

DEEP NETWORK APPROXIMATION CHARACTERIZED BY NUMBER OF NEURONS*

ZUOWEI SHEN[†], HAIZHAO YANG[‡], AND SHIJUN ZHANG[§]

Abstract. This paper quantitatively characterizes the approximation power of deep feed-forward neural networks (FNNs) in terms of the number of neurons, i.e., the product of the network width and depth. It is shown by construction that ReLU FNNs with width $\max\{8d\lfloor N^{1/d} \rfloor + 4d, 12N + 14\}$ and depth $9L + 12$ can approximate an arbitrary Hölder continuous function of order α with a Lipschitz constant ν on $[0, 1]^d$ with a tight approximation rate $5(8\sqrt{d})^\alpha \nu N^{-2\alpha/d} L^{-2\alpha/d}$ for any given $N, L \in \mathbb{N}^+$. The constructive approximation is a corollary of a more general result for an arbitrary continuous function f in terms of its modulus of continuity $\omega_f(\cdot)$. In particular, the approximation rate of ReLU FNNs with width $\max\{8d\lfloor N^{1/d} \rfloor + 4d, 12N + 14\}$ and depth $9L + 12$ for a general continuous function f is $5\omega_f(8\sqrt{d}N^{-2/d}L^{-2/d})$. We also extend our analysis to the case when the domain of f is irregular or localized in an ε -neighborhood of a $d_{\mathcal{M}}$ -dimensional smooth manifold $\mathcal{M} \subseteq [0, 1]^d$ with $d_{\mathcal{M}} \ll d$. Especially, in the case of an essentially low-dimensional domain, we show an approximation rate $3\omega_f(\frac{4\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}}) + 5\omega_f(\frac{16d}{(1-\delta)\sqrt{d_\delta}}N^{-2/d_\delta}L^{-2/d_\delta})$ for ReLU FNNs to approximate f in the ε -neighborhood, where $d_\delta = \mathcal{O}(d_{\mathcal{M}} \frac{\ln(d/\delta)}{\delta^2})$ for any given $\delta \in (0, 1)$. Our analysis provides a general guide for selecting the width and the depth of ReLU FNNs to approximate continuous functions especially in parallel computing.

Key words. Deep ReLU Neural Networks, Hölder Continuity, Modulus of Continuity, Approximation Theory, Low-Dimensional Manifold, Parallel Computing.

1. Introduction. The approximation theory of neural networks has been an active research topic in the past few decades. Previously, as a special kind of ridge function approximation, shallow neural networks with one hidden layer and various activation functions (e.g., wavelets pursuits [42, 10], adaptive splines [19, 50], radial basis functions [8, 48, 18, 23, 60], sigmoid functions [27, 41, 34, 7, 38, 35, 14, 13, 15]) were widely discussed and admit good approximation properties, e.g., the universal approximation property [16, 28, 27], overcoming the curse of dimensionality [4], and providing attractive approximation rate in nonlinear approximation [19, 42, 10, 50, 18, 23, 60].

The introduction of deep networks with more than one hidden layers has made significant impacts in many fields in computer science and engineering including computer vision [32] and natural language processing [1]. New scientific computing tools based on deep networks have also emerged and facilitated large-scale and high-dimensional problems that were impractical previously [22, 20]. The design of deep ReLU networks is the key of such a revolution. These breakthroughs have stimulated broad research topics from different points of views to study the power of deep ReLU networks, e.g. in terms of combinatorics [46], topology [6], Vapnik-Chervonenkis (VC) dimension [5, 53, 25], fat-shattering dimension [31, 2], information theory [49], classical approximation theory [16, 28, 4, 62, 57, 57], optimization [29, 47, 30] etc.

