Let $\mu(\cdot)$ denote the Lebesgue measure of a set. Then, for each $\chi \in \mathcal{X}$, there exists
486 $\mathcal{H}_\chi \subseteq [0, 1]^d$ with $\mu(\mathcal{H}_\chi) = 0$ such that

$$487 \quad |\phi_\chi(\mathbf{x}) - f_\chi(\mathbf{x})| \leq \frac{3}{2}\varepsilon \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \mathcal{H}_\chi.$$

488 Set $\mathcal{H} = \bigcup_{\chi \in \mathcal{X}} \mathcal{H}_\chi$; then we have $\mu(\mathcal{H}) = 0$ and

$$489 \quad |\phi_\chi(\mathbf{x}) - f_\chi(\mathbf{x})| \leq \frac{3}{2}\varepsilon \quad \text{for any } \chi \in \mathcal{X} \text{ and } \mathbf{x} \in [0, 1]^d \setminus \mathcal{H}. \quad (2.3)$$

490 Clearly, there exists $r \in (0, 1)$ such that

$$491 \quad g_\beta(\mathbf{x}) \geq \frac{1}{2}g_\beta(\mathbf{x}_{Q_\beta}) > 0 \quad \text{for any } \mathbf{x} \in rQ_\beta,$$

492 where \mathbf{x}_{Q_β} is the center of Q_β .

493 Note that $(rQ_\beta) \setminus \mathcal{H}$ is not empty, since $\mu((rQ_\beta) \setminus \mathcal{H}) > 0$ for each β . Then, for any
494 $\chi \in \mathcal{X}$ and $\beta \in \{0, 1, \dots, K-1\}^d$, there exists $\mathbf{x}_\beta \in (rQ_\beta) \setminus \mathcal{H}$ such that

$$495 \quad |f_\chi(\mathbf{x}_\beta)| = |g_\beta(\mathbf{x}_\beta)| \geq \frac{1}{2}|g_\beta(\mathbf{x}_{Q_\beta})| = \frac{1}{2}K^{-s}g(\mathbf{0}) = \frac{1}{2}K^{-s}/\tilde{C}_{s,d} \geq 2\varepsilon, \quad (2.4)$$

496 where the last inequality is attained by setting $K = \lfloor (4\varepsilon\tilde{C}_{s,d})^{-1/s} \rfloor$. Note that it is
497 necessary to verify $K \neq 0$; we do this later in the proof.

498 By Equations (2.3) and (2.4), we have, for each $\boldsymbol{\beta} \in \{0, 1, \dots, K-1\}^d$ and each $\chi \in \mathcal{X}$,

499
$$|f_\chi(\mathbf{x}_\beta)| \geq 2\varepsilon > \frac{3}{2}\varepsilon \geq |f_\chi(\mathbf{x}_\beta) - \phi_\chi(\mathbf{x}_\beta)|,$$

500 implying $f_\chi(\mathbf{x}_\beta)$ and $\phi_\chi(\mathbf{x}_\beta)$ have the same sign. Then $\{\phi_\chi : \chi \in \mathcal{X}\}$ shatters $\{\mathbf{x}_\beta : \beta \in$
 501 $\{0, 1, \dots, K-1\}^d\}$ since $\{f_\chi : \chi \in \mathcal{X}\}$ shatters $\{\mathbf{x}_\beta : \beta \in \{0, 1, \dots, K-1\}^d\}$. Hence,

502
$$\text{VCDim}(\mathcal{F}) \geq \text{VCDim}(\{\phi_\chi : \chi \in \mathcal{X}\}) \geq K^d = \lfloor (4\varepsilon\tilde{C}_{s,d})^{-1/s} \rfloor^d \geq 2^{-d}(4\varepsilon\tilde{C}_{s,d})^{-d/s},$$

503 where the last inequality comes from the fact that $\lfloor x \rfloor \geq x/2$ for any $x \in [1, \infty)$.

504 Finally, by setting

505
$$C_{s,d} = 2^{-d}(4\tilde{C}_{s,d})^{-d/s} = 2^{-d}(4\|\tilde{g}\|_{C^s([0,1]^d)})^{-d/s},$$

506 we have

507
$$\text{VCDim}(\mathcal{F}) \geq 2^{-d}(4\varepsilon\tilde{C}_{s,d})^{-d/s} = 2^{-d}(4\tilde{C}_{s,d})^{-d/s}\varepsilon^{-d/s} = C_{s,d}\varepsilon^{-d/s}$$

508 and

509
$$K = \lfloor (4\varepsilon\tilde{C}_{s,d})^{-1/s} \rfloor = \lfloor \varepsilon^{-1/s}(4\tilde{C}_{s,d})^{-1/s} \rfloor = \lfloor \varepsilon^{-1/s}(2^d C_{s,d})^{1/d} \rfloor \geq 1,$$

510 where the last inequality comes from the assumption $\varepsilon \leq (2^d C_{s,d})^{s/d}$. Such an assumption
 511 is reasonable since $\varepsilon > (2^d C_{s,d})^{s/d}$ is a trivial case, which implies

512
$$\text{VCDim}(\mathcal{F}) \geq 1 \geq 2^{-d} = C_{s,d} \left((2^d C_{s,d})^{s/d} \right)^{-d/s} > C_{s,d}\varepsilon^{-d/s}.$$

513 So we finish the proof. □

514 3 Proof of Theorem 2.1

515 Intuitively speaking, Theorem 2.1 shows that if a ReLU FNN can implement a
 516 function g approximating the target function f well except for the trifling region, then
 517 we can design a new ReLU network with a similar size to approximate f well on the
 518 whole domain. For example, if g approximates a one-dimensional continuous function
 519 f well except for a region in \mathbb{R} with a sufficiently small measure δ , then $\text{mid}(g(x +$
 520 $\delta), g(x), g(x - \delta))$ can approximate f well on the whole domain, where $\text{mid}(\cdot, \cdot, \cdot)$ is a
 521 function returning the middle value of three inputs and can be implemented via a ReLU
 522 FNN as shown in Lemma 3.1. This key idea is called the horizontal shift (translation)
 523 of g in this paper.

524 **Lemma 3.1.** *The middle value function $\text{mid}(x_1, x_2, x_3)$ can be implemented by a ReLU*
 525 *FNN with width 14 and depth 2.*

526 *Proof.* Recall the fact that

527
$$x = \sigma(x) - \sigma(-x) \quad \text{and} \quad |x| = \sigma(x) + \sigma(-x) \quad \text{for any } x \in \mathbb{R}. \quad (3.1)$$

528 Therefore,

529
$$\begin{aligned} \max(x, y) &= \frac{x + y + |x - y|}{2} \\ &= \frac{1}{2}\sigma(x + y) - \frac{1}{2}\sigma(-x - y) + \frac{1}{2}\sigma(x - y) + \frac{1}{2}\sigma(-x + y), \end{aligned} \quad (3.2)$$

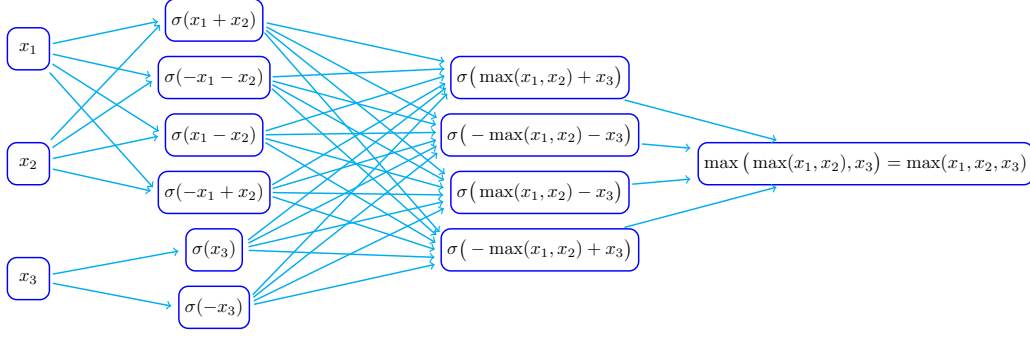


Figure 3: An illustration of the network architecture implementing $\max(x_1, x_2, x_3)$ based on Equations (3.1) and (3.2).

530 for any $x, y \in \mathbb{R}$. Thus, $\max(x_1, x_2, x_3)$ can be implemented by the network shown in
 531 Figure 3.

532 Clearly,

533
$$\max(x_1, x_2, x_3) \in \mathcal{NN}(\#input = 3; \text{widthvec} = [6, 4]).$$

534 Similarly, we have

535
$$\min(x_1, x_2, x_3) \in \mathcal{NN}(\#input = 3; \text{widthvec} = [6, 4]).$$

536 It is easy to check that

537
$$\begin{aligned} \text{mid}(x_1, x_2, x_3) &= x_1 + x_2 + x_3 - \max(x_1, x_2, x_3) - \min(x_1, x_2, x_3) \\ &= \sigma(x_1 + x_2 + x_3) - \sigma(-x_1 - x_2 - x_3) - \max(x_1, x_2, x_3) - \min(x_1, x_2, x_3). \end{aligned}$$

538 Hence,

539
$$\text{mid}(x_1, x_2, x_3) \in \mathcal{NN}(\#input = 3; \text{widthvec} = [14, 10]).$$

540 That is, $\text{mid}(x_1, x_2, x_3)$ can be implemented by a ReLU FNN with width 14 and depth
 541 2. So we finish the proof. \square

542 The next lemma shows a simple but useful property of the $\text{mid}(x_1, x_2, x_3)$ function
 543 that helps to exclude poor approximation in the trifling region.

544 **Lemma 3.2.** *For any $\varepsilon > 0$, if at least two elements of $\{x_1, x_2, x_3\}$ are in $\mathcal{B}(y, \varepsilon)$, then*
 545 $\text{mid}(x_1, x_2, x_3) \in \mathcal{B}(y, \varepsilon)$.

546 *Proof.* Without loss of generality, we may assume $x_1, x_2 \in \mathcal{B}(y, \varepsilon)$ and $x_1 \leq x_2$. Then the
 547 proof can be divided into three cases.

548 1. If $x_3 < x_1$, then $x_3 < x_1 \leq x_2$, implying $\text{mid}(x_1, x_2, x_3) = x_1 \in \mathcal{B}(y, \varepsilon)$.

549 2. If $x_1 \leq x_3 \leq x_2$, then $\text{mid}(x_1, x_2, x_3) = x_3 \in \mathcal{B}(y, \varepsilon)$ since $y - \varepsilon \leq x_1 \leq x_3 \leq x_2 \leq y + \varepsilon$.

550 3. If $x_2 < x_3$, then $x_1 \leq x_2 < x_3$, implying $\text{mid}(x_1, x_2, x_3) = x_2 \in \mathcal{B}(y, \varepsilon)$.

551 So we finish the proof. \square

552 Next, given a function g approximating f well on $[0, 1]$ except for the trifling region,
 553 Lemma 3.3 below shows how to use the $\text{mid}(x_1, x_2, x_3)$ function to construct a new
 554 function ϕ uniformly approximating f well on $[0, 1]$, leveraging the useful property of
 555 $\text{mid}(x_1, x_2, x_3)$ in Lemma 3.2.

556 **Lemma 3.3.** *Given any $\varepsilon > 0$, $K \in \mathbb{N}^+$, and $\delta \in (0, \frac{1}{3K}]$, assume $f \in C([0, 1])$ and
 557 $g : \mathbb{R} \rightarrow \mathbb{R}$ is a general function with*

$$558 \quad |g(x) - f(x)| \leq \varepsilon, \text{ i.e., } g(x) \in \mathcal{B}(f(x), \varepsilon) \quad \text{for any } x \in [0, 1] \setminus \Omega([0, 1], K, \delta). \quad (3.3)$$

559 Then

$$560 \quad |\phi(x) - f(x)| \leq \varepsilon + \omega_f(\delta) \quad \text{for any } x \in [0, 1],$$

561 where

$$562 \quad \phi(x) := \text{mid}(g(x - \delta), g(x), g(x + \delta)) \quad \text{for any } x \in \mathbb{R}.$$

563 *Proof.* Divide $[0, 1]$ into K small intervals denoted by $Q_k = [\frac{k}{K}, \frac{k+1}{K}]$ for $k = 0, 1, \dots, K-1$.
 564 For each $k \in \{0, 1, \dots, K-1\}$, we further divide Q_k into four small closed intervals as
 565 shown in Figure 4, i.e.,

$$566 \quad Q_k = Q_{k,1} \cup Q_{k,2} \cup Q_{k,3} \cup Q_{k,4},$$

567 where $Q_{k,1} = [\frac{k}{K}, \frac{k}{K} + \delta]$, $Q_{k,2} = [\frac{k}{K} + \delta, \frac{k+1}{K} - 2\delta]$, $Q_{k,3} = [\frac{k+1}{K} - 2\delta, \frac{k+1}{K} - \delta]$, and $Q_{k,4} =$
 568 $[\frac{k+1}{K} - \delta, \frac{k+1}{K}]$.

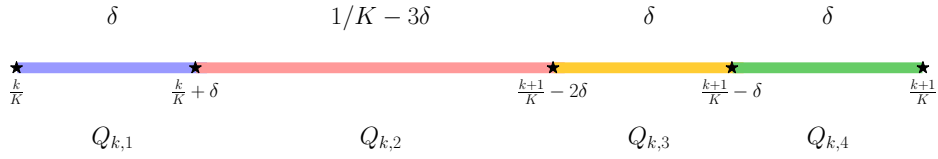


Figure 4: An illustration of $Q_{k,i}$ for $i = 1, 2, 3, 4$.

569 It is easy to verify that

- 570 • $Q_{k,i} \subseteq [0, 1] \setminus \Omega([0, 1], K, \delta)$ for $k = 0, 1, \dots, K-1$ and $i = 1, 2, 3$;
- 571 • $Q_{K-1,4} \subseteq [0, 1] \setminus \Omega([0, 1], K, \delta)$.

572 To estimate the difference between $\phi(x)$ and $f(x)$, we consider the following four
 573 cases of x in $[0, 1]$ for each $k \in \{0, 1, \dots, K-1\}$.

574 **Case 1:** $x \in Q_{k,1}$.

575 If $x \in Q_{k,1}$, then $x \in [0, 1] \setminus \Omega([0, 1], K, \delta)$ and

$$576 \quad x + \delta \in Q_{k,2} \cup Q_{k,3} \subseteq [0, 1] \setminus \Omega([0, 1], K, \delta).$$

577 It follows from Equation (3.3) that

$$578 \quad g(x) \in \mathcal{B}(f(x), \varepsilon) \subseteq \mathcal{B}(f(x), \varepsilon + \omega_f(\delta))$$

579 and

$$580 \quad g(x + \delta) \in \mathcal{B}(f(x + \delta), \varepsilon) \subseteq \mathcal{B}(f(x), \varepsilon + \omega_f(\delta)).$$

581 By Lemma 3.2, we get

$$582 \quad \text{mid}(g(x - \delta), g(x), g(x + \delta)) \in \mathcal{B}(f(x), \varepsilon + \omega_f(\delta)).$$

583 **Case 2:** $x \in Q_{k,2}$.

584 If $x \in Q_{k,2}$, then

$$585 \quad x - \delta, x, x + \delta \in Q_{k,1} \cup Q_{k,2} \cup Q_{k,3} \subseteq [0, 1] \setminus \Omega([0, 1], K, \delta).$$

586 It follows from Equation (3.3) that

$$587 \quad g(x - \delta) \in \mathcal{B}(f(x - \delta), \varepsilon) \subseteq \mathcal{B}(f(x), \varepsilon + \omega_f(\delta)),$$

588

$$589 \quad g(x) \in \mathcal{B}(f(x), \varepsilon) \subseteq \mathcal{B}(f(x), \varepsilon + \omega_f(\delta)),$$

590 and

$$591 \quad g(x + \delta) \in \mathcal{B}(f(x + \delta), \varepsilon) \subseteq \mathcal{B}(f(x), \varepsilon + \omega_f(\delta)).$$

592 Then, by Lemma 3.2, we have

$$593 \quad \text{mid}(g(x - \delta), g(x), g(x + \delta)) \in \mathcal{B}(f(x), \varepsilon + \omega_f(\delta)).$$

594 **Case 3:** $x \in Q_{k,3}$.

595 If $x \in Q_{k,3}$, then $x \in [0, 1] \setminus \Omega([0, 1], K, \delta)$ and

$$596 \quad x - \delta \in Q_{k,1} \cup Q_{k,2} \subseteq [0, 1] \setminus \Omega([0, 1], K, \delta).$$

597 It follows from Equation (3.3) that

$$598 \quad g(x) \in \mathcal{B}(f(x), \varepsilon) \subseteq \mathcal{B}(f(x), \varepsilon + \omega_f(\delta))$$

599 and

$$600 \quad g(x - \delta) \in \mathcal{B}(f(x - \delta), \varepsilon) \subseteq \mathcal{B}(f(x), \varepsilon + \omega_f(\delta)).$$

601 By Lemma 3.2, we get

$$602 \quad \text{mid}(g(x - \delta), g(x), g(x + \delta)) \in \mathcal{B}(f(x), \varepsilon + \omega_f(\delta)).$$

603 **Case 4:** $x \in Q_{k,4}$.

604 If $x \in Q_{k,4}$, we can divide this case into two sub-cases.

- 605 • If $k \in \{0, 1, \dots, K - 2\}$, then $x - \delta \in Q_{k,3} \in [0, 1] \setminus \Omega([0, 1], K, \delta)$ and $x + \delta \in Q_{k+1,1} \subseteq$
606 $[0, 1] \setminus \Omega([0, 1], K, \delta)$. It follows from Equation (3.3) that

$$607 \quad g(x - \delta) \in \mathcal{B}(f(x - \delta), \varepsilon) \subseteq \mathcal{B}(f(x), \varepsilon + \omega_f(\delta))$$

608 and

$$609 \quad g(x + \delta) \in \mathcal{B}(f(x + \delta), \varepsilon) \subseteq \mathcal{B}(f(x), \varepsilon + \omega_f(\delta)).$$

610 By Lemma 3.2, we get

$$611 \quad \text{mid}(g(x - \delta), g(x), g(x + \delta)) \in \mathcal{B}(f(x), \varepsilon + \omega_f(\delta)).$$

612 • If $k = K - 1$, then $x \in Q_{k,4} = Q_{K-1,4} \subseteq [0, 1] \setminus \Omega([0, 1], K, \delta)$ and $x - \delta \in Q_{k,3} \subseteq$
 613 $[0, 1] \setminus \Omega([0, 1], K, \delta)$. It follows from Equation (3.3) that

$$614 \quad g(x) \in \mathcal{B}(f(x), \varepsilon) \subseteq \mathcal{B}(f(x), \varepsilon + \omega_f(\delta))$$

615 and

$$616 \quad g(x - \delta) \in \mathcal{B}(f(x - \delta), \varepsilon) \subseteq \mathcal{B}(f(x), \varepsilon + \omega_f(\delta)).$$

617 By Lemma 3.2, we get

$$618 \quad \text{mid}(g(x - \delta), g(x), g(x + \delta)) \in \mathcal{B}(f(x), \varepsilon + \omega_f(\delta)).$$

619 Since $[0, 1] = \bigcup_{k=0}^{K-1} \left(\bigcup_{i=1}^4 Q_{k,i} \right)$, we have

$$620 \quad \text{mid}(g(x - \delta), g(x), g(x + \delta)) \in \mathcal{B}(f(x), \varepsilon + \omega_f(\delta)) \quad \text{for any } x \in [0, 1].$$

621 Recall that $\phi(x) = \text{mid}(g(x - \delta), g(x), g(x + \delta))$. Then we have

$$622 \quad |\phi(x) - f(x)| \leq \varepsilon + \omega_f(\delta) \quad \text{for any } x \in [0, 1].$$

623 So we finish the proof. □

624 The next lemma below extend Lemma 3.3 to the multidimensional case.

625 **Lemma 3.4.** *Given any $\varepsilon > 0$, $K \in \mathbb{N}^+$, and $\delta \in (0, \frac{1}{3K}]$, assume $f \in C([0, 1]^d)$ and*
 626 *$g : \mathbb{R}^d \rightarrow \mathbb{R}$ is a general function with*

$$627 \quad |g(\mathbf{x}) - f(\mathbf{x})| \leq \varepsilon, \text{ i.e., } g(\mathbf{x}) \in \mathcal{B}(f(\mathbf{x}), \varepsilon) \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta).$$

628 *Then*

$$629 \quad |\phi(\mathbf{x}) - f(\mathbf{x})| \leq \varepsilon + d \cdot \omega_f(\delta) \quad \text{for any } \mathbf{x} \in [0, 1]^d,$$

630 *where $\phi = \phi_d$ is defined by induction through*

$$631 \quad \phi_{i+1}(\mathbf{x}) := \text{mid}(\phi_i(\mathbf{x} - \delta \mathbf{e}_{i+1}), \phi_i(\mathbf{x}), \phi_i(\mathbf{x} + \delta \mathbf{e}_{i+1})) \quad \text{for } i = 0, 1, \dots, d-1, \quad (3.4)$$

632 *where $\phi_0 = g$ and $\{\mathbf{e}_i\}_{i=1}^d$ is the standard basis in \mathbb{R}^d .*

633 *Proof.* For $\ell = 0, 1, \dots, d$, we define

$$634 \quad E_\ell := \left\{ \mathbf{x} = [x_1, x_2, \dots, x_d]^T : x_i \in \begin{cases} [0, 1], & \text{if } i \leq \ell, \\ [0, 1] \setminus \Omega([0, 1], K, \delta), & \text{if } i > \ell \end{cases} \right\}.$$

635 Clearly, $E_0 = [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)$ and $E_d = [0, 1]^d$. See Figure 5 for the illustrations of
 636 E_ℓ for $\ell = 0, 1, \dots, d$ when $K = 4$ and $d = 2$.

637 We would like to construct a sequence of functions $\phi_0, \phi_1, \dots, \phi_d$ by induction, based
 638 on Equation (3.4), such that, for each $\ell \in \{0, 1, \dots, d\}$,

$$639 \quad \phi_\ell(\mathbf{x}) \in \mathcal{B}(f(\mathbf{x}), \varepsilon + \ell \cdot \omega_f(\delta)) \quad \text{for any } \mathbf{x} \in E_\ell. \quad (3.5)$$

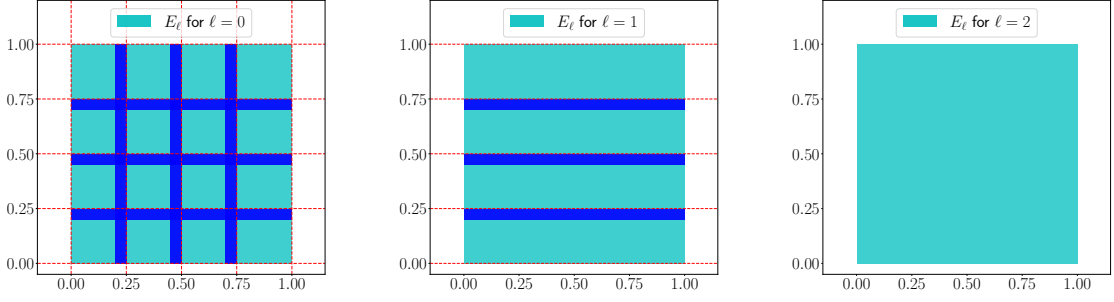


Figure 5: Illustrations of E_ℓ for $\ell = 0, 1, 2$ when $K = 4$ and $d = 2$.

640 Let us first consider the case $\ell = 0$. Note that $\phi_0 = g$, $E_0 = [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)$,
 641 and $|g(\mathbf{x}) - f(\mathbf{x})| \leq \varepsilon$ for any $\mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)$. Then we have

$$642 \quad \phi_0(\mathbf{x}) = g(\mathbf{x}) \in \mathcal{B}(f(\mathbf{x}), \varepsilon) \quad \text{for any } \mathbf{x} \in E_0.$$

643 That is, Equation (3.5) is true for $\ell = 0$.

644 Now assume Equation (3.5) is true for $\ell = i$. We will prove that it also holds for
 645 $\ell = i + 1$. By the hypothesis of induction, we have

$$646 \quad \phi_i(x_1, \dots, x_i, t, x_{i+2}, \dots, x_d) \in \mathcal{B}(f(x_1, \dots, x_i, t, x_{i+2}, \dots, x_d), \varepsilon + i \cdot \omega_f(\delta)) \quad (3.6)$$

647 for any $x_1, \dots, x_i \in [0, 1]$ and $t, x_{i+2}, \dots, x_d \in [0, 1] \setminus \Omega([0, 1], K, \delta)$.

648 For fixed $x_1, \dots, x_i \in [0, 1]$ and $x_{i+2}, \dots, x_d \in [0, 1] \setminus \Omega([0, 1], K, \delta)$, denote

$$649 \quad \mathbf{x}^{[i]} := [x_1, \dots, x_i, x_{i+2}, \dots, x_d]^T \in [0, 1]^{d-1}.$$

650 Then define

$$651 \quad \psi_{\mathbf{x}^{[i]}}(t) := \phi_i(x_1, \dots, x_i, t, x_{i+2}, \dots, x_d) \quad \text{for any } t \in \mathbb{R}$$

652 and

$$653 \quad f_{\mathbf{x}^{[i]}}(t) := f(x_1, \dots, x_i, t, x_{i+2}, \dots, x_d) \quad \text{for any } t \in \mathbb{R}.$$

654 It follows from Equation (3.6) that

$$655 \quad \psi_{\mathbf{x}^{[i]}}(t) \in \mathcal{B}(f_{\mathbf{x}^{[i]}}(t), \varepsilon + i \cdot \omega_f(\delta)) \quad \text{for any } t \in [0, 1] \setminus \Omega([0, 1], K, \delta).$$

656 Then by Lemma 3.3 (set $g = \psi_{\mathbf{x}^{[i]}}$ and $f = f_{\mathbf{x}^{[i]}}$ therein), we get, for any $t \in [0, 1]$,

$$657 \quad \begin{aligned} \text{mid}(\psi_{\mathbf{x}^{[i]}}(t - \delta), \psi_{\mathbf{x}^{[i]}}(t), \psi_{\mathbf{x}^{[i]}}(t + \delta)) &\in \mathcal{B}(f_{\mathbf{x}^{[i]}}(t), \varepsilon + i \cdot \omega_f(\delta) + \omega_{f_{\mathbf{x}^{[i]}}}(\delta)) \\ &\subseteq \mathcal{B}(f_{\mathbf{x}^{[i]}}(t), \varepsilon + (i + 1)\omega_f(\delta)). \end{aligned}$$

658 That is, for any $x_{i+1} = t \in [0, 1]$,

$$659 \quad \begin{aligned} &\text{mid}(\phi_i(x_1, \dots, x_i, x_{i+1} - \delta, x_{i+2}, \dots, x_d), \phi_i(x_1, \dots, x_d), \phi_i(x_1, \dots, x_i, x_{i+1} + \delta, x_{i+2}, \dots, x_d)) \\ &\in \mathcal{B}(f(x_1, \dots, x_d), \varepsilon + (i + 1)\omega_f(\delta)). \end{aligned}$$

660 Note that $x_1, \dots, x_i \in [0, 1]$, $x_{i+1} = t \in [0, 1]$, and $x_{i+2}, \dots, x_d \in [0, 1] \setminus \Omega([0, 1], K, \delta)$ are
 661 arbitrary. Thus, for any $\mathbf{x} \in E_{i+1}$, we have

$$662 \quad \text{mid}(\phi_i(\mathbf{x} - \delta \mathbf{e}_{i+1}), \phi_i(\mathbf{x}), \phi_i(\mathbf{x} + \delta \mathbf{e}_{i+1})) \in \mathcal{B}(f(\mathbf{x}), \varepsilon + (i+1)\omega_f(\delta)),$$

663 which implies

$$664 \quad \phi_{i+1}(\mathbf{x}) \in \mathcal{B}(f(\mathbf{x}), \varepsilon + (i+1)\omega_f(\delta)) \quad \text{for any } \mathbf{x} \in E_{i+1}.$$

665 So Equation (3.5) is true for $\ell = i+1$, which means we finish the process of mathematical
 666 induction.

667 By the principle of induction, we have

$$668 \quad \phi(\mathbf{x}) := \phi_d(\mathbf{x}) \in \mathcal{B}(f(\mathbf{x}), \varepsilon + d \cdot \omega_f(\delta)) \quad \text{for any } \mathbf{x} \in E_d = [0, 1]^d.$$

669 Therefore,

$$670 \quad |\phi(\mathbf{x}) - f(\mathbf{x})| \leq \varepsilon + d \cdot \omega_f(\delta) \quad \text{for any } \mathbf{x} \in [0, 1]^d,$$

671 which means we finish the proof. □

672 With Lemma 3.4 in hand, we are ready to prove Theorem 2.1.

673 *Proof of Theorem 2.1.* Set $\phi_0 = \tilde{\phi}$ and define ϕ_i for $i \in \{1, 2, \dots, d\}$ by induction as follows:

$$674 \quad \phi_{i+1}(\mathbf{x}) := \text{mid}(\phi_i(\mathbf{x} - \delta \mathbf{e}_{i+1}), \phi_i(\mathbf{x}), \phi_i(\mathbf{x} + \delta \mathbf{e}_{i+1})) \quad \text{for } i = 0, 1, \dots, d-1,$$

675 where $\{\mathbf{e}_i\}_{i=1}^d$ is the standard basis in \mathbb{R}^d . Then by Lemma 3.4 with $\phi = \phi_d$, we have

$$676 \quad |\phi(\mathbf{x}) - f(\mathbf{x})| \leq \varepsilon + d \cdot \omega_f(\delta) \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

677 It remains to determine the network architecture implementing $\phi = \phi_d$. Clearly, $\phi_0 = \tilde{\phi} \in$
 678 $\mathcal{NN}(\text{width} \leq N; \text{depth} \leq L)$ implies

$$679 \quad \phi_0(\cdot - \delta \mathbf{e}_1), \phi_0(\cdot), \phi_0(\cdot + \delta \mathbf{e}_1) \in \mathcal{NN}(\text{width} \leq N; \text{depth} \leq L).$$

680 By defining a vector-valued function $\Phi_0 : \mathbb{R}^d \rightarrow \mathbb{R}^3$ as

$$681 \quad \Phi_0(\mathbf{x}) := (\phi_0(\mathbf{x} - \delta \mathbf{e}_1), \phi_0(\mathbf{x}), \phi_0(\mathbf{x} + \delta \mathbf{e}_1)) \quad \text{for any } \mathbf{x} \in \mathbb{R}^d,$$

682 we have $\Phi_0 \in \mathcal{NN}(\#\text{input} = d; \text{width} \leq 3N; \text{depth} \leq L; \#\text{output} = 3)$. Recall that
 683 $\text{mid}(\cdot, \cdot, \cdot) \in \mathcal{NN}(\text{width} \leq 14; \text{depth} \leq 2)$ by Lemma 3.1. Therefore, $\phi_1 = \min(\cdot, \cdot, \cdot) \circ \Phi_0$
 684 can be implemented by a ReLU FNN with width $\max\{3N, 14\} \leq 3(N+4)$ and depth
 685 $L+2$. Similarly, $\phi = \phi_d$ can be implemented by a ReLU FNN with width $3^d(N+4)$ and
 686 depth $L+2d$. So we finish the proof. □

687 4 Proof of Theorem 2.2

688 In this section, we prove Theorem 2.2, a weaker version of the main theorem of
 689 this paper (Theorem 1.1) targeting a ReLU FNN constructed to approximate a smooth
 690 function outside the trifling region. The main idea is to construct ReLU FNNs through
 691 Taylor expansions of smooth functions. We first discuss the proof sketch in Section 4.1
 692 and give the detailed proof in Section 4.2.

693 4.1 Proof sketch of Theorem 2.2

694 Set $K = \mathcal{O}(N^{2/d}L^{2/d})$ and let $\Omega([0,1]^d, K, \delta)$ partition $[0,1]^d$ into K^d cubes Q_β
695 for $\beta \in \{0, 1, \dots, K-1\}^d$. As we shall see later, the introduction of the trifling region
696 $\Omega([0,1]^d, K, \delta)$ can reduce the difficulty in constructing ReLU FNNs to achieve the op-
697 timal approximation error simultaneously in width and depth, since it is only required
698 to uniformly control the approximation error outside the trifling region and there is
699 no requirement for the ReLU FNN inside the trifling region. In particular, for each
700 $\beta = [\beta_1, \beta_2, \dots, \beta_d]^T \in \{0, 1, \dots, K-1\}^d$, we define $\mathbf{x}_\beta := \beta/K$ and

$$701 \quad Q_\beta := \{\mathbf{x} = [x_1, x_2, \dots, x_d]^T : x_i \in [\frac{\beta_i}{K}, \frac{\beta_i+1}{K} - \delta \cdot \mathbf{1}_{\{\beta_i \leq K-2\}}]\text{ for } i = 1, 2, \dots, d\}.$$

702 Clearly, $[0,1]^d = \Omega([0,1]^d, K, \delta) \cup (\cup_{\beta \in \{0,1,\dots,K-1\}^d} Q_\beta)$ and \mathbf{x}_β is the vertex of Q_β with
minimum $\|\cdot\|_1$ norm. See Figure 6 for the illustrations of Q_β and \mathbf{x}_β .

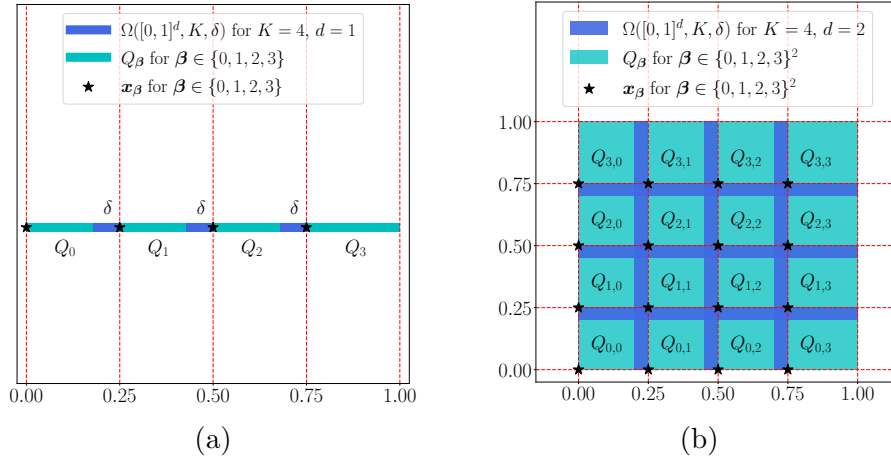


Figure 6: Illustrations of $\Omega([0,1]^d, K, \delta)$, Q_β , and \mathbf{x}_β for $\beta \in \{0, 1, \dots, K-1\}^d$. (a) $K = 4$ and $d = 1$. (b) $K = 4$ and $d = 2$.