Particularly in approximation theory, **non-quantitative and asymptotic** approximation rates of ReLU FNNs have been proposed for various types of functions. For example, smooth functions [39, 36, 61, 21], piecewise smooth functions [49], band-limited functions [45], continuous functions [62]. However, to the best of our knowledge, existing theories [39, 45, 61, 36, 44, 58, 49, 62, 21, 17] can only provide implicit formulas in the sense that the approximation error contains an unknown prefactor, or work only for sufficiently large N and L larger than some unknown numbers. For example, [62] estimated an approximation rate $c(d)L^{-2\alpha/d}$ via a narrow and deep ReLU FNN, where $c(d)$ is an unknown number depending on d and L is required to be larger than a sufficiently large unknown number \mathcal{L} . For another example, given an approximation error ε , [49] proved the existence of a ReLU FNN with a constant but still unknown number of layers approximating a C^β function within the target error. These works can be divided into two cases: 1) FNNs with varying width

*Submitted to the editors June 13, 2019.

Funding: H. Yang was supported by the startup of the Department of Mathematics at the National University of Singapore and Ministry of Education in Singapore under the grant MOE2018-T2-2-147.

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

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

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

and only one hidden layer [18, 23, 37, 60] (visualized by the region in ■ in Figure 1); 2) FNNs with a fixed width of $\mathcal{O}(d)$ and a varying depth larger than an unknown number \mathcal{L} [40, 62] (represented by the region in ■ in Figure 1).

As far as we know, the first **quantitative and non-asymptotic** approximation rate of deep ReLU FNNs was obtained in [57]. Specifically, [57] identified an explicit formulas of the approximation rate

$$(1.1) \quad \begin{cases} 2\nu N^{-2\alpha}, & \text{when } L \geq 2 \text{ and } d = 1, \\ 2(2\sqrt{d})^\alpha \nu N^{-2\alpha/d}, & \text{when } L \geq 3 \text{ and } d \geq 2, \end{cases}$$

for ReLU FNNs with an arbitrary width $N \in \mathbb{N}^+$ and a fixed depth $L \in \mathbb{N}^+$ to approximate a Hölder continuous function f of order α with a constant ν (visualized in the region shown by ■ in Figure 1). The approximation rate $\mathcal{O}(N^{-2\alpha/d})$ is tight in terms of N and increasing L cannot improve the approximation rate in N . The success of deep FNNs in a broad range of applications has motivated a well-known conjecture that the depth L has an important role in improving the approximation power of deep FNNs. In particular, a very important question in practice would be, given an arbitrary L and N , what is the explicit formula to characterize the approximation error so as to see whether the network is large enough to meet the accuracy requirement. Due to the highly nonlinear structure of deep FNNs, it is still a challenging open problem to characterize N and L simultaneously in the approximation rate.

To answer the question just above, we establish the first framework that is able to quantify the approximation power of deep ReLU FNNs with arbitrary width N and depth L , achieving the optimal approximation rate, $5\omega_f(8\sqrt{d}N^{-2/d}L^{-2/d})$, for continuous functions $f \in C([0, 1]^d)$. Our result is based on new analysis techniques merely based on the structure of FNNs and a modified bit extraction technique inspired by [5], instead of designing FNNs to approximate traditional approximation basis like polynomials and splines as in the existing literature [52, 36, 51, 61, 62, 44, 24, 39, 49, 55, 58]. The approximation rate obtained here admits an explicit formula to compute the prefactor when $\omega_f(\cdot)$ is known. For example, in the case of Hölder continuous functions of order α with a Lipschitz constant ν , $\omega_f(r) \leq \nu r^\alpha$ for $r \geq 0$, resulting in the approximation rate $5(8\sqrt{d})^\alpha \nu N^{-2\alpha/d} L^{-2\alpha/d}$ as mentioned previously. As a consequence, existing works for the function class $C([0, 1]^d)$ are special cases of our result (see Figure 1 for a comparison).

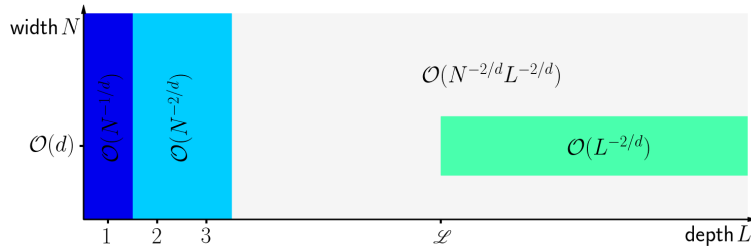


FIG. 1. A summary of existing and our new results on the approximation rate of ReLU FNNs for continuous functions. Existing results [18, 23, 37, 60, 40, 62, 57] are applicable in the areas in ■, ■, and ■; our new result is suitable for almost all areas when $L \geq 2$.