703 For any $\beta \in \{0, 1, \dots, K-1\}^d$ and $\mathbf{x} \in Q_\beta$, there exists $\xi_{\mathbf{x}} \in (0, 1)$ such that
704

$$705 \quad f(\mathbf{x}) = \underbrace{\sum_{\|\alpha\|_1 \leq s-1} \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \mathbf{h}^\alpha}_{\mathcal{T}_1} + \underbrace{\sum_{\|\alpha\|_1 = s} \frac{\partial^\alpha f(\mathbf{x}_\beta + \xi_{\mathbf{x}} \mathbf{h})}{\alpha!} \mathbf{h}^\alpha}_{\mathcal{T}_2} =: \mathcal{T}_1 + \mathcal{T}_2, \quad \textcircled{6} \quad (4.1)$$

706 where $\mathbf{h}(\mathbf{x}) = \mathbf{x} - \mathbf{x}_\beta = \mathbf{x} - \beta/K$. Clearly, the magnitude of \mathcal{T}_2 is bounded by $\mathcal{O}(K^{-s}) =$
707 $\mathcal{O}(N^{-2s/d}L^{-2s/d})$. So we only need to construct a ReLU FNN with width $\mathcal{O}(N \ln N)$ and
708 depth $\mathcal{O}(L \ln L)$ to approximate

$$709 \quad \mathcal{T}_1 = \sum_{\|\alpha\|_1 \leq s-1} \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \mathbf{h}^\alpha$$

710 within an error $\mathcal{O}(N^{-2s/d}L^{-2s/d})$. To approximate \mathcal{T}_1 well by ReLU FNNs, we need three
711 key steps as follows.

$\textcircled{6} \sum_{\|\alpha\|_1 = s}$ is short for $\sum_{\|\alpha\|_1 = s, \alpha \in \mathbb{N}^d}$. The same notation is used throughout this paper.

- 712 (i) Construct a ReLU FNN to implement a function $P_{\alpha} : \mathbb{R}^d \rightarrow \mathbb{R}$ approximating the
713 polynomial \mathbf{h}^{α} well for each $\alpha \in \mathbb{N}^d$ with $\|\alpha\|_1 \leq s - 1$.
- 714 (ii) Construct a ReLU FNN to implement a vector-valued function $\Psi : \mathbb{R}^d \rightarrow \mathbb{R}^d$ pro-
715 jecting the whole cube Q_{β} to a point $\mathbf{x}_{\beta} = \frac{\beta}{K}$, i.e., $\Psi(\mathbf{x}) = \mathbf{x}_{\beta}$ for any $\mathbf{x} \in Q_{\beta}$ and
716 each $\beta \in \{0, 1, \dots, K - 1\}^d$.
- 717 (iii) Construct a ReLU FNN to implement a function $\phi_{\alpha} : \mathbb{R}^d \rightarrow \mathbb{R}$ approximating $\partial^{\alpha} f$
718 via solving a point fitting problem, i.e., ϕ_{α} should fit $\partial^{\alpha} f$ well at all points in
719 $\{\mathbf{x}_{\beta} : \beta \in \{0, 1, \dots, K - 1\}^d\}$ for each $\alpha \in \mathbb{N}^d$ with $\|\alpha\|_1 \leq s - 1$. That is, for each
720 $\alpha \in \mathbb{N}^d$ with $\|\alpha\|_1 \leq s - 1$, we need to design ϕ_{α} satisfying

$$721 \quad |\phi_{\alpha}(\mathbf{x}_{\beta}) - \partial^{\alpha} f(\mathbf{x}_{\beta})| \leq \mathcal{O}(N^{-2s/d} L^{-2s/d}) \quad \text{for any } \beta \in \{0, 1, \dots, K - 1\}^d. \quad (4.2)$$

722 We will establish three propositions corresponding to these three steps above. They
723 will be applied to support the construction of the desired ReLU FNNs. Their proofs will
724 be available in Section 5.

725 First, we establish a general proposition, Proposition 4.1 below, showing how to use
726 ReLU FNNs to approximate multivariate polynomials. With Proposition 4.1 in hand,
727 Step (i) is straightforward.

728 **Proposition 4.1.** *Assume $P(\mathbf{x}) = \mathbf{x}^{\alpha} = x_1^{\alpha_1} x_2^{\alpha_2} \dots x_d^{\alpha_d}$ for $\alpha \in \mathbb{N}^d$ with $\|\alpha\|_1 \leq k \in \mathbb{N}^+$.
729 For any $N, L \in \mathbb{N}^+$, there exists a function ϕ implemented by a ReLU FNN with width
730 $9(N + 1) + k - 1$ and depth $7k^2 L$ such that*

$$731 \quad |\phi(\mathbf{x}) - P(\mathbf{x})| \leq 9k(N + 1)^{-7kL} \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

732 Proposition 4.1 shows that ReLU FNNs with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ are
733 able to approximate polynomials with an error $\mathcal{O}(N^{-L})$. This reveals the power of
734 depth in ReLU FNNs for approximating polynomials, from the perspective of function
735 compositions. The starting point of a good approximation of functions is to approximate
736 polynomials with high accuracy. In classical approximation theory, the approximation
737 power of any numerical scheme depends on the degree of polynomials that can be locally
738 reproduced. Being able to approximate polynomials by ReLU FNNs with high accuracy
739 plays a vital role in the proof of Theorem 1.1. It is interesting to study whether there
740 is any other function space with reasonable size, besides polynomial space, having an
741 exponential error $\mathcal{O}(N^{-L})$ when approximated by ReLU FNNs. Obviously, the space of
742 smooth functions is too big due to the optimality of Theorem 1.1 as shown in Section 2.3.

743 Proposition 4.1 can be generalized to the case of polynomials defined on an arbitrary
744 hypercube $[a, b]^d$. Let us give an example for the polynomial xy below. Its proof will be
745 provided later in Section 5.1.

746 **Lemma 4.2.** *For any $N, L \in \mathbb{N}^+$ and $a, b \in \mathbb{R}$ with $a < b$, there exists a function ϕ
747 implemented by a ReLU FNN with width $9N + 1$ and depth L such that*

$$748 \quad |\phi(x, y) - xy| \leq 6(b - a)^2 N^{-L} \quad \text{for any } x, y \in [a, b].$$

749 Second, our goal is to construct a step function Ψ mapping $\mathbf{x} \in Q_\beta$ to $\mathbf{x}_\beta = \frac{\beta}{K}$ for any
750 $\beta \in \{0, 1, \dots, K-1\}^d$. We only need to approximate one-dimensional step functions, be-
751 cause in the multidimensional case we can simply set $\Psi(\mathbf{x}) = [\psi(x_1), \psi(x_2), \dots, \psi(x_d)]^T$,
752 where ψ is a one-dimensional step function. Therefore, to implement Step (ii), we
753 need to construct ReLU FNNs with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ to approximate one-
754 dimensional step functions with $\mathcal{O}(K) = \mathcal{O}(N^{2/d}L^{2/d})$ “steps” as shown in Proposition 4.3
755 below.

756 **Proposition 4.3.** *For any $N, L, d \in \mathbb{N}^+$ and $\delta \in (0, \frac{1}{3K}]$ with $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor$, there
757 exists a one-dimensional function ϕ implemented by a ReLU FNN with width $4\lfloor N^{1/d} \rfloor + 3$
758 and depth $4L + 5$ such that*

$$759 \quad \phi(x) = k \quad \text{if } x \in \left[\frac{k}{K}, \frac{k+1}{K} - \delta \cdot \mathbb{1}_{\{k \leq K-2\}} \right] \quad \text{for } k = 0, 1, \dots, K-1.$$

760 Next, the aim of Step (iii) is to construct ϕ_α implemented by a ReLU FNN such that
761 Equation (4.2) holds for each α . To this end, we establish a proposition, Proposition 4.4
762 below, to show that ReLU FNNs with width $\mathcal{O}(sN \ln N)$ and depth $\mathcal{O}(L \ln L)$ can be
763 constructed to fit $N^2 L^2$ points within an error $N^{-2s} L^{-2s}$.

764 **Proposition 4.4.** *Given any $N, L, s \in \mathbb{N}^+$ and $\xi_i \in [0, 1]$ for $i = 0, 1, \dots, N^2 L^2 - 1$, there
765 exists a function ϕ implemented by a ReLU FNN with width $16s(N+1) \log_2(8N)$ and
766 depth $5(L+2) \log_2(4L)$ such that*

$$767 \quad (i) \quad |\phi(i) - \xi_i| \leq N^{-2s} L^{-2s} \quad \text{for } i = 0, 1, \dots, N^2 L^2 - 1;$$

$$768 \quad (ii) \quad 0 \leq \phi(x) \leq 1 \quad \text{for any } x \in \mathbb{R}.$$

769 The proofs of Propositions 4.1, 4.3, and 4.4 can be found in Sections 5.1, 5.2, and
770 5.3, respectively. The main ideas of proving Theorem 1.1 are summarized in Table 2.

Table 2: A list of sub-networks for approximating smooth functions. Recall that $\mathbf{h} = \mathbf{x} - \Psi(\mathbf{x}) = \mathbf{x} - \mathbf{x}_\beta$ for $\mathbf{x} \in Q_\beta$.

| target function | function implemented by network | width | depth | approximation error |
|--|---|------------------------|------------------------|---|
| step function | $\Psi(\mathbf{x})$ | $\mathcal{O}(N)$ | $\mathcal{O}(L)$ | no error outside $\Omega([0, 1]^d, K, \delta)$ |
| $x_1 x_2$ | $\varphi(x_1, x_2)$ | $\mathcal{O}(N)$ | $\mathcal{O}(L)$ | $\mathcal{E}_1 = 216(N+1)^{-2s(L+1)}$ |
| \mathbf{h}^α | $P_\alpha(\mathbf{h})$ | $\mathcal{O}(N)$ | $\mathcal{O}(L)$ | $\mathcal{E}_2 = 9s(N+1)^{-7sL}$ |
| $\partial^\alpha f(\Psi(\mathbf{x}))$ | $\phi_\alpha(\Psi(\mathbf{x}))$ | $\mathcal{O}(N \ln N)$ | $\mathcal{O}(L \ln L)$ | $\mathcal{E}_3 = 2N^{-2s} L^{-2s}$ |
| $\sum_{\ \alpha\ \leq s-1} \frac{\partial^\alpha f(\Psi(\mathbf{x}))}{\alpha!} \mathbf{h}^\alpha$ | $\sum_{\ \alpha\ \leq s-1} \varphi\left(\frac{\phi_\alpha(\Psi(\mathbf{x}))}{\alpha!}, P_\alpha(\mathbf{h})\right)$ | $\mathcal{O}(N \ln N)$ | $\mathcal{O}(L \ln L)$ | $\mathcal{O}(\mathcal{E}_1 + \mathcal{E}_2 + \mathcal{E}_3)$ |
| $f(\mathbf{x})$ | $\phi(\mathbf{x}) := \sum_{\ \alpha\ \leq s-1} \varphi\left(\frac{\phi_\alpha(\Psi(\mathbf{x}))}{\alpha!}, P_\alpha(\mathbf{x} - \Psi(\mathbf{x}))\right)$ | $\mathcal{O}(N \ln N)$ | $\mathcal{O}(L \ln L)$ | $\mathcal{O}(\ \mathbf{h}\ _2^{-s} + \mathcal{E}_1 + \mathcal{E}_2 + \mathcal{E}_3)$ $\leq \mathcal{O}(K^{-s}) = \mathcal{O}(N^{-2s/d} L^{-2s/d})$ |

771 Finally, we would like to compare our analysis with that in [46]. Both [46] and our
772 analysis rely on local Taylor expansions as in Equation (4.1) to approximate the target
773 function f . Both analysis methods construct ReLU FNNs to approximate polynomials
774 and encode the Taylor expansion coefficients into ReLU FNNs. However, the way to lo-
775 calize the Taylor expansion (i.e., defining the local neighborhood such that the expansion
776 is valid) and the approach to constructing ReLU FNNs are different. We will discuss the
777 details as follows.

778 **Localization.** In [46], a “two-scale” partition procedure and a standard triangulation divide $[0, 1]^d$ into simplexes and a partition of unity is constructed using compactly supported functions that are linear on each simplex, which implies that these functions
779 in the partition of unity can be represented by ReLU FNNs. Taylor expansions of f
780 are constructed within each support of the functions in the partition of unity. In this
781 paper, we simply divide the domain into small hypercubes of uniform size as visualized
782 in Figure 6. Taylor expansions of f are constructed within each hypercube. The reader
783 can understand our approach as a simple way to construct a partition of unity using
784 piecewise constant functions with binary values. The introduction of the trifling region
785 allows us to simply construct ReLU FNNs to approximate these piecewise constant functions
786 without caring about the approximation error within the trifling region. Hence, our
787 construction can be much simplified and makes it easy to estimate all constant prefactors
788 in our error estimates, which is challenging in [46].

791 **ReLU FNNs for Taylor expansions.** In [46], very deep ReLU FNNs with width
792 $\mathcal{O}(1)$ are constructed to approximate polynomials in local Taylor expansions, and hence,
793 the optimal approximation error in width was not explored in [46]. In this paper, we
794 construct ReLU FNNs with arbitrary width and depth to approximate polynomials in
795 local Taylor expansions using Proposition 4.1, which allows us to explore the optimal
796 approximation error in width and is more challenging. In [46], the coefficients of adjacent
797 local Taylor expansions, i.e., $\partial^\alpha f$ in Equation (4.1), are encoded into ReLU FNNs via bit
798 extraction, which is the key to achieving a better approximation error of ReLU FNNs to
799 approximate f than the original local Taylor expansions, since the number of coefficients
800 can be significantly reduced via encoding. Actually, the error in depth by bit extraction
801 is nearly optimal. In this paper, the approximation to $\partial^\alpha f$ is reduced to a point fitting
802 problem that can be solved by constructing ReLU FNNs using bit extraction as sketched
803 out in the previous paragraphs. Hence, we can also achieve the optimal approximation
804 error in depth. The key to achieving the optimal approximation error in width in the
805 above approximation is the application of Lemma 5.4 that essentially fits $\mathcal{O}(N^2)$ samples
806 with ReLU FNNs of width $\mathcal{O}(N)$ and depth 2. Due to the simplicity of our analysis, we
807 can construct ReLU FNNs with arbitrary width and depth to approximate f and specify
808 all constant prefactors in our approximation error.

809 4.2 Constructive proof

810 According to the key ideas of proving Theorem 2.2 summarized in Section 4.1, let
811 us present the detailed proof.

812 *Proof of Theorem 2.2.* The detailed proof can be divided into four steps as follows.

813 **Step 1:** Set up.

814 Set $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor$ and let $\Omega([0, 1]^d, K, \delta)$ partition $[0, 1]^d$ into K^d cubes Q_β for
815 $\beta \in \{0, 1, \dots, K-1\}^d$. In particular, for each $\beta = [\beta_1, \beta_2, \dots, \beta_d]^T \in \{0, 1, \dots, K-1\}^d$, we
816 define $\mathbf{x}_\beta := \beta/K$ and

$$817 \quad Q_\beta := \left\{ \mathbf{x} = [x_1, x_2, \dots, x_d]^T : x_i \in \left[\frac{\beta_i}{K}, \frac{\beta_i+1}{K} - \delta \cdot \mathbf{1}_{\{\beta_i \leq K-2\}} \right] \text{ for } i = 1, 2, \dots, d \right\}.$$

818 Clearly, $[0, 1]^d = \Omega([0, 1]^d, K, \delta) \cup \left(\cup_{\beta \in \{0, 1, \dots, K-1\}^d} Q_\beta \right)$ and \mathbf{x}_β is the vertex of Q_β with
819 minimum $\|\cdot\|_1$ norm. See Figure 6 for the illustrations of Q_β and \mathbf{x}_β .

820 By Proposition 4.3, there exists $\psi \in \mathcal{NN}$ (width $\leq 4N + 3$; depth $\leq 4N + 5$) such that

$$821 \quad \psi(x) = k \quad \text{if } x \in \left[\frac{k}{K}, \frac{k+1}{K} - \delta \cdot \mathbb{1}_{\{k \leq K-2\}} \right] \quad \text{for } k = 0, 1, \dots, K-1.$$

822 Then for each $\beta \in \{0, 1, \dots, K-1\}^d$, $\psi(x_i) = \beta_i$ for all $\mathbf{x} \in Q_\beta$ for $i = 1, 2, \dots, d$.

823 Define

$$824 \quad \Psi(\mathbf{x}) := [\psi(x_1), \psi(x_2), \dots, \psi(x_d)]^T / K \quad \text{for any } \mathbf{x} \in [0, 1]^d,$$

825 then

$$826 \quad \Psi(\mathbf{x}) = \beta / K = \mathbf{x}_\beta \quad \text{if } \mathbf{x} \in Q_\beta \quad \text{for } \beta \in \{0, 1, \dots, K-1\}^d.$$

827 For any $\mathbf{x} \in Q_\beta$ and $\beta \in \{0, 1, \dots, K-1\}^d$, by the Taylor expansion, there exists
828 $\xi_{\mathbf{x}} \in (0, 1)$ such that

$$829 \quad f(\mathbf{x}) = \sum_{\|\alpha\|_1 \leq s-1} \frac{\partial^\alpha f(\Psi(\mathbf{x}))}{\alpha!} \mathbf{h}^\alpha + \sum_{\|\alpha\|_1 = s} \frac{\partial^\alpha f(\Psi(\mathbf{x}) + \xi_{\mathbf{x}} \mathbf{h})}{\alpha!} \mathbf{h}^\alpha, \quad \text{where } \mathbf{h} = \mathbf{x} - \Psi(\mathbf{x}).$$

830 **Step 2:** Construct the desired function ϕ .

831 By Lemma 4.2, there exists

$$832 \quad \varphi \in \mathcal{NN}(\text{width} \leq 9(N+1) + 1; \text{depth} \leq 2s(L+1))$$

833 such that

$$834 \quad |\varphi(x_1, x_2) - x_1 x_2| \leq 216(N+1)^{-2s(L+1)} =: \mathcal{E}_1 \quad \text{for any } x_1, x_2 \in [-3, 3]. \quad (4.3)$$

835 For each $\alpha \in \mathbb{N}^d$ with $\|\alpha\|_1 \leq s$, by Proposition 4.1, there exists

$$836 \quad P_\alpha \in \mathcal{NN}(\text{width} \leq 9(N+1) + s - 1; \text{depth} \leq 7s^2 L)$$

837 such that

$$838 \quad |P_\alpha(\mathbf{x}) - \mathbf{x}^\alpha| \leq 9s(N+1)^{-7sL} =: \mathcal{E}_2 \quad \text{for any } \mathbf{x} \in [0, 1]^d. \quad (4.4)$$

839 For each $i \in \{0, 1, \dots, K^d - 1\}$, define

$$840 \quad \boldsymbol{\eta}(i) = [\eta_1, \eta_2, \dots, \eta_d]^T \in \{0, 1, \dots, K-1\}^d$$

841 such that $\sum_{j=1}^d \eta_j K^{j-1} = i$. Such a map $\boldsymbol{\eta}$ is a bijection from $\{0, 1, \dots, K^d - 1\}$ to $\{0, 1, \dots, K-1\}^d$. For each $\alpha \in \mathbb{N}^d$ with $\|\alpha\|_1 \leq s-1$, define

$$843 \quad \xi_{\alpha, i} = \left(\partial^\alpha f\left(\frac{\boldsymbol{\eta}(i)}{K}\right) + 1 \right) / 2 \quad \text{for } i \in \{0, 1, \dots, K^d - 1\}.$$

844 Then $\|\partial^\alpha f\|_{L^\infty([0,1]^d)} \leq 1$ implies $\xi_{\alpha, i} \in [0, 1]$ for $i = 0, 1, \dots, K^d - 1$ and each α . Note that
845 $K^d = \left(\lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor \right)^d \leq N^2 L^2$. By Proposition 4.4, there exists

$$846 \quad \tilde{\varphi}_\alpha \in \mathcal{NN}(\text{width} \leq 16s(N+1) \log_2(8N); \text{depth} \leq 5(L+2) \log_2(4L))$$

847 such that, for each $\alpha \in \mathbb{N}^d$ with $\|\alpha\|_1 \leq s-1$, we have

$$848 \quad |\tilde{\varphi}_\alpha(i) - \xi_{\alpha, i}| \leq N^{-2s} L^{-2s} \quad \text{for } i = 0, 1, \dots, K^d - 1.$$

849 For each $\alpha \in \mathbb{N}^d$ with $\|\alpha\|_1 \leq s-1$, define

$$850 \quad \phi_\alpha(\mathbf{x}) := 2\tilde{\phi}_\alpha\left(\sum_{j=1}^d x_j K^{j-1}\right) - 1 \quad \text{for any } \mathbf{x} = [x_1, x_2, \dots, x_d]^T \in \mathbb{R}^d.$$

851 It is easy to verify that

$$852 \quad \phi_\alpha \in \mathcal{NN}(\text{width} \leq 16s(N+1)\log_2(8N); \text{depth} \leq 5(L+2)\log_2(4L)).$$

853 Then, for each $\alpha \in \mathbb{N}^d$ with $\|\alpha\|_1 \leq s-1$ and each $\boldsymbol{\eta} = \boldsymbol{\eta}(i) = [\eta_1, \eta_2, \dots, \eta_d]^T \in \{0, 1, \dots, K-1\}^d$ corresponding to $i = \sum_{j=1}^d \eta_j K^{j-1} \in \{0, 1, \dots, K^d - 1\}$, we have

$$855 \quad \begin{aligned} \left| \phi_\alpha\left(\frac{\boldsymbol{\eta}}{K}\right) - \partial^\alpha f\left(\frac{\boldsymbol{\eta}}{K}\right) \right| &= \left| 2\tilde{\phi}_\alpha\left(\sum_{j=1}^d \eta_j K^{j-1}\right) - 1 - (2\xi_{\alpha,i} - 1) \right| \\ &= 2|\tilde{\phi}_\alpha(i) - \xi_{\alpha,i}| \leq 2N^{-2s}L^{-2s}. \end{aligned}$$

856 Therefore, for each $\beta \in \{0, 1, \dots, K-1\}^d$ and each $\alpha \in \mathbb{N}^d$ with $\|\alpha\|_1 \leq s-1$, we have

$$857 \quad \left| \phi_\alpha(\mathbf{x}_\beta) - \partial^\alpha f(\mathbf{x}_\beta) \right| = \left| \phi_\alpha\left(\frac{\beta}{K}\right) - \partial^\alpha f\left(\frac{\beta}{K}\right) \right| \leq 2N^{-2s}L^{-2s} =: \mathcal{E}_3. \quad (4.5)$$

858 Now we can construct the desired function ϕ as

$$859 \quad \phi(\mathbf{x}) := \sum_{\|\alpha\|_1 \leq s-1} \varphi\left(\frac{\phi_\alpha(\Psi(\mathbf{x}))}{\alpha!}, P_\alpha(\mathbf{x} - \Psi(\mathbf{x}))\right) \quad \text{for any } \mathbf{x} \in \mathbb{R}^d. \quad (4.6)$$

860 It remains to estimate the approximation error and determine the size of the network
861 implementing ϕ .

862 **Step 3:** Estimate approximation error.

863 Fix $\beta \in \{0, 1, \dots, K-1\}^d$, let us estimate the approximation error for a fixed $\mathbf{x} \in Q_\beta$.
864 See Table 2 for a summary of the approximation errors. Recall that $\Psi(\mathbf{x}) = \mathbf{x}_\beta$ and
865 $\mathbf{h} = \mathbf{x} - \Psi(\mathbf{x}) = \mathbf{x} - \mathbf{x}_\beta$. It is easy to check that $|f(\mathbf{x}) - \phi(\mathbf{x})|$ is bounded by

$$866 \quad \begin{aligned} &\left| \sum_{\|\alpha\|_1 \leq s-1} \frac{\partial^\alpha f(\Psi(\mathbf{x}))}{\alpha!} \mathbf{h}^\alpha + \sum_{\|\alpha\|_1 = s} \frac{\partial^\alpha f(\Psi(\mathbf{x}) + \xi_{\mathbf{x}} \mathbf{h})}{\alpha!} \mathbf{h}^\alpha - \sum_{\|\alpha\|_1 \leq s-1} \varphi\left(\frac{\phi_\alpha(\Psi(\mathbf{x}))}{\alpha!}, P_\alpha(\mathbf{x} - \Psi(\mathbf{x}))\right) \right| \\ &\leq \underbrace{\sum_{\|\alpha\|_1 = s} \left| \frac{\partial^\alpha f(\mathbf{x}_\beta + \xi_{\mathbf{x}} \mathbf{h})}{\alpha!} \mathbf{h}^\alpha \right|}_{\mathcal{I}_1} + \underbrace{\sum_{\|\alpha\|_1 \leq s-1} \left| \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \mathbf{h}^\alpha - \varphi\left(\frac{\phi_\alpha(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) \right|}_{\mathcal{I}_2} =: \mathcal{I}_1 + \mathcal{I}_2. \end{aligned}$$

867 Recall the fact that

$$868 \quad \sum_{\|\alpha\|_1 = s} 1 = |\{\alpha \in \mathbb{N}^d : \|\alpha\|_1 = s\}| \leq (s+1)^{d-1} \quad \textcircled{7}$$

869 and

$$870 \quad \sum_{\|\alpha\|_1 \leq s-1} 1 = \sum_{i=0}^{s-1} \left(\sum_{\|\alpha\|_1 = i} 1 \right) \leq \sum_{i=0}^{s-1} (i+1)^{d-1} \leq s \cdot (s-1+1)^{d-1} = s^d.$$

$\textcircled{7}$ In fact, we have $|\{\alpha \in \mathbb{N}^d : \|\alpha\|_1 = s\}| = \binom{s+d-1}{d-1}$, implying $(s/d+1)^{d-1} \leq \sum_{\|\alpha\|_1 = s} 1 \leq (s+1)^{d-1}$. Thus, the lower bound of the estimate is still exponentially large in d . To the best of our knowledge, we cannot avoid a constant prefactor that is exponentially large in d when Taylor expansion is used in the analysis.

871 For the first part \mathcal{I}_1 , we have

$$872 \quad \mathcal{I}_1 = \sum_{\|\alpha\|_1=s} \left| \frac{\partial^\alpha f(\mathbf{x}_\beta + \xi \mathbf{h})}{\alpha!} \mathbf{h}^\alpha \right| \leq \sum_{\|\alpha\|_1=s} \left| \frac{1}{\alpha!} \mathbf{h}^\alpha \right| \leq (s+1)^{d-1} K^{-s}.$$

873 For the second part \mathcal{I}_2 , we have

$$874 \quad \mathcal{I}_2 = \sum_{\|\alpha\|_1 \leq s-1} \underbrace{\left| \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \mathbf{h}^\alpha - \varphi\left(\frac{\phi_\alpha(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) \right|}_{\mathcal{I}_2(\alpha)} =: \sum_{\|\alpha\|_1 \leq s-1} \mathcal{I}_2(\alpha).$$

875 Fix $\alpha \in \mathbb{N}^d$ with $\|\alpha\|_1 \leq s-1$, we have

$$\begin{aligned} \mathcal{I}_2(\alpha) &= \left| \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \mathbf{h}^\alpha - \varphi\left(\frac{\phi_\alpha(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) \right| \\ 876 \quad &\leq \underbrace{\left| \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \mathbf{h}^\alpha - \varphi\left(\frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) \right|}_{\mathcal{I}_{2,1}(\alpha)} + \underbrace{\left| \varphi\left(\frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) - \varphi\left(\frac{\phi_\alpha(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) \right|}_{\mathcal{I}_{2,2}(\alpha)} \\ &=: \mathcal{I}_{2,1}(\alpha) + \mathcal{I}_{2,2}(\alpha). \end{aligned}$$

877 Note that $\mathcal{E}_2 = 9s(N+1)^{-7sL} \leq 9s(2)^{-7s} \leq 2$. By $\mathbf{h}^\alpha \in [0, 1]$ and Equation (4.4), we
878 have $P_\alpha(\mathbf{h}) \in [-2, 3] \subseteq [-3, 3]$. Then by $\partial^\alpha f(\mathbf{x}_\beta) \in [-1, 1]$ and Equations (4.3) and (4.4),
879 we have

$$\begin{aligned} \mathcal{I}_{2,1}(\alpha) &= \left| \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \mathbf{h}^\alpha - \varphi\left(\frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) \right| \\ &\leq \left| \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \mathbf{h}^\alpha - \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} P_\alpha(\mathbf{h}) \right| + \underbrace{\left| \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} P_\alpha(\mathbf{h}) - \varphi\left(\frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) \right|}_{\leq \mathcal{E}_1 \text{ by Eq. (4.3)}} \\ 880 \quad &\leq \frac{1}{\alpha!} \underbrace{\left| \mathbf{h}^\alpha - P_\alpha(\mathbf{h}) \right|}_{\leq \mathcal{E}_2 \text{ by Eq. (4.4)}} + \mathcal{E}_1 \leq \frac{1}{\alpha!} \mathcal{E}_2 + \mathcal{E}_1 \leq \mathcal{E}_1 + \mathcal{E}_2. \end{aligned}$$

881 To estimate $\mathcal{I}_{2,2}(\alpha)$, we need the following fact derived from Equation (4.3):

$$\begin{aligned} |\varphi(x_1, x_2) - \varphi(\tilde{x}_1, x_2)| &\leq \underbrace{|\varphi(x_1, x_2) - x_1 x_2|}_{\leq \mathcal{E}_1 \text{ by Eq. (4.3)}} + \underbrace{|\varphi(\tilde{x}_1, x_2) - \tilde{x}_1 x_2|}_{\leq \mathcal{E}_1 \text{ by Eq. (4.3)}} + |x_1 x_2 - \tilde{x}_1 x_2| \\ 882 \quad &\leq 2\mathcal{E}_1 + 3|x_1 - \tilde{x}_1|, \end{aligned} \tag{4.7}$$

883 for any $x_1, \tilde{x}_1, x_2 \in [-3, 3]$.

884 Since $\mathcal{E}_3 = 2N^{-2s}L^{-2s} \leq 2$ and $\partial^\alpha f(\mathbf{x}_\beta) \in [-1, 1]$, we have $\phi_\alpha(\mathbf{x}_\beta) \in [-3, 3]$ by
885 Equation (4.5). Then by $P_\alpha(\mathbf{h}) \in [-3, 3]$ and Equations (4.7) and (4.5), we have

$$\begin{aligned} \mathcal{I}_{2,2}(\alpha) &= \left| \varphi\left(\frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) - \varphi\left(\frac{\phi_\alpha(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) \right| \\ 886 \quad &\leq 2\mathcal{E}_1 + 3 \underbrace{\left| \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} - \frac{\phi_\alpha(\mathbf{x}_\beta)}{\alpha!} \right|}_{\leq \mathcal{E}_3 \text{ by Eq. (4.5)}} \leq 2\mathcal{E}_1 + 3\mathcal{E}_3. \end{aligned}$$

887 Therefore, we get

$$\begin{aligned}
|f(\mathbf{x}) - \phi(\mathbf{x})| &\leq \mathcal{I}_1 + \mathcal{I}_2 \leq \mathcal{I}_1 + \sum_{\|\alpha\|_1 \leq s-1} \mathcal{I}_2(\alpha) \leq \mathcal{I}_1 + \sum_{\|\alpha\|_1 \leq s-1} \left(\mathcal{I}_{2,1}(\alpha) + \mathcal{I}_{2,2}(\alpha) \right) \\
&\leq (s+1)^{d-1} K^{-s} + s^d \left((\mathcal{E}_1 + \mathcal{E}_2) + (2\mathcal{E}_1 + 3\mathcal{E}_3) \right) \\
&\leq (s+1)^d (K^{-s} + 3\mathcal{E}_1 + \mathcal{E}_2 + 3\mathcal{E}_3).
\end{aligned}$$

889 Since $\beta \in \{0, 1, \dots, K-1\}^d$ and $\mathbf{x} \in Q_\beta$ are arbitrary and

$$890 \quad [0, 1]^d = \Omega([0, 1]^d, K, \delta) \cup \left(\cup_{\beta \in \{0, 1, \dots, K-1\}^d} Q_\beta \right),$$

891 we have, for any $\mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)$,

$$892 \quad |f(\mathbf{x}) - \phi(\mathbf{x})| \leq (s+1)^d (K^{-s} + 3\mathcal{E}_1 + \mathcal{E}_2 + 3\mathcal{E}_3).$$

893 Recall that $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor \geq \frac{N^{2/d} L^{2/d}}{8}$ and

$$894 \quad (N+1)^{-7sL} \leq (N+1)^{-2s(L+1)} \leq (N+1)^{-2s} 2^{-2sL} \leq N^{-2s} L^{-2s}.$$

895 Then we have

$$\begin{aligned}
&(s+1)^d (K^{-s} + 3\mathcal{E}_1 + \mathcal{E}_2 + 3\mathcal{E}_3) \\
&= (s+1)^d \left(K^{-s} + 648(N+1)^{-2s(L+1)} + 9s(N+1)^{-7sL} + 6N^{-2s} L^{-2s} \right) \\
&\leq (s+1)^d \left(8^s N^{-2s/d} L^{-2s/d} + (654 + 9s) N^{-2s} L^{-2s} \right) \\
&\leq (s+1)^d (8^s + 654 + 9s) N^{-2s/d} L^{-2s/d} \leq 84(s+1)^d 8^s N^{-2s/d} L^{-2s/d}.
\end{aligned}$$

897 **Step 4:** Determine the size of the network implementing ϕ .