Our main result, Theorem 1.1 below, shows that ReLU FNNs with width $\max\{8d\lfloor N^{1/d} \rfloor + 4d, 12N + 14\}$ and depth $9L + 12$ can approximate f with an approximation rate $5\omega_f(8\sqrt{d}N^{-2/d}L^{-2/d})$, where $\omega_f(\cdot)$ is the modulus of continuity of f defined via

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

where $|\cdot|$ is the absolute value of a scalar or the length of a vector.

THEOREM 1.1. *Let f be a given function in $C([0, 1]^d)$. For any arbitrary $L \in \mathbb{N}^+$, $N \in \mathbb{N}^+$, and $\eta > 0$, there exists a ReLU FNN ϕ with width $\max\{8d\lfloor N^{1/d} \rfloor + 4d, 12N + 14\}$ and depth $9L + 12$ such that*

$$\|f - \phi\|_{L^\infty([0, 1]^d \setminus \mathcal{H})} \leq 5\omega_f(8\sqrt{d}N^{-2/d}L^{-2/d}), \quad \text{and} \quad \|\phi\|_{L^\infty([0, 1]^d)} \leq |f(0)| + \omega_f(\sqrt{d}),$$

where $\mathcal{H} \subseteq [0, 1]^d$ with $\mu(\mathcal{H}) \leq \eta$.

Note that ϕ in Theorem 1.1 is uniformly bounded, the measure of \mathcal{H} can be arbitrarily small and is independent of f , N , and L . Hence, the L^p -norm of the approximation error can also be controlled as follows.

COROLLARY 1.2. *Let f be a given function in $C([0, 1]^d)$. For arbitrary $L \in \mathbb{N}^+$, $N \in \mathbb{N}^+$, and $p \in [1, \infty)$, there exists a ReLU FNN ϕ with width $\max\{8d\lfloor N^{1/d} \rfloor + 4d, 12N + 14\}$ and depth $9L + 12$ such that*

$$\|f - \phi\|_{L^p([0, 1]^d)} \leq 5\omega_f(8\sqrt{d}N^{-2/d}L^{-2/d}).$$

When Theorem 1.1 and Corollary 1.2 are applied to a Hölder continuous function f of order α with a Lipschitz constant ν , i.e., $f \in \text{Lip}(\nu, \alpha, d)$, the approximation rate is $5(8\sqrt{d})^\alpha \nu N^{-2\alpha/d} L^{-2\alpha/d}$, because $\omega_f(r) \leq \nu r^\alpha$ for any $r \geq 0$. An immediate question following the constructive approximation is how much we can improve the approximation rate. In fact, the approximation rate of $f \in \text{Lip}(\nu, \alpha, d)$ is asymptotically tight based on VC dimension as we shall see later.

In most real applications of neural networks, though the target function f is defined in a high dimensional domain, e.g., $[0, 1]^d$, where d could be tens of thousands or even millions, only the approximation error of f in a neighborhood of a $d_{\mathcal{M}}$ -dimensional manifold \mathcal{M} with $d_{\mathcal{M}} \ll d$ is concerned. Hence, we extend Theorem 1.1 to the case when the domain of f is localized in an ε -neighborhood of a compact $d_{\mathcal{M}}$ -dimensional Riemannian submanifold $\mathcal{M} \subseteq [0, 1]^d$ having condition number $1/\tau$, volume V , and geodesic covering regularity \mathcal{R} . The ε -neighborhood is defined as

$$(1.2) \quad \mathcal{M}_\varepsilon := \{\mathbf{x} \in [0, 1]^d : \inf\{|\mathbf{x} - \mathbf{y}| : \mathbf{y} \in \mathcal{M}\} \leq \varepsilon\}, \quad \text{for } \varepsilon \in (0, 1).$$