898 It remains to estimate the width and depth of the network implementing ϕ . Recall
899 that, for $\alpha \in \mathbb{N}^d$ with $\|\alpha\|_1 \leq s-1$,

$$900 \quad \begin{cases} \Psi \in \mathcal{NN}(\text{width} \leq d(4N+3); \text{depth} \leq 4L+5), \\ \phi_\alpha \in \mathcal{NN}(\text{width} \leq 16s(N+1) \log_2(8N); \text{depth} \leq 5(L+2) \log_2(4L)), \\ P_\alpha \in \mathcal{NN}(\text{width} \leq 9(N+1) + s-1; \text{depth} \leq 7s^2L), \\ \varphi \in \mathcal{NN}(\text{width} \leq 9(N+1) + 1; \text{depth} \leq 2s(L+1)). \end{cases}$$

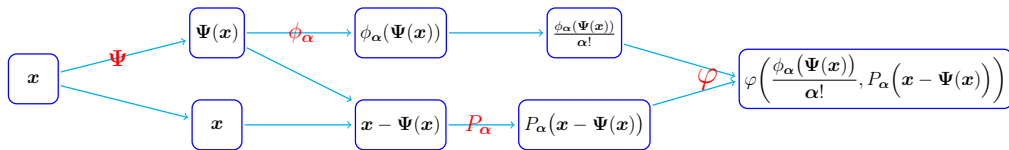


Figure 7: An illustration of the sub-network architecture implementing each component of ϕ , $\varphi\left(\frac{\phi_\alpha(\Psi(\mathbf{x}))}{\alpha!}, P_\alpha(\mathbf{x} - \Psi(\mathbf{x}))\right)$ for each $\alpha \in \mathbb{N}^d$ with $\|\alpha\| \leq s-1$.

901 By Equation (4.6) and Figure 7, it easy to verify that ϕ can be implemented by a
 902 ReLU FNN with width

$$903 \quad \sum_{\|\alpha\|_1 \leq s-1} 16sd(N+2)\log_2(8N) \leq s^d \cdot 16sd(N+2)\log_2(8N)$$

$$= 16s^{d+1}d(N+2)\log_2(8N)$$

904 and depth

$$905 \quad (4L+5) + 2s(L+1) + 7s^2L + 5(L+2)\log_2(4L) + 3 \leq 18s^2(L+2)\log_2(4L)$$

906 as desired. So we finish the proof. □

907 5 Proofs of Propositions in Section 4.1

908 In this section, we will prove all propositions in Section 4.1.

909 5.1 Proof of Proposition 4.1 for polynomial approximation

910 To prove Proposition 4.1, we will construct ReLU FNNs to approximate multivariate
 911 polynomials following the four steps below.

- 912 • $f(x) = x^2$. We approximate $f(x) = x^2$ by the combinations and compositions of
 913 “sawtooth” functions as shown in Figures 8 and 9.
- 914 • $f(x, y) = xy$. To approximate $f(x, y) = xy$, we use the result of the previous step
 915 and the fact that $xy = 2\left(\left(\frac{x+y}{2}\right)^2 - \left(\frac{x}{2}\right)^2 - \left(\frac{y}{2}\right)^2\right)$.
- 916 • $f(x_1, x_2, \dots, x_k) = x_1x_2 \cdots x_k$. We approximate $f(x_1, x_2, \dots, x_k) = x_1x_2 \cdots x_k$ for any
 917 $k \geq 2$ via mathematical induction based on the result of the previous step.
- 918 • A general polynomial $P(\mathbf{x}) = \mathbf{x}^\alpha = x_1^{\alpha_1}x_2^{\alpha_2} \cdots x_d^{\alpha_d}$ with $\|\alpha\|_1 \leq k$. Any one-term
 919 polynomial of degree $\leq k$ can be written as $Cz_1z_2 \cdots z_k$ with some entries equaling
 920 1, where C is a constant and $\mathbf{z} = [z_1, z_2, \dots, z_k]^T$ can be attained via an affine linear
 921 map with \mathbf{x} as the input. Then use the result of the previous step.

922 The idea of using “sawtooth” functions (see Figure 8) was first raised in [44] for
 923 approximating x^2 using FNNs with width 6 and depth $\mathcal{O}(L)$ and achieving an error
 924 $\mathcal{O}(2^{-L})$; our construction is different from and more general than that in [44], working
 925 for ReLU FNNs of width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ for any N and L , and achieving an
 926 error $\mathcal{O}(N^{-L})$. As discussed below Proposition 4.1, this $\mathcal{O}(N^{-L})$ approximation error of
 927 polynomial functions shows the power of depth in ReLU FNNs via function composition.

928 First, let us show how to construct ReLU FNNs to approximate $f(x) = x^2$.

929 **Lemma 5.1.** *For any $N, L \in \mathbb{N}^+$, there exists a function ϕ implemented by a ReLU FNN
 930 with width $3N$ and depth L such that*

$$931 \quad |\phi(x) - x^2| \leq N^{-L} \quad \text{for any } x \in [0, 1].$$

932 *Proof.* Define a set of “sawtooth” functions $T_i : [0, 1] \rightarrow [0, 1]$ by induction as follows.
 933 Set

934
$$T_1(x) = \begin{cases} 2x, & \text{if } x \in [0, \frac{1}{2}], \\ 2(1-x), & \text{if } x \in (\frac{1}{2}, 1], \end{cases}$$

935 and

936
$$T_i = T_{i-1} \circ T_1 \quad \text{for } i = 2, 3, \dots.$$

937 It is easy to check that T_i has 2^{i-1} “sawteeth” and

938
$$T_{m+n} = T_m \circ T_n \quad \text{for any } m, n \in \mathbb{N}^+.$$

See Figure 8 for illustrations of T_i for $i = 1, 2, 3, 4$.

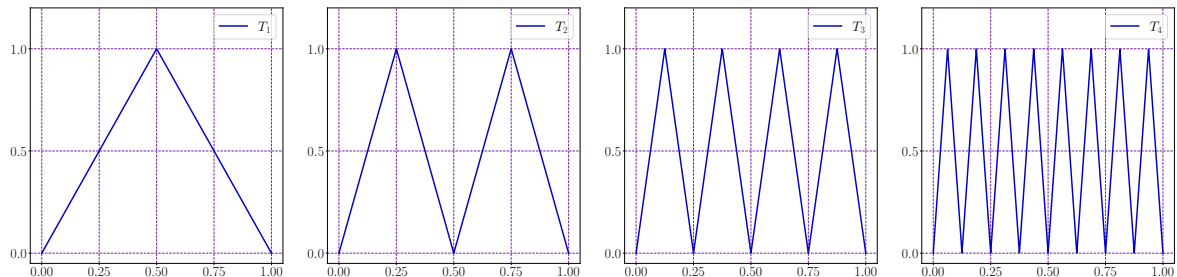


Figure 8: Examples of “sawtooth” functions T_1 , T_2 , T_3 , and T_4 .

939

940 Define piecewise linear functions $f_s : [0, 1] \rightarrow [0, 1]$ for $s \in \mathbb{N}^+$ satisfying the following
 941 two requirements (see Figure 9 for several examples of f_s).

- 942 • $f_s(\frac{j}{2^s}) = (\frac{j}{2^s})^2$ for $j = 0, 1, 2, \dots, 2^s$.
- 943 • $f_s(x)$ is linear between any two adjacent points of $\{\frac{j}{2^s} : j = 0, 1, 2, \dots, 2^s\}$.

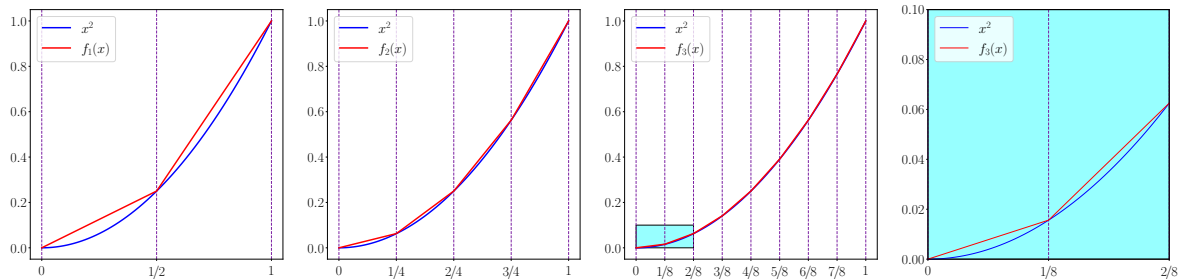


Figure 9: Illustrations of f_1 , f_2 , and f_3 for approximating x^2 .

944 Recall the fact

945
$$0 \leq tx_1^2 + (1-t)x_2^2 - (tx_1 + (1-t)x_2)^2 \leq \frac{(x_2 - x_1)^2}{4} \quad \text{for any } t, x_1, x_2 \in [0, 1].$$

946 Thus, we have

947
$$0 \leq f_s(x) - x^2 \leq \frac{(2^{-s})^2}{4} = 2^{-2(s+1)} \quad \text{for any } x \in [0, 1] \text{ and } s \in \mathbb{N}^+. \quad (5.1)$$

948 Note that $f_{i-1}(x) = f_i(x) = x^2$ for $x \in \{\frac{j}{2^{i-1}} : j = 0, 1, 2, \dots, 2^{i-1}\}$ and the graph of $f_{i-1} - f_i$
 949 is a symmetric “sawtooth” between any two adjacent points of $\{\frac{j}{2^{i-1}} : j = 0, 1, 2, \dots, 2^{i-1}\}$.
 950 It is easy to verify that

$$951 \quad f_{i-1}(x) - f_i(x) = \frac{T_i(x)}{2^{2i}} \quad \text{for any } x \in [0, 1] \text{ and } i = 2, 3, \dots.$$

952 Therefore, for any $x \in [0, 1]$ and $s \in \mathbb{N}^+$, we have

$$953 \quad f_s(x) = f_1(x) + \sum_{i=2}^s (f_i - f_{i-1}) = x - (x - f_1(x)) - \sum_{i=2}^s \frac{T_i(x)}{2^{2i}} = x - \sum_{i=2}^s \frac{T_i(x)}{2^{2i}}.$$

954 Given $N \in \mathbb{N}^+$, there exists a unique $k \in \mathbb{N}^+$ such that $(k-1)2^{k-1} + 1 \leq N \leq k2^k$.
 955 For this k , using $s = Lk$, we can construct a ReLU FNN as shown in Figure 10 to
 956 implement a function $\phi = f_{Lk}$ approximating x^2 well. Note that T_i can be implemented
 957 by a one-hidden-layer ReLU FNN with width 2^i . Hence, the network in Figure 10 has
 958 width $k2^k + 1 \leq 3N$ [Ⓢ] and depth $2L$.

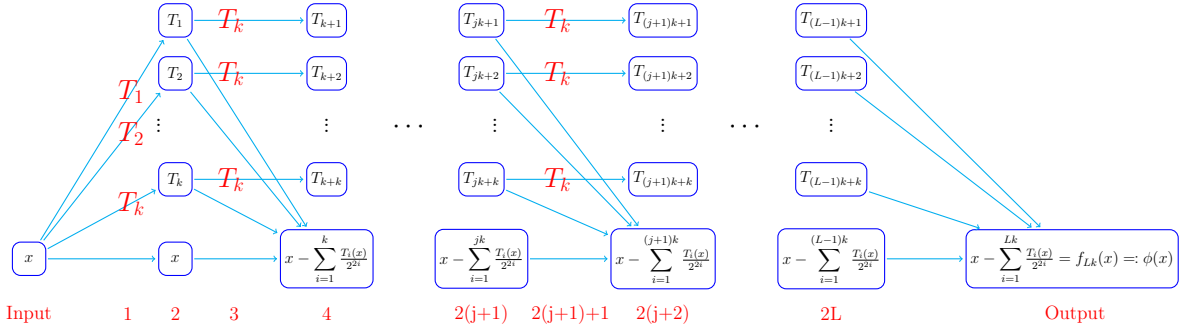


Figure 10: An illustration of the target network architecture for approximating x^2 on $[0, 1]$. T_i can be implemented by a one-hidden-layer ReLU FNN with width 2^i for $i = 1, 2, \dots, K$. The red numbers below the architecture indicate the order of hidden layers.

959 As shown in Figure 10, the (2ℓ) -th hidden layer of the network has the identify
 960 function as activation functions for $\ell = 1, 2, \dots, L$. Thus, the network in Figure 10 can
 961 be interpreted as a ReLU FNN with width $3N$ and depth L . In fact, if all activation
 962 functions in a certain hidden layer are identity maps, the depth can be reduced by one via
 963 combining two adjacent linear transforms into one. For example, suppose $\mathbf{W}_1 \in \mathbb{R}^{N_1 \times N_2}$,
 964 $\mathbf{W}_2 \in \mathbb{R}^{N_2 \times N_3}$, and ϱ is an identity map that can be applied to vectors or matrices
 965 elementwisely; then $\mathbf{W}_1 \varrho(\mathbf{W}_2 \mathbf{x}) = \mathbf{W}_3 \mathbf{x}$ for any $\mathbf{x} \in \mathbb{R}^{N_3}$, where $\mathbf{W}_3 = \mathbf{W}_1 \cdot \mathbf{W}_2 \in \mathbb{R}^{N_1 \times N_3}$.

966 It remains to estimate the approximation error of $\phi(x) \approx x^2$. By Equation (5.1), for
 967 any $x \in [0, 1]$, we have

$$968 \quad |\phi(x) - x^2| = |f_{Lk}(x) - x^2| \leq 2^{-2(Lk+1)} \leq 2^{-2Lk} \leq N^{-L},$$

969 where the last inequality comes from $N \leq k2^k \leq 2^{2k}$. So we finish the proof. \square

[Ⓢ]This inequality is clear for $k = 1, 2, 3, 4$. In the case $k \geq 5$, we have $k2^k + 1 \leq \frac{k2^k + 1}{N} N \leq \frac{(k+1)2^k}{(k-1)2^{k-1}} N \leq 2 \frac{k+1}{k-1} N \leq 3N$.

970 We have constructed a ReLU FNN to approximate $f(x) = x^2$. By the fact that
 971 $xy = 2\left(\left(\frac{x+y}{2}\right)^2 - \left(\frac{x}{2}\right)^2 - \left(\frac{y}{2}\right)^2\right)$, it is easy to construct a new ReLU FNN to approximate
 972 $f(x, y) = xy$ as follows.

973 **Lemma 5.2.** *For any $N, L \in \mathbb{N}^+$, there exists a function ϕ implemented by a ReLU FNN
 974 with width $9N$ and depth L such that*

$$975 \quad |\phi(x, y) - xy| \leq 6N^{-L} \quad \text{for any } x, y \in [0, 1].$$

976 *Proof.* By Lemma 5.1, there exists a function ψ implemented by a ReLU FNN with
 977 width $3N$ and depth L such that

$$978 \quad |x^2 - \psi(x)| \leq N^{-L} \quad \text{for any } x \in [0, 1].$$

979 Inspired by the fact

$$980 \quad xy = 2\left(\left(\frac{x+y}{2}\right)^2 - \left(\frac{x}{2}\right)^2 - \left(\frac{y}{2}\right)^2\right) \quad \text{for any } x, y \in \mathbb{R},$$

981 we construct the desired function ϕ as

$$982 \quad \phi(x, y) := 2\left(\psi\left(\frac{x+y}{2}\right) - \psi\left(\frac{x}{2}\right) - \psi\left(\frac{y}{2}\right)\right) \quad \text{for any } x, y \in \mathbb{R}. \quad (5.2)$$

983 Then ϕ can be implemented by the network architecture in Figure 11.

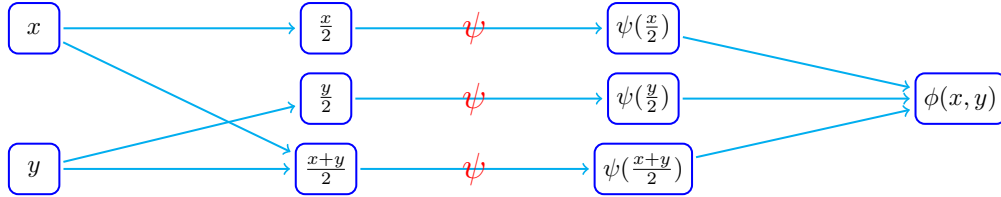


Figure 11: An illustration of the network architecture implementing ϕ for approximating xy on $[0, 1]^2$.

984 It follows from $\psi \in \mathcal{NN}$ (width $\leq 3N$; depth $\leq L$) that the network in Figure 11 is
 985 with width $9N$ and depth $L + 2$. Similar to the discussion in the proof of Lemma 5.1,
 986 the network in Figure 11 can be interpreted as a ReLU FNN with width $9N$ and depth
 987 L , since two of the hidden layers have the identity function as their activation functions.
 988 Moreover, for any $x, y \in [0, 1]$,

$$989 \quad \begin{aligned} |xy - \phi(x, y)| &= \left| 2\left(\left(\frac{x+y}{2}\right)^2 - \left(\frac{x}{2}\right)^2 - \left(\frac{y}{2}\right)^2\right) - 2\left(\psi\left(\frac{x+y}{2}\right) - \psi\left(\frac{x}{2}\right) - \psi\left(\frac{y}{2}\right)\right) \right| \\ &\leq 2\left|\left(\frac{x+y}{2}\right)^2 - \psi\left(\frac{x+y}{2}\right)\right| + 2\left|\left(\frac{x}{2}\right)^2 - \psi\left(\frac{x}{2}\right)\right| + 2\left|\left(\frac{y}{2}\right)^2 - \psi\left(\frac{y}{2}\right)\right| \leq 6N^{-L}. \end{aligned}$$

990 Therefore, we have finished the proof. \square

991 Now let us prove Lemma 4.2, which shows how to construct a ReLU FNN to approx-
 992 imate $f(x, y) = xy$ on $[a, b]^2$ with arbitrary $a < b$, i.e., a rescaled version of Lemma 5.2.

993 *Proof of Lemma 4.2.* By Lemma 5.2, there exists a function ψ implemented by a ReLU
 994 FNN with width $9N$ and depth L such that

$$995 \quad |\psi(\tilde{x}, \tilde{y}) - \tilde{x}\tilde{y}| \leq 6N^{-L} \quad \text{for any } \tilde{x}, \tilde{y} \in [0, 1].$$

996 By setting $\tilde{x} = \frac{x-a}{b-a}$ and $\tilde{y} = \frac{y-a}{b-a}$ for any $x, y \in [a, b]$, we have $\tilde{x}, \tilde{y} \in [0, 1]$, implying

$$997 \quad \left| \psi\left(\frac{x-a}{b-a}, \frac{y-a}{b-a}\right) - \frac{x-a}{b-a} \frac{y-a}{b-a} \right| \leq 6N^{-L} \quad \text{for any } x, y \in [a, b].$$

998 It follows that, for any $x, y \in [a, b]$,

$$999 \quad \left| (b-a)^2 \psi\left(\frac{x-a}{b-a}, \frac{y-a}{b-a}\right) + a(x+y) - a^2 - xy \right| \leq 6(b-a)^2 N^{-L}.$$

1000 Define, for any $x, y \in \mathbb{R}$,

$$1001 \quad \phi(x, y) := (b-a)^2 \psi\left(\frac{x-a}{b-a}, \frac{y-a}{b-a}\right) + a \cdot \sigma(x+y+2|a|) - a^2 - 2a|a|.$$

1002 Then ϕ can be implemented by the network architecture in Figure 12.

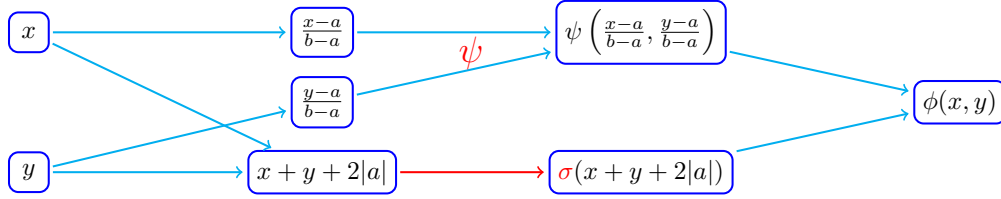


Figure 12: An illustration of the network architecture implementing ϕ for approximating xy on $[a, b]^2$. Two of the hidden layers have the identify function as their activation functions, since the red “ σ ” comes from the red arrow “ \rightarrow ”, where the red arrow “ \rightarrow ” is a ReLU FNN with width 1 and depth L .

1003 It follows from $\psi \in \mathcal{NN}$ (width $\leq 9N$; depth $\leq L$) that the network in Figure 12 is
 1004 with width $9N+1$ and depth $L+2$. Similar to the discussion in the proof of Lemma 5.1,
 1005 the network in Figure 12 can be interpreted as a ReLU FNN with width $9N+1$ and
 1006 depth L , since two of the hidden layers have the identify function as their activation
 1007 functions.

1008 Note that $x+y+2|a| \geq 0$ for any $x, y \in [a, b]$, implying

$$1009 \quad \phi(x, y) = (b-a)^2 \psi\left(\frac{x-a}{b-a}, \frac{y-a}{b-a}\right) + a(x+y) - a^2 \quad \text{for any } x, y \in [a, b].$$

1010 Hence,

$$1011 \quad |\phi(x, y) - xy| \leq 6(b-a)^2 N^{-L} \quad \text{for any } x, y \in [a, b].$$

1012 So we finish the proof. □

1013 The next lemma shows how to construct a ReLU FNN to approximate a multivariate
 1014 function $f(x_1, x_2, \dots, x_k) = x_1 x_2 \dots x_k$ on $[0, 1]^k$.

1015 **Lemma 5.3.** *For any $N, L, k \in \mathbb{N}^+$ with $k \geq 2$, there exists a function ϕ implemented by*
 1016 *a ReLU FNN with width $9(N+1) + k - 1$ and depth $7kL(k-1)$ such that*

$$1017 \quad |\phi(\mathbf{x}) - x_1 x_2 \dots x_k| \leq 9(k-1)(N+1)^{-7kL} \quad \text{for any } \mathbf{x} = [x_1, x_2, \dots, x_k]^T \in [0, 1]^k.$$

1018 *Proof.* By Lemma 4.2, there exists a function ϕ_1 implemented by a ReLU FNN with
 1019 width $9(N+1)+1$ and depth $7kL$ such that

$$1020 \quad |\phi_1(x, y) - xy| \leq 6(1.2)^2(N+1)^{-7kL} \leq 9(N+1)^{-7kL} \quad \text{for any } x, y \in [-0.1, 1.1]. \quad (5.3)$$

1021 Next, we construct a sequence of functions $\phi_i : [0, 1]^{i+1} \rightarrow [0, 1]$ for $i \in \{1, 2, \dots, k-1\}$ by
 1022 induction such that

1023 (i) ϕ_i can be implemented by a ReLU FNN with width $9(N+1)+i$ and depth $7kLi$
 1024 for each $i \in \{1, 2, \dots, k-1\}$.

1025 (ii) For any $i \in \{1, 2, \dots, k-1\}$ and $x_1, x_2, \dots, x_{i+1} \in [0, 1]$, it holds that

$$1026 \quad |\phi_i(x_1, \dots, x_{i+1}) - x_1 x_2 \cdots x_{i+1}| \leq 9i(N+1)^{-7kL}. \quad (5.4)$$

1027 First, let us consider the case $i = 1$, it is obvious that the two required conditions
 1028 are true: 1) $9(N+1)+i = 9(N+1)+1$ and $7kLi = 7kL$ if $i = 1$; 2) Equation (5.3) implies
 1029 Equation (5.4) for $i = 1$.

1030 Now assume ϕ_i has been defined; we then define

$$1031 \quad \phi_{i+1}(x_1, \dots, x_{i+2}) := \phi_1(\phi_i(x_1, \dots, x_{i+1}), \sigma(x_{i+2})) \quad \text{for any } x_1, \dots, x_{i+2} \in \mathbb{R}.$$

1032 Note that $\phi_i \in \mathcal{NN}$ (width $\leq 9(N+1)+i$; depth $\leq 7kLi$) and $\phi_1 \in \mathcal{NN}$ (width $\leq 9(N+1)+$
 1033 1 ; depth $\leq 7kL$). Then ϕ_{i+1} can be implemented via a ReLU FNN with width

$$1034 \quad \max\{9(N+1)+i+1, 9(N+1)+1\} = 9(N+1)+(i+1)$$

1035 and depth $7kLi+7kL = 7kL(i+1)$.

1036 By the hypothesis of induction, we have

$$1037 \quad |\phi_i(x_1, \dots, x_{i+1}) - x_1 x_2 \cdots x_{i+1}| \leq 9i(N+1)^{-7kL}. \quad (5.5)$$

1038 Recall the fact that $9i(N+1)^{-7kL} \leq 9k2^{-7k} \leq 9k\frac{2^{-7}}{k} \leq 0.1$ for any $N, L, k \in \mathbb{N}^+$ and
 1039 $i \in \{1, 2, \dots, k-1\}$. It follows that

$$1040 \quad \phi_i(x_1, \dots, x_{i+1}) \in [-0.1, 1.1] \quad \text{for any } x_1, \dots, x_{i+1} \in [0, 1].$$

1041 Therefore, by Equations (5.3) and (5.5), we have

$$\begin{aligned} & |\phi_{i+1}(x_1, \dots, x_{i+2}) - x_1 x_2 \cdots x_{i+2}| \\ &= |\phi_1(\phi_i(x_1, \dots, x_{i+1}), \sigma(x_{i+2})) - x_1 x_2 \cdots x_{i+2}| \\ 1042 &\leq |\phi_1(\phi_i(x_1, \dots, x_{i+1}), x_{i+2}) - \phi_i(x_1, \dots, x_{i+1})x_{i+2}| + |\phi_i(x_1, \dots, x_{i+1})x_{i+2} - x_1 x_2 \cdots x_{i+2}| \\ &\leq 9(N+1)^{-7kL} + 9i(N+1)^{-7kL} = 9(i+1)(N+1)^{-7kL}, \end{aligned}$$

1043 for any $x_1, x_2, \dots, x_{i+2} \in [0, 1]$, which means we finish the process of induction.

1044 Now let $\phi := \phi_{k-1}$, by the principle of induction, we have

$$1045 \quad |\phi(x_1, \dots, x_k) - x_1 x_2 \cdots x_k| \leq 9(k-1)(N+1)^{-7kL} \quad \text{for any } x_1, \dots, x_k \in [0, 1].$$

1046 So ϕ is the desired function implemented by a ReLU FNN with width $9(N+1)+k-1$
 1047 and depth $7kL(k-1)$, which means we finish the proof. \square

1048 With Lemma 5.3 in hand, we are ready to prove Proposition 4.1 for approximating
 1049 general multivariate polynomials by ReLU FNNs.

1050 *Proof of Proposition 4.1.* The case $k = 1$ is trivial, so we assume $k \geq 2$ below. Set
 1051 $\tilde{k} = \|\boldsymbol{\alpha}\|_1 \leq k$, denote $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_d]^T$, and let $[z_1, z_2, \dots, z_{\tilde{k}}]^T \in \mathbb{R}^{\tilde{k}}$ be the vector such
 1052 that

$$1053 \quad z_\ell = x_j \quad \text{if} \quad \sum_{i=1}^{j-1} \alpha_i < \ell \leq \sum_{i=1}^j \alpha_i \quad \text{for } j = 1, 2, \dots, d.$$

1054 That is,

$$1055 \quad [z_1, z_2, \dots, z_{\tilde{k}}]^T = \left[\overbrace{x_1, \dots, x_1}^{\alpha_1 \text{ times}}, \overbrace{x_2, \dots, x_2}^{\alpha_2 \text{ times}}, \dots, \overbrace{x_d, \dots, x_d}^{\alpha_d \text{ times}} \right]^T \in \mathbb{R}^{\tilde{k}}.$$

1056 Then we have $P(\mathbf{x}) = \mathbf{x}^\alpha = z_1 z_2 \dots z_{\tilde{k}}$.

1057 We construct the target ReLU FNN in two steps. First, there exists an affine linear
 1058 map $\mathcal{L} : \mathbb{R}^d \rightarrow \mathbb{R}^k$ that duplicates \mathbf{x} to form a new vector $[z_1, z_2, \dots, z_{\tilde{k}}, 1, \dots, 1]^T \in \mathbb{R}^k$,
 1059 i.e., $\mathcal{L}(\mathbf{x}) = [z_1, z_2, \dots, z_{\tilde{k}}, 1, \dots, 1]^T \in \mathbb{R}^k$. Second, by Lemma 5.3, there exists a function
 1060 $\psi : \mathbb{R}^k \rightarrow \mathbb{R}$ implemented by a ReLU FNN with width $9(N+1) + k - 1$ and depth
 1061 $7kL(k-1)$ such that ψ maps $[z_1, z_2, \dots, z_{\tilde{k}}, 1, \dots, 1]^T \in \mathbb{R}^k$ to $z_1 z_2 \dots z_{\tilde{k}}$ within an error
 1062 $9(k-1)(N+1)^{-7kL}$. Hence, we can construct the desired function via $\phi := \psi \circ \mathcal{L}$. Then ϕ can
 1063 be implemented by a ReLU FNN with width $9(N+1) + k - 1$ and depth $7kL(k-1) \leq 7k^2L$,
 1064 and

$$\begin{aligned} |\phi(\mathbf{x}) - P(\mathbf{x})| &= |\phi(\mathbf{x}) - \mathbf{x}^\alpha| = |\psi \circ \mathcal{L}(\mathbf{x}) - x_1^{\alpha_1} x_2^{\alpha_2} \dots x_d^{\alpha_d}| \\ 1065 \quad &= |\psi(z_1, z_2, \dots, z_{\tilde{k}}, 1, \dots, 1) - z_1 z_2 \dots z_{\tilde{k}}| \\ &\leq 9(k-1)(N+1)^{-7kL} \leq 9k(N+1)^{-7kL} \end{aligned}$$

1066 for any $x_1, x_2, \dots, x_d \in [0, 1]$. So, we finish the proof. \square

1067 5.2 Proof of Proposition 4.3 for step function approximation

1068 To prove Proposition 4.3 in this sub-section, we will discuss how to pointwisely
 1069 approximate step functions by ReLU FNNs except for the trifling region. Before proving
 1070 Proposition 4.3, let us first introduce a basic lemma about fitting $\mathcal{O}(N_1 N_2)$ samples
 1071 using a two-hidden-layer ReLU FNN with $\mathcal{O}(N_1 + N_2)$ neurons.

1072 **Lemma 5.4.** *For any $N_1, N_2 \in \mathbb{N}^+$, given $N_1(N_2 + 1) + 1$ samples $(x_i, y_i) \in \mathbb{R}^2$ with
 1073 $x_0 < x_1 < \dots < x_{N_1(N_2+1)}$ and $y_i \geq 0$ for $i = 0, 1, \dots, N_1(N_2+1)$, there exists $\phi \in \mathcal{NN}(\#input =$
 1074 $1; \text{widthvec} = [2N_1, 2N_2 + 1])$ satisfying the following conditions:*

- 1075 1. $\phi(x_i) = y_i$ for $i = 0, 1, \dots, N_1(N_2 + 1)$.
- 1076 2. ϕ is linear on each interval $[x_{i-1}, x_i]$ for $i \notin \{(N_2 + 1)j : j = 1, 2, \dots, N_1\}$.

1077 The above lemma is Lemma 2.2 of [40]; and the reader is referred to [40] for its
 1078 proof. Essentially, this lemma shows the equivalence of one-hidden-layer ReLU FNNs of
 1079 size $\mathcal{O}(N^2)$ and two-hidden-layer ones of size $\mathcal{O}(N)$ to fit $\mathcal{O}(N^2)$ samples.

1080 The next lemma below shows that special shallow and wide ReLU FNNs can be
 1081 represented by deep and narrow ones. This lemma was proposed as Proposition 2.2
 1082 in [41].

1083 **Lemma 5.5.** For any $N, L, d \in \mathbb{N}^+$, it holds that

$$1084 \quad \begin{aligned} & \mathcal{NN}(\#\text{input} = d; \text{widthvec} = [N, NL]; \#\text{output} = 1) \\ & \subseteq \mathcal{NN}(\#\text{input} = d; \text{width} \leq 2N + 2; \text{depth} \leq L + 1; \#\text{output} = 1). \end{aligned}$$

1085 With Lemmas 5.4 and 5.5 in hand, let us present the detailed proof of Proposi-
1086 tion 4.3.

1087 *Proof of Proposition 4.3.* We divide the proof into two cases: $d = 1$ and $d \geq 2$.

1088 **Case 1:** $d = 1$.

1089 In this case, $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor = N^2 L^2$. Denote $M = N^2 L$ and consider the sample
1090 set

$$1091 \quad \begin{aligned} & \{(1, M - 1), (2, 0)\} \cup \left\{ \left(\frac{m}{M}, m \right) : m = 0, 1, \dots, M - 1 \right\} \\ & \cup \left\{ \left(\frac{m+1}{M} - \delta, m \right) : m = 0, 1, \dots, M - 2 \right\}. \end{aligned}$$

1092 Its size is $2M + 1 = N \cdot ((2NL - 1) + 1) + 1$. By Lemma 5.4 (set $N_1 = N$ and $N_2 = 2NL - 1$
1093 therein), there exists

$$1094 \quad \begin{aligned} & \phi_1 \in \mathcal{NN}(\text{widthvec} = [2N, 2(2NL - 1) + 1]) \\ & = \mathcal{NN}(\text{widthvec} = [2N, 4NL - 1]) \end{aligned}$$

1095 such that

- 1096 • $\phi_1\left(\frac{M-1}{M}\right) = \phi_1(1) = M - 1$ and $\phi_1\left(\frac{m}{M}\right) = \phi_1\left(\frac{m+1}{M} - \delta\right) = m$ for $m = 0, 1, \dots, M - 2$;
- 1097 • ϕ_1 is linear on $\left[\frac{M-1}{M}, 1\right]$ and each interval $\left[\frac{m}{M}, \frac{m+1}{M} - \delta\right]$ for $m = 0, 1, \dots, M - 2$.