Let $d_\delta = \mathcal{O}\left(\frac{d_{\mathcal{M}} \ln(dV\mathcal{R}\tau^{-1}\delta^{-1})}{\delta^2}\right) = \mathcal{O}\left(d_{\mathcal{M}} \frac{\ln(d/\delta)}{\delta^2}\right)$ be an integer for any $\delta \in (0, 1)$ such that $d_{\mathcal{M}} \leq d_\delta \leq d$. We show an approximation rate $3\omega_f\left(\frac{4\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}}\right) + 5\omega_f\left(\frac{16d}{(1-\delta)\sqrt{d_\delta}}N^{-2/d_\delta}L^{-2/d_\delta}\right)$ for ReLU FNNs to approximate f in a $L^p([0, 1]^d, \mu_\varrho)$ -norm, where ϱ is a probability density function (i.e., $\int \varrho(\mathbf{x})d\mathbf{x} = 1$) supported in \mathcal{M}_ε , and $\mu_\varrho(\cdot)$ is the corresponding measure of ϱ defined via $\mu_g(E) := \int_E g(\mathbf{x})d\mathbf{x}$ for any measurable set $E \subseteq \mathbb{R}^d$. The key ideas of the proof is the application of Theorem 3.1 in [3], which provides a nearly isometric projection $\mathbf{A} \in \mathbb{R}^{d_\delta \times d}$ that maps points in $\mathcal{M} \subseteq [0, 1]^d$ to a d_δ -dimensional domain with

$$(1 - \delta)|\mathbf{x}_1 - \mathbf{x}_2| \leq |\mathbf{A}\mathbf{x}_1 - \mathbf{A}\mathbf{x}_2| \leq (1 + \delta)|\mathbf{x}_1 - \mathbf{x}_2| \quad \text{for any } \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{M},$$

and the application of Theorem 1.1 in this paper, which constructs the desired ReLU FNN with a size depending on d_δ instead of d to lessen the curse of dimensionality. When δ is closer to 1, d_δ is closer to $d_{\mathcal{M}}$ but the isometric property of the projection is weakened; when δ is closer to 0, the isometric property becomes better but d_δ could be larger than d , in which case we can simply enforce $d_\delta = d$ and choose the identity map as the projection. Hence, $\delta \in (0, 1)$ is a parameter to make a balance between isometry and dimension reduction.

THEOREM 1.3. *Let f be a continuous function on $[0, 1]^d$ and $\mathcal{M} \subseteq [0, 1]^d$ be a compact $d_{\mathcal{M}}$ -dimensional Riemannian submanifold. For any $\varepsilon \in (0, 1)$, let ϱ be a probability density function supported on \mathcal{M}_ε defined in (1.2) with $\mu_\varrho(\cdot)$ as its corresponding measure. For any $N \in \mathbb{N}^+$, $L \in \mathbb{N}^+$, $p \in [1, \infty)$, $\delta \in (0, 1)$, there exists a ReLU FNN ϕ with width $\max\{8d_\delta\lfloor N^{1/d_\delta} \rfloor + 4d_\delta, 12N + 14\}$ and depth $9L + 12$ such that*

$$(1.3) \quad \|f - \phi\|_{L^p([0, 1]^d, \mu_\varrho)} \leq 3\omega_f\left(\frac{4\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}}\right) + 5\omega_f\left(\frac{16d}{(1-\delta)\sqrt{d_\delta}}N^{-2/d_\delta}L^{-2/d_\delta}\right).$$

The approximation rate of deep neural networks for functions defined precisely on low-dimensional smooth manifolds has been studied in [56] for C^2 functions and in [9, 11] for Lipschitz continuous functions. Considering that it might be more reasonable to assume data located in a small neighborhood of low-dimensional smooth manifold in real applications, we introduce the ε -neighborhood of the manifold \mathcal{M} in Theorem 1.3. In general, existing results are again asymptotic and they cannot be applied to estimate the approximation accuracy of a ReLU FNN with arbitrarily given width N and depth L , since there is no explicit formula without unknown constants to specify the exact error bound. For example, [9] provides an approximation rate $c_1 (NL)^{-c_2/d_\delta}$ with unknown constants (e.g., c_1 and c_2) and requires NL greater than an unknown large number. The demand of an explicit error estimation motivates Theorem 1.3 in this paper. When data are concentrating around \mathcal{M} , ε is very small and the dominant term of the approximation error in (1.3) is $5\omega_f\left(\frac{16d}{(1-\delta)\sqrt{d_\delta}}N^{-2/d_\delta}L^{-2/d_\delta}\right)$ implying that the approximation via deep ReLU FNNs can lessen the curse of dimensionality.

The analysis above provides a general guide for selecting the width and depth of ReLU FNNs to approximate continuous functions, especially when the computation is conducted with parallel computing, which is usually the case in real applications [54, 12]. As we shall see later, when the approximation accuracy and the parallel computing efficiency are considered together, very deep FNNs become less attractive than those with $\mathcal{O}(1)$ depth. Besides, the width requirement $\max\{8d\lfloor N^{1/d} \rfloor + 4d, 12N + 14\}$ in Theorem 1.1 implies that a minimum width of $\mathcal{O}(d)$ is required to maintain a tight approximation rate agreeing with the observation in [39]. In the case of data in an ε -neighborhood of a low-dimensional smooth manifold, the width requirement is reduced to $\max\{8d_\delta\lfloor N^{1/d_\delta} \rfloor + 4d_\delta, 12N + 14\}$.

The approximation theories in this paper assume that the target function f is fully accessible, making it possible to estimate the approximation error and identify an asymptotically optimal ReLU FNN with a given budget of neurons to minimize the approximation error. In real applications, usually only a limited number of possibly noisy observations of f is available, resulting in a regression problem in statistics. In the latter case, the problem is usually formulated in a stochastic setting with randomly generated noisy observations and the regression error contains mainly two components: bias and variance. The bias is the difference of the expectation of an estimated function and its ground truth f . The approximation theories in this paper play an important role in characterizing the power of neural networks when they are applied to solve regression problems by providing a lower bound of the regression bias.

The rest of this paper is organized as follows. A constructive proof of Theorem 1.1 will be shown in Section 2. An asymptotic analysis will be presented in Section 3 to show the tightness of Theorem 1.1. In Section 4, three aspects of neural networks in practice will be discussed: 1) neural network approximation in a high-dimensional irregular domain; 2) neural network approximation in the case of a low-dimensional data structure; 3) the optimal ReLU FNN in parallel computation. Finally, Section 5 concludes this paper with a short discussion.

2. Constructive Approximation Rate $\mathcal{O}(\omega_f(N^{-2/d}L^{-2/d}))$. In this section, we prove the quantitative and constructive approximation rate of ReLU FNNs in Theorem 1.1. Notations throughout the proof will be summarized in Section 2.1. The main ideas of the proof are summarized in Section 2.2. The proof of Theorem 1.1 in the one and multi dimensional cases are presented in Section 2.3 and 2.4, respectively. Finally, several auxiliary lemmas required in the proof of Theorem 3.6 are proved in Section 2.5.

2.1. Notations. Let us summarize all basic notations used in this paper as follows.

- Matrices are denoted by bold uppercase letters. For instance, $\mathbf{A} \in \mathbb{R}^{m \times n}$ is a real matrix of size $m \times n$, and \mathbf{A}^T denotes the transpose of \mathbf{A} .
- Vectors are denoted as bold lowercase letters. For example, $\mathbf{v} \in \mathbb{R}^n$ is a column vector of size n .

Correspondingly, $\mathbf{v}(i)$ is the i -th element of \mathbf{v} . $\mathbf{v} = [v_1, \dots, v_n]^T = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$ are vectors consisting of numbers $\{v_i\}$ with $\mathbf{v}(i) = v_i$.

- Let $\mu(\cdot)$ be the Lebesgue measure.
- The set difference of two sets A and B is denoted by $A \setminus B := \{x : x \in A, x \notin B\}$.