1098 Then

$$1099 \quad \phi_1(x) = m \quad \text{if } x \in \left[\frac{m}{M}, \frac{m+1}{M} - \delta \cdot \mathbf{1}_{\{m \leq M-2\}} \right] \quad \text{for } m = 0, 1, \dots, M - 1. \quad (5.6)$$

1100 Now consider another sample set

$$1101 \quad \begin{aligned} & \left\{ \left(\frac{1}{M}, L - 1 \right), (2, 0) \right\} \cup \left\{ \left(\frac{\ell}{ML}, \ell \right) : \ell = 0, 1, \dots, L - 1 \right\} \\ & \cup \left\{ \left(\frac{\ell+1}{ML} - \delta, \ell \right) : \ell = 0, 1, \dots, L - 2 \right\}. \end{aligned}$$

1102 Its size is $2L + 1 = 1 \cdot ((2L - 1) + 1) + 1$. By Lemma 5.4 (set $N_1 = 1$ and $N_2 = 2L - 1$
1103 therein), there exists

$$1104 \quad \begin{aligned} & \phi_2 \in \mathcal{NN}(\text{widthvec} = [2, 2(2L - 1) + 1]) \\ & = \mathcal{NN}(\text{widthvec} = [2, 4L - 1]) \end{aligned}$$

1105 such that

- 1106 • $\phi_2\left(\frac{L-1}{ML}\right) = \phi_2\left(\frac{1}{M}\right) = L - 1$ and $\phi_2\left(\frac{\ell}{ML}\right) = \phi_2\left(\frac{\ell+1}{ML} - \delta\right) = \ell$ for $\ell = 0, 1, \dots, L - 2$;
- 1107 • ϕ_2 is linear on $\left[\frac{L-1}{ML}, \frac{1}{M}\right]$ and each interval $\left[\frac{\ell}{ML}, \frac{\ell+1}{ML} - \delta\right]$ for $\ell = 0, 1, \dots, L - 2$.

1108 It follows that, for $m = 0, 1, \dots, M - 1$ and $\ell = 0, 1, \dots, L - 1$,

$$1109 \quad \phi_2\left(x - \frac{m}{M}\right) = \ell \quad \text{for } x \in \left[\frac{mL+\ell}{ML}, \frac{mL+\ell+1}{ML} - \delta \cdot \mathbb{1}_{\{\ell \leq L-2\}}\right]. \quad (5.7)$$

1110 $K = ML$ implies that any $k \in \{0, 1, \dots, K - 1\}$ can be unique represented by $k = mL + \ell$
 1111 for $m \in \{0, 1, \dots, M - 1\}$ and $\ell \in \{0, 1, \dots, L - 1\}$. Then the desired function ϕ can be
 1112 implemented by ReLU FNN as shown in Figure 13.

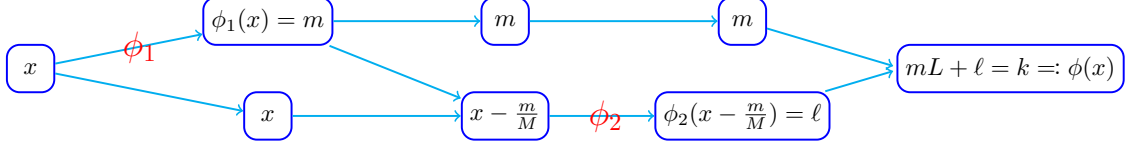


Figure 13: An illustration of the network architecture implementing ϕ based on Equations (5.6) and (5.7) with $x \in \left[\frac{k}{K}, \frac{k+1}{K} - \delta \cdot \mathbb{1}_{\{k \leq K-2\}}\right] = \left[\frac{mL+\ell}{ML}, \frac{mL+\ell+1}{ML} - \delta \cdot \mathbb{1}_{\{m \leq M-2 \text{ or } \ell \leq L-2\}}\right]$, where $k = mL + \ell$ for $m = 0, 1, \dots, M - 1$ and $\ell = 0, 1, \dots, L - 1$.

1113 Clearly,

$$1114 \quad \phi(x) = k \quad \text{if } x \in \left[\frac{k}{K}, \frac{k+1}{K} - \delta \cdot \mathbb{1}_{\{k \leq K-2\}}\right] \quad \text{for } k \in \{0, 1, \dots, K - 1\}.$$

1115 By Lemma 5.5, $\phi_1 \in \mathcal{NN}(\text{widthvec} = [2N, 4NL - 1]) \subseteq \mathcal{NN}(\text{width} \leq 4N + 2; \text{depth} \leq$
 1116 $2L + 1)$ and $\phi_2 \in \mathcal{NN}(\text{widthvec} = [2, 4L - 1]) \subseteq \mathcal{NN}(\text{width} \leq 6; \text{depth} \leq 2L + 1)$, implying
 1117 $\phi \in \mathcal{NN}(\text{width} \leq \max\{4N + 2 + 1, 6 + 1\} = 4N + 3; \text{depth} \leq (2L + 1) + 2 + (2L + 1) + 1 = 4L + 5)$.
 1118 So we finish the proof for the case $d = 1$

1119 **Case 2: $d \geq 2$.**

1120 Now we consider the case when $d \geq 2$. Consider the sample set

$$1121 \quad \left\{ (1, K - 1), (2, 0) \right\} \cup \left\{ \left(\frac{k}{K}, k \right) : k = 0, 1, \dots, K - 1 \right\} \\ \cup \left\{ \left(\frac{k+1}{K} - \delta, k \right) : k = 0, 1, \dots, K - 2 \right\},$$

1122 whose size is $2K + 1 = \lfloor N^{1/d} \rfloor ((2 \lfloor N^{1/d} \rfloor \lfloor L^{2/d} \rfloor - 1) + 1) + 1$. By Lemma 5.4 (set $N_1 = \lfloor N^{1/d} \rfloor$
 1123 and $N_2 = 2 \lfloor N^{1/d} \rfloor \lfloor L^{2/d} \rfloor - 1$ therein), there exists

$$1124 \quad \phi \in \mathcal{NN}(\text{widthvec} = [2 \lfloor N^{1/d} \rfloor, 2(2 \lfloor N^{1/d} \rfloor \lfloor L^{2/d} \rfloor - 1) + 1]) \\ = \mathcal{NN}(\text{widthvec} = [2 \lfloor N^{1/d} \rfloor, 4 \lfloor N^{1/d} \rfloor \lfloor L^{2/d} \rfloor - 1])$$

1125 such that

- 1126 • $\phi\left(\frac{K-1}{K}\right) = \phi(1) = K - 1$, and $\phi\left(\frac{k}{K}\right) = \phi\left(\frac{k+1}{K} - \delta\right) = k$ for $k = 0, 1, \dots, K - 2$;
- 1127 • ϕ is linear on $\left[\frac{K-1}{K}, 1\right]$ and each interval $\left[\frac{k}{K}, \frac{k+1}{K} - \delta\right]$ for $k = 0, 1, \dots, K - 2$.

1128 Then

$$1129 \quad \phi(x) = k \quad \text{if } x \in \left[\frac{k}{K}, \frac{k+1}{K} - \delta \cdot \mathbb{1}_{\{k \leq K-2\}}\right] \quad \text{for } k = 0, 1, \dots, K - 1.$$

1130 By Lemma 5.5,

$$1131 \quad \phi \in \mathcal{NN}(\text{widthvec} = [2 \lfloor N^{1/d} \rfloor, 4 \lfloor N^{1/d} \rfloor \lfloor L^{2/d} \rfloor - 1]) \\ \subseteq \mathcal{NN}(\text{width} \leq 4 \lfloor N^{1/d} \rfloor + 2; \text{depth} \leq 2 \lfloor L^{2/d} \rfloor + 1) \\ \subseteq \mathcal{NN}(\text{width} \leq 4 \lfloor N^{1/d} \rfloor + 3; \text{depth} \leq 4L + 5).$$

1132 which means we have finished the proof for the case $d \geq 2$. □

1133 5.3 Proof of Proposition 4.4 for point fitting

1134 In this sub-section, we will discuss how to use ReLU FNNs to fit a collection of points
 1135 in \mathbb{R}^2 .^⑨ It is trivial to fit n points via one-hidden-layer ReLU FNNs with $\mathcal{O}(n)$ param-
 1136 eters. However, to prove Proposition 4.4, we need to fit $\mathcal{O}(n)$ points with much fewer
 1137 parameters, which is the main difficulty of our proof. Our proof below is mainly based
 1138 on the “bit extraction” technique and the composition architecture of neural networks.

1139 Let us first introduce a basic lemma based on the “bit extraction” technique, which
 1140 is actually Lemma 2.6 of [41].

1141 **Lemma 5.6.** *For any $N, L \in \mathbb{N}^+$, any $\theta_{m,\ell} \in \{0, 1\}$ for $m = 0, 1, \dots, M-1$ and $\ell = 0, 1, \dots, L-$
 1142 1 , where $M = N^2L$, there exists a function ϕ implemented by a ReLU FNN with width
 1143 $4N + 3$ and depth $3L + 3$ such that*

$$1144 \quad \phi(m, \ell) = \sum_{j=0}^{\ell} \theta_{m,j} \quad \text{for } m = 0, 1, \dots, M-1 \text{ and } \ell = 0, 1, \dots, L-1.$$

1145 Next, let us introduce Lemma 5.7, a variant of Lemma 5.6 for a different mapping
 1146 for the “bit extraction”. Its proof is based on Lemmas 5.4, 5.5, and 5.6.

1147 **Lemma 5.7.** *For any $N, L \in \mathbb{N}^+$ and any $\theta_i \in \{0, 1\}$ for $i = 0, 1, \dots, N^2L^2 - 1$, there exists
 1148 a function ϕ implemented by a ReLU FNN with width $8N + 6$ and depth $5L + 7$ such that*

$$1149 \quad \phi(i) = \theta_i \quad \text{for } i = 0, 1, \dots, N^2L^2 - 1.$$

1150 *Proof.* The case $L = 1$ is clear. We assume $L \geq 2$ below.

1151 Denote $M = N^2L$, for each $i \in \{0, 1, \dots, N^2L^2 - 1\}$, there exists a unique representation
 1152 $i = mL + \ell$ for $m \in \{0, 1, \dots, M-1\}$ and $\ell \in \{0, 1, \dots, L-1\}$. Thus, we can define, for
 1153 $m = 0, 1, \dots, M-1$ and $\ell = 0, 1, \dots, L-1$,

$$1154 \quad a_{m,\ell} := \theta_i, \quad \text{where } i = mL + \ell.$$

1155 Then, for $m = 0, 1, \dots, M-1$, we set $b_{m,0} = 0$ and $b_{m,\ell} = a_{m,\ell-1}$ for $\ell = 1, 2, \dots, L-1$.

1156 By Lemma 5.6, there exist $\phi_1, \phi_2 \in \mathcal{NN}$ (width $\leq 4N + 3$; depth $\leq 3L + 3$) such that

$$1157 \quad \phi_1(m, \ell) = \sum_{j=0}^{\ell} a_{m,j} \quad \text{and} \quad \phi_2(m, \ell) = \sum_{j=0}^{\ell} b_{m,j}$$

1158 for $m = 0, 1, \dots, M-1$ and $\ell = 0, 1, \dots, L-1$.

1159 We consider the sample set

$$1160 \quad \{(mL, m) : m = 0, 1, \dots, M\} \cup \{((m+1)L - 1, m) : m = 0, 1, \dots, M-1\}.$$

1161 Its size is $2M + 1 = N \cdot ((2NL - 1) + 1) + 1$. By Lemma 5.4 (set $N_1 = N$ and $N_2 = 2NL - 1$
 1162 therein), there exists

$$1163 \quad \begin{aligned} \psi &\in \mathcal{NN}(\text{widthvec} = [2N, 2(2NL - 1) + 1]) \\ &= \mathcal{NN}(\text{widthvec} = [2N, 4NL - 1]) \end{aligned}$$

1164 such that

^⑨Fitting a collection of points $\{(x_i, y_i)\}_i$ in \mathbb{R}^2 means that the target ReLU FNN takes a value close to y_i at the location x_i .

- 1165 • $\psi(ML) = M$ and $\psi(mL) = \psi((m+1)L - 1) = m$ for $m = 0, 1, \dots, M-1$;
 1166 • ψ is linear on each interval $[mL, (m+1)L - 1]$ for $m = 0, 1, \dots, M-1$.

1167 It follows that

1168
$$\psi(x) = m \quad \text{if } x \in [mL, (m+1)L - 1] \quad \text{for } m = 0, 1, \dots, M-1,$$

1169 implying

1170
$$\psi(mL + \ell) = m \quad \text{for } m = 0, 1, \dots, M-1 \text{ and } \ell = 0, 1, \dots, L-1.$$

1171 For $i = 0, 1, \dots, N^2L^2 - 1$, by representing $i = mL + \ell$ for $m = 0, 1, \dots, M-1$ and
 1172 $\ell = 0, 1, \dots, L-1$, we have $\psi(i) = \psi(mL + \ell) = m$ and $i - L\psi(i) = \ell$, from which we deduce

1173
$$\begin{aligned} & \phi_1(\psi(i), i - L\psi(i)) - \phi_2(\psi(i), i - L\psi(i)) \\ &= \phi_1(m, \ell) - \phi_2(m, \ell) = \sum_{j=0}^{\ell} a_{m,j} - \sum_{j=0}^{\ell} b_{m,j} \\ &= \sum_{j=0}^{\ell} a_{m,j} - \sum_{j=1}^{\ell} a_{m,j-1} - b_0 = a_{m,\ell} = \theta_i. \end{aligned} \tag{5.8}$$

1174 Therefore, the desired function ϕ can be implemented by the network architecture
 1175 described in Figure 14.

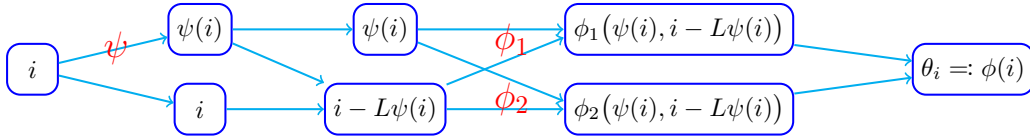


Figure 14: An illustration of the network architecture implementing the desired function ϕ based on Equation (5.8).

1176 Note that

1177
$$\phi_1, \phi_2 \in \mathcal{NN}(\text{width} \leq 4N + 3; \text{depth} \leq 3L + 3).$$

1178 And by Lemma 5.5,

1179
$$\begin{aligned} \psi &\in \mathcal{NN}(\text{widthvec} = [2N, 4NL - 1]) \\ &\subseteq \mathcal{NN}(\text{width} \leq 4N + 2; \text{depth} \leq 2L + 1). \end{aligned}$$

1180 Hence, the network architecture shown in Figure 14 is with width $\max\{4L + 2 + 1, 2(4L +$
 1181 $3)\} = 8N + 6$ and depth $(2L + 1) + 2 + (3L + 3) + 1 = 5L + 7$, implying $\phi \in \mathcal{NN}(\text{width} \leq$
 1182 $8N + 6; \text{depth} \leq 5L + 7)$. So we finish the proof. \square

1183 With Lemma 5.7 in hand, we are now ready to prove Proposition 4.4.

1184 *Proof of Proposition 4.4.* Set $J = \lceil 2s \log_2(NL + 1) \rceil \in \mathbb{N}^+$. For each $\xi_i \in [0, 1]$, there exist
 1185 $\xi_{i,1}, \xi_{i,2}, \dots, \xi_{i,J} \in \{0, 1\}$ such that

1186
$$|\xi_i - \text{bin}_{0.\xi_{i,1}\xi_{i,2}\dots\xi_{i,J}}| \leq 2^{-J} \quad \text{for } i = 0, 1, \dots, N^2L^2 - 1.$$

1187 By Lemma 5.7, there exist

1188
$$\phi_1, \phi_2, \dots, \phi_J \in \mathcal{NN}(\text{width} \leq 8N + 6; \text{depth} \leq 5L + 7)$$

1189 such that

1190
$$\phi_j(i) = \xi_{i,j} \quad \text{for } i = 0, 1, \dots, N^2L^2 - 1 \text{ and } j = 1, 2, \dots, J.$$

1191 Define

1192
$$\tilde{\phi}(x) := \sum_{j=1}^J 2^{-j} \phi_j(x) \quad \text{for any } x \in \mathbb{R}.$$

1193 It follows that, for $i = 0, 1, \dots, N^2L^2 - 1$,

1194
$$\begin{aligned} |\tilde{\phi}(i) - \xi_i| &= \left| \sum_{j=1}^J 2^{-j} \phi_j(i) - \xi_i \right| = \left| \sum_{j=1}^J 2^{-j} \xi_{i,j} - \xi_i \right| \\ &= \left| \text{bin} 0.\xi_{i,1}\xi_{i,2}\dots\xi_{i,J} - \xi_i \right| \leq 2^{-J} \leq N^{-2s} L^{-2s}, \end{aligned}$$

1195 where the last inequality comes from

1196
$$2^{-J} = 2^{-\lceil 2s \log_2(NL+1) \rceil} \leq 2^{-2s \log_2(NL+1)} = (NL+1)^{-2s} \leq N^{-2s} L^{-2s}.$$

1197 Now let us estimate the width and depth of the network implementing $\tilde{\phi}$. Recall
1198 that

1199
$$\begin{aligned} J = \lceil 2s \log_2(NL+1) \rceil &\leq 2s(1 + \log_2(NL+1)) \leq 2s(1 + \log_2(2N) + \log_2 L) \\ &\leq 2s(1 + \log_2(2N))(1 + \log_2 L) \leq 2s \lceil \log_2(4N) \rceil \lceil \log_2(2L) \rceil, \end{aligned}$$

1200 and $\phi_j \in \mathcal{NN}(\text{width} \leq 8N + 6; \text{depth} \leq 5L + 7)$ for each j .