- For any $\xi \in \mathbb{R}$, let $\lfloor \xi \rfloor := \max\{i : i \leq \xi, i \in \mathbb{N}\}$ and $\lceil \xi \rceil := \min\{i : i \geq \xi, i \in \mathbb{N}\}$.
- Assume $\mathbf{n} \in \mathbb{N}^n$, 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 $+\infty$.
- Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ denote the rectified linear unit (ReLU), i.e. $\sigma(x) = \max\{0, x\}$. With the abuse of notations, 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 $\mathbf{x} = [x_1, \dots, x_d]^T \in \mathbb{R}^d$.
- We will use NN as a ReLU neural network for short and use Python-type notations to specify a class of NN's, e.g., $\text{NN}(c_1; c_2; \dots; c_m)$ is a set of ReLU FNNs satisfying m conditions given by $\{c_i\}_{1 \leq i \leq m}$, each of which may specify the number of inputs (#input), the total number of nodes in all hidden layers (#node), the number of hidden layers (#layer), the number of total parameters (#parameter), and the width in each hidden layer (widthvec), the maximum width of all hidden layers (maxwidth), etc. For example, if $\phi \in \text{NN}(\text{\#input} = 2; \text{widthvec} = [100, 100])$, then ϕ satisfies
 - ϕ maps from \mathbb{R}^2 to \mathbb{R} .
 - ϕ has two hidden layers and the number of nodes in each hidden layer is 100.
- $[n]^L$ is short for $[n, n, \dots, n] \in \mathbb{N}^L$. For example,

$$\text{NN}(\text{\#input} = d; \text{widthvec} = [100, 100]) = \text{NN}(\text{\#input} = d; \text{widthvec} = [100]^2).$$

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

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

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

$$\mathbf{h}_i := \mathbf{W}_i \tilde{\mathbf{h}}_{i-1} + \mathbf{b}_i, \quad \text{for } i = 1, \dots, L+1,$$

and

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

- The expression, an FNN with width N and depth L , means
 - The maximum width of this FNN for all hidden layers less than or equal to N .
 - The number of hidden layers of this FNN less than or equal to L .
- For $x \in [0, 1]$, suppose its binary representation is $x = \sum_{\ell=1}^{\infty} x_{\ell} 2^{-\ell}$ with $x_{\ell} \in \{0, 1\}$, we introduce a special notation $\text{bin}0.x_1 x_2 \dots x_L$ to denote the L -term binary representation of x , i.e., $\sum_{\ell=1}^L x_{\ell} 2^{-\ell}$.

2.2. Main ideas. We will show that an almost piecewise constant ReLU FNN ϕ is enough to achieve the desired approximation rate in Theorem 1.1. Given an arbitrary $f \in C([0, 1]^d)$, we introduce a piecewise constant function $f_p \approx f$ serving as an intermediate approximant in our construction in the sense that

$$f \approx f_p \approx \phi.$$

The approximation in $f \approx f_p$ is a simple and standard technique in constructive approximation. For example, given arbitrary N and L , uniformly partition $[0, 1]^d$ into $\mathcal{O}(N^2 L^2)$ pieces and define f_p using this partition. Then the approximation error of $f_p \approx f$ scales like $\mathcal{O}(N^{-2/d} L^{-2/d})$. We will address the approximation in $f_p \approx \phi$ with the same error scaling and a limited budget of the FNN size, e.g., $\mathcal{O}(NL)$ neurons, based on the fact that f_p is almost surely a ReLU FNN in $[0, 1]^d \setminus \mathcal{H}$, where \mathcal{H} is a “don’t-care” region near the discontinuous locations of f_p with an arbitrarily small Lebesgue measure (see Figure 2 for an illustration). The introduction of the “don’t-care” region is to ease the construction of a deep ReLU FNN ϕ , which is a piecewise linear and continuous function, to approximate the discontinuous function f_p by removing the difficulty near discontinuous points, essentially smoothing f_p by restricting the approximation domain in $[0, 1]^d \setminus \mathcal{H}$.