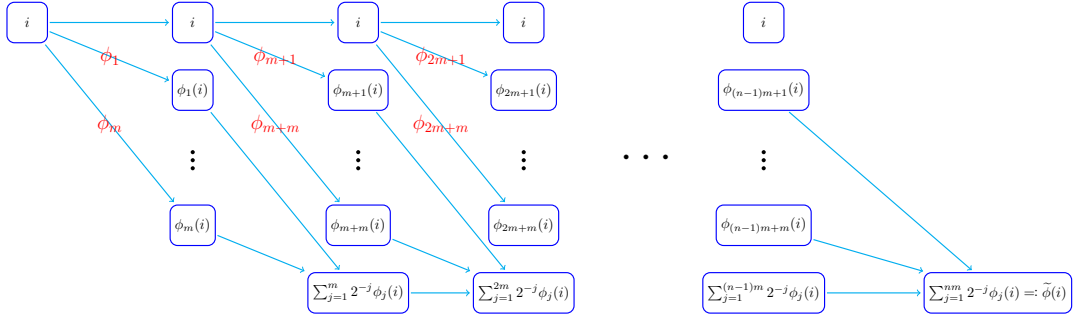


Figure 15: An illustration of the network architecture implementing $\tilde{\phi} = \sum_{j=1}^J 2^{-j} \phi_j$ for any $i \in \{0, 1, \dots, N^2L^2 - 1\}$. We assume $J = mn$, where $m = 2s \lceil \log_2(4N) \rceil$ and $n = \lceil \log_2(2L) \rceil$, since we can set $\phi_{J+1} = \dots = \phi_{nm} = 0$ if $J < nm$.

1201 As we can see from Figure 15, $\tilde{\phi} = \sum_{j=1}^J 2^{-j} \phi_j$ can be implemented by a ReLU FNN
1202 with width

1203
$$\begin{aligned} (8N + 6)m + (1 + m + 1) &= (8N + 6)2s \lceil \log_2(4N) \rceil + 2s \lceil \log_2(4N) \rceil + 2 \\ &\leq 16s(N + 1) \log_2(8N) \end{aligned}$$

1204 and depth

1205
$$((5L + 7) + 1)n = (5L + 8) \lceil \log_2(2L) \rceil \leq (5N + 8) \log_2(4L).$$

1206 Finally, we define

1207
$$\phi(x) := \min \{ \sigma(\tilde{\phi}(x)), 1 \} = \min \{ \max\{0, \tilde{\phi}(x)\}, 1 \} \quad \text{for any } x \in \mathbb{R}.$$

1208 Then $0 \leq \phi(x) \leq 1$ for any $x \in \mathbb{R}$ and ϕ can be implemented by a ReLU FNN with width
 1209 $16s(N+1)\log_2(8N)$ and depth $(5L+8)\log_2(4L)+3 \leq 5(L+2)\log_2(4L)$. See Figure 16
 1210 for the network architecture implementing ϕ . Note that

1211
$$\tilde{\phi}(i) = \sum_{j=1}^J 2^{-j} \phi_j(i) = \sum_{j=1}^J 2^{-j} \xi_{i,j} \in [0, 1] \quad \text{for } i = 0, 1, \dots, N^2 L^2 - 1.$$

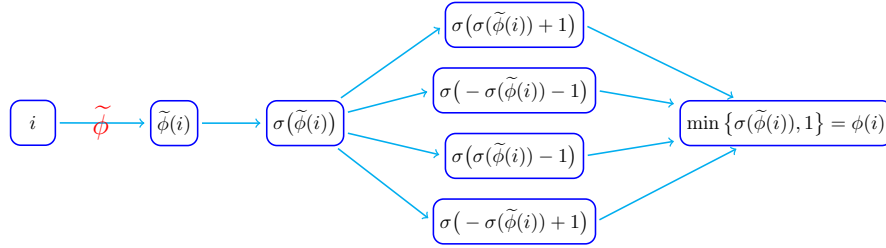


Figure 16: An illustration of the network architecture implementing the desired function ϕ based on the fact that $\min\{x_1, x_2\} = \frac{x_1+x_2-|x_1-x_2|}{2} = \frac{\sigma(x_1+x_2)-\sigma(-x_1-x_2)-\sigma(x_1-x_2)-\sigma(-x_1+x_2)}{2}$.

1212 It follows that

1213
$$|\phi(i) - \xi_i| = \left| \min \{ \max\{0, \tilde{\phi}(i)\}, 1 \} - \xi_i \right| = |\tilde{\phi}(i) - \xi_i| \leq N^{-2s} L^{-2s},$$

1214 for $i = 0, 1, \dots, N^2 L^2 - 1$. The proof is complete. □

1215 6 Conclusions

1216 This paper has established a nearly optimal approximation error of ReLU FNNs
 1217 in terms of both width and depth to approximate smooth functions. It is shown that
 1218 ReLU FNNs with width $\mathcal{O}(N \ln N)$ and depth $\mathcal{O}(L \ln L)$ can approximate functions in
 1219 the unit ball of $C^s([0, 1]^d)$ with an approximation error $\mathcal{O}(N^{-2s/d} L^{-2s/d})$. Through VC-
 1220 dimension, it is also proved that this approximation error is asymptotically nearly tight
 1221 for the closed unit ball of $C^s([0, 1]^d)$.

1222 We would like to remark that our analysis is for the fully connected feed-forward
 1223 neural networks with the ReLU activation function. It would be an interesting direction
 1224 for further study to generalize our results to neural networks with other architectures
 1225 (e.g., convolutional neural networks and ResNet) and activation functions (e.g., tanh
 1226 and sigmoid functions). These will be subjects of future work.

1227 Acknowledgments

1228 The work of J. Lu is supported in part by the National Science Foundation via
 1229 grants DMS-1415939, CCF-1934964, and DMS-2012286. Z. Shen is supported by Tan
 1230 Chin Tuan Centennial Professorship. H. Yang H. Yang was partially supported by the
 1231 National Science Foundation under award DMS-1945029.

1232 References

- 1233 [1] Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in
1234 overparameterized neural networks, going beyond two layers. *arXiv e-prints*, page
1235 arXiv:1811.04918, November 2018.
- 1236 [2] Martin Anthony and Peter L. Bartlett. *Neural Network Learning: Theoretical Foun-*
1237 *ndations*. Cambridge University Press, New York, NY, USA, 1st edition, 2009.
- 1238 [3] Sanjeev Arora, Simon S. Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained
1239 analysis of optimization and generalization for overparameterized two-layer neural
1240 networks. In *ICML*, 2019.
- 1241 [4] Chenglong Bao, Qianxiao Li, Zuwei Shen, Cheng Tai, Lei Wu, and Xueshuang
1242 Xiang. Approximation analysis of convolutional neural networks. 2019.
- 1243 [5] A. R. Barron. Universal approximation bounds for superpositions of a sigmoidal
1244 function. *IEEE Transactions on Information Theory*, 39(3):930–945, May 1993.
- 1245 [6] Andrew R. Barron and Jason M. Klusowski. Approximation and estimation for
1246 high-dimensional deep learning networks. *arXiv e-prints*, page arXiv:1809.03090,
1247 September 2018.
- 1248 [7] Peter Bartlett, Vitaly Maiorov, and Ron Meir. Almost linear VC-dimension bounds
1249 for piecewise polynomial networks. *Neural Computation*, 10:2159–2173, 1998.
- 1250 [8] M. Bianchini and F. Scarselli. On the complexity of neural network classifiers: A
1251 comparison between shallow and deep architectures. *IEEE Transactions on Neural*
1252 *Networks and Learning Systems*, 25(8):1553–1565, Aug 2014.
- 1253 [9] Helmut. Bölcskei, Philipp. Grohs, Gitta. Kutyniok, and Philipp. Petersen. Optimal
1254 approximation with sparsely connected deep neural networks. *SIAM Journal on*
1255 *Mathematics of Data Science*, 1(1):8–45, 2019.
- 1256 [10] Yuan Cao and Quanquan Gu. Generalization bounds of stochastic gradient descent
1257 for wide and deep neural networks. *CoRR*, abs/1905.13210, 2019.
- 1258 [11] Liang Chen and Congwei Wu. A note on the expressive power of deep rectified linear
1259 unit networks in high-dimensional spaces. *Mathematical Methods in the Applied*
1260 *Sciences*, 42(9):3400–3404, 2019.
- 1261 [12] Minshuo Chen, Haoming Jiang, Wenjing Liao, and Tuo Zhao. Efficient approx-
1262 imation of deep ReLU networks for functions on low dimensional manifolds. In
1263 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett,
1264 editors, *Advances in Neural Information Processing Systems 32*, pages 8174–8184.
1265 Curran Associates, Inc., 2019.
- 1266 [13] Zixiang Chen, Yuan Cao, Difan Zou, and Quanquan Gu. How much
1267 over-parameterization is sufficient to learn deep ReLU networks? *CoRR*,
1268 arXiv:1911.12360, 2019.

- 1269 [14] Charles K. Chui, Shao-Bo Lin, and Ding-Xuan Zhou. Construction of neural net-
1270 works for realization of localized deep learning. *Frontiers in Applied Mathematics*
1271 *and Statistics*, 4:14, 2018.
- 1272 [15] George Cybenko. Approximation by superpositions of a sigmoidal function. *MCSSS*,
1273 2:303–314, 1989.
- 1274 [16] Ronald A. Devore. Optimal nonlinear approximation. *Manuskripta Math*, pages
1275 469–478, 1989.
- 1276 [17] Weinan E, Chao Ma, and Qingcan Wang. A priori estimates of the population risk
1277 for residual networks. *ArXiv*, abs/1903.02154, 2019.
- 1278 [18] Weinan E, Chao Ma, and Lei Wu. A priori estimates of the population risk for two-
1279 layer neural networks. *Communications in Mathematical Sciences*, 17(5):1407–1425,
1280 2019.
- 1281 [19] Weinan E and Qingcan Wang. Exponential convergence of the deep neural network
1282 approximation for analytic functions. *CoRR*, abs/1807.00297, 2018.
- 1283 [20] Rémi Gribonval, Gitta Kutyniok, Morten Nielsen, and Felix Voigtlaender. Approx-
1284 imation spaces of deep neural networks. *arXiv e-prints*, page arXiv:1905.01208, May
1285 2019.
- 1286 [21] Ingo Gühring, Gitta Kutyniok, and Philipp Petersen. Error bounds for approx-
1287 imations with deep ReLU neural networks in $W^{s,p}$ norms. *arXiv e-prints*, page
1288 arXiv:1902.07896, Feb 2019.
- 1289 [22] Nick Harvey, Christopher Liaw, and Abbas Mehrabian. Nearly-tight VC-dimension
1290 bounds for piecewise linear neural networks. In Satyen Kale and Ohad Shamir,
1291 editors, *Proceedings of the 2017 Conference on Learning Theory*, volume 65 of *Pro-*
1292 *ceedings of Machine Learning Research*, pages 1064–1068, Amsterdam, Netherlands,
1293 07–10 Jul 2017. PMLR.
- 1294 [23] Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural*
1295 *Networks*, 4(2):251–257, 1991.
- 1296 [24] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward
1297 networks are universal approximators. *Neural Networks*, 2(5):359–366, 1989.
- 1298 [25] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Con-
1299 vergence and generalization in neural networks. *CoRR*, abs/1806.07572, 2018.
- 1300 [26] Ziwei Ji and Matus Telgarsky. Polylogarithmic width suffices for gradient de-
1301 scent to achieve arbitrarily small test error with shallow ReLU networks. *ArXiv*,
1302 abs/1909.12292, 2020.
- 1303 [27] Michael J. Kearns and Robert E. Schapire. Efficient distribution-free learning of
1304 probabilistic concepts. *J. Comput. Syst. Sci.*, 48(3):464–497, June 1994.

- 1305 [28] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification
1306 with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou,
1307 and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*,
1308 volume 25, pages 1097–1105. Curran Associates, Inc., 2012.
- 1309 [29] Qianxiao Li, Ting Lin, and Zuowei Shen. Deep learning via dynamical systems: An
1310 approximation perspective. *Journal of European Mathematical Society*, to appear.
- 1311 [30] Shiyu Liang and R. Srikant. Why deep neural networks? *CoRR*, abs/1610.04161,
1312 2016.
- 1313 [31] Hadrien Montanelli and Qiang Du. New error bounds for deep networks using sparse
1314 grids. *SIAM Journal on Mathematics of Data Science*, 1(1):78–92, 2019.
- 1315 [32] Hadrien Montanelli and Haizhao Yang. Error bounds for deep ReLU networks using
1316 the Kolmogorov–Arnold superposition theorem. *Neural Networks*, 129:1–6, 2020.
- 1317 [33] Hadrien Montanelli, Haizhao Yang, and Qiang Du. Deep ReLU networks over-
1318 come the curse of dimensionality for bandlimited functions. *arXiv e-prints*, page
1319 arXiv:1903.00735, March 2019.
- 1320 [34] Guido F Montufar, Razvan Pascanu, Kyunghyun Cho, and Yoshua Bengio. On the
1321 number of linear regions of deep neural networks. In Z. Ghahramani, M. Welling,
1322 C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural
1323 Information Processing Systems 27*, pages 2924–2932. Curran Associates, Inc., 2014.
- 1324 [35] Ryumei Nakada and Masaaki Imaizumi. Adaptive approximation and generalization
1325 of deep neural network with intrinsic dimensionality. *Journal of Machine Learning
1326 Research*, 21(174):1–38, 2020.
- 1327 [36] J. A. A. Opschoor, Ch. Schwab, and J. Zech. Exponential ReLU DNN expression
1328 of holomorphic maps in high dimension. Technical Report 2019-35, Seminar for
1329 Applied Mathematics, ETH Zürich, Switzerland., 2019.
- 1330 [37] Philipp Petersen and Felix Voigtlaender. Optimal approximation of piecewise
1331 smooth functions using deep ReLU neural networks. *Neural Networks*, 108:296–
1332 330, 2018.
- 1333 [38] T. Poggio, H. N. Mhaskar, L. Rosasco, B. Miranda, and Q. Liao. Why and when
1334 can deep—but not shallow—networks avoid the curse of dimensionality: A review.
1335 *International Journal of Automation and Computing*, 14:503–519, 2017.
- 1336 [39] Akito Sakurai. Tight bounds for the VC-dimension of piecewise polynomial net-
1337 works. In *Advances in Neural Information Processing Systems*, pages 323–329.
1338 Neural information processing systems foundation, 1999.
- 1339 [40] Zuowei Shen, Haizhao Yang, and Shijun Zhang. Nonlinear approximation via com-
1340 positions. *Neural Networks*, 119:74–84, 2019.

- 1341 [41] Zuowei Shen, Haizhao Yang, and Shijun Zhang. Deep network approximation
1342 characterized by number of neurons. *Communications in Computational Physics*,
1343 28(5):1768–1811, 2020.
- 1344 [42] Zuowei Shen, Haizhao Yang, and Shijun Zhang. Optimal approximation rate of
1345 ReLU networks in terms of width and depth. *Journal de Mathématiques Pures et*
1346 *Appliquées*, to appear.
- 1347 [43] Taiji Suzuki. Adaptivity of deep ReLU network for learning in Besov and mixed
1348 smooth Besov spaces: optimal rate and curse of dimensionality. In *International*
1349 *Conference on Learning Representations*, 2019.
- 1350 [44] Dmitry Yarotsky. Error bounds for approximations with deep ReLU networks.
1351 *Neural Networks*, 94:103–114, 2017.
- 1352 [45] Dmitry Yarotsky. Optimal approximation of continuous functions by very deep
1353 ReLU networks. In Sébastien Bubeck, Vianney Perchet, and Philippe Rigollet,
1354 editors, *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Pro-*
1355 *ceedings of Machine Learning Research*, pages 639–649. PMLR, 06–09 Jul 2018.
- 1356 [46] Dmitry Yarotsky and Anton Zhevnerchuk. The phase diagram of approximation
1357 rates for deep neural networks. In H. Larochelle, M. Ranzato, R. Hadsell, M. F.
1358 Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*,
1359 volume 33, pages 13005–13015. Curran Associates, Inc., 2020.
- 1360 [47] Ding-Xuan Zhou. Universality of deep convolutional neural networks. *Applied and*
1361 *Computational Harmonic Analysis*, 48(2):787–794, 2020.