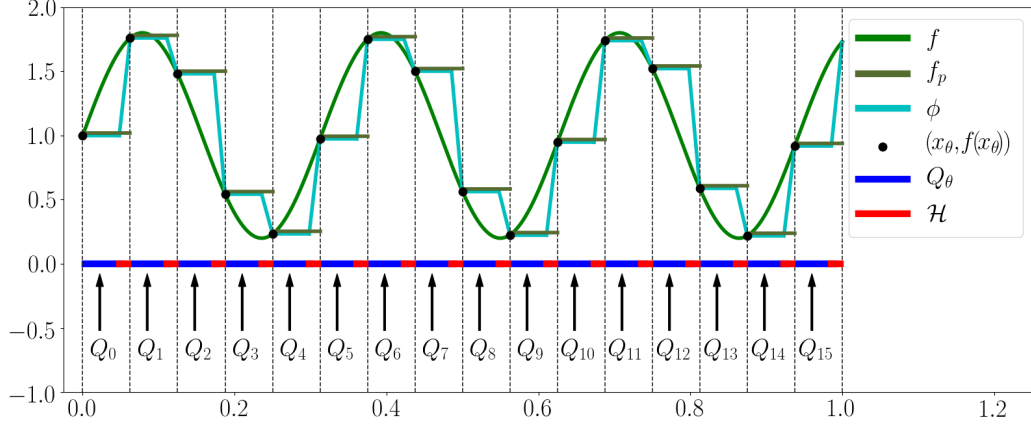


FIG. 2. An illustration of f , f_p , ϕ , x_θ , Q_θ , and the “don’t-care” region \mathcal{H} in the one-dimensional case for $\theta = 0, 1, \dots, N^2L^2 - 1$, where $N = 2$ and $L = 2$. f is the target function; f_p is the piecewise constant function approximating f ; ϕ is the ReLU FNN approximating f ; and x_θ is a representative of Q_θ . The measure of \mathcal{H} can be arbitrarily small as we shall see in the proof of Theorem 1.1.

In the one-dimensional case, a naive way to construct a ReLU FNN ϕ to approximate f_p is to choose ϕ as a piecewise linear function satisfying: 1) matching f_p in $[0, 1] \setminus \mathcal{H}$; 2) having one piece of linear region in each isolated contiguous region of \mathcal{H} (see ϕ and \mathcal{H} in Figure 2 for an illustration). Such a ϕ has $\mathcal{O}(N^2L^2)$ pieces of linear regions. It can be represented as a ReLU FNN with one hidden layer and $\mathcal{O}(N^2L^2)$ neurons, which is far beyond our budget, e.g., $\mathcal{O}(NL)$ neurons. This motivates our design of a deep ReLU FNN via function compositions for the limited budget.

Let us use the one-dimensional case and Figure 2 to illustrate the construction of $\phi(x) \approx f_p(x)$ for $x \in [0, 1] \setminus \mathcal{H}$ as a warm-up. First of all, we simplify the function approximation $\phi(x) \approx f_p(x)$ for $x \in [0, 1] \setminus \mathcal{H}$ in the L^∞ -norm to a regression problem with N^2L^2 samples $\{x_\theta, f_p(x_\theta)\}_{\theta \in \{0, 1, \dots, N^2L^2 - 1\}}$, since in the one-dimensional case f_p has N^2L^2 pieces $\{Q_\theta\}_\theta$ and constant values $\{f_p(x_\theta)\}_\theta$, where x_θ is a point representative in Q_θ , e.g., an end point of Q_θ as in Figure 2, and $\{Q_\theta\}_\theta$ are the intervals separated by \mathcal{H} .

We design a ReLU FNN $\phi = \tilde{\phi}_2 \circ \phi_1$ to approximate f_p , where ϕ_1 and $\tilde{\phi}_2$ are two ReLU FNNs such that $\phi_1(x) = x_\theta$ for $x \in Q_\theta$ and $f_p(x_\theta) = \tilde{\phi}_2(x_\theta)$ for $\theta \in \{0, 1, \dots, N^2L^2 - 1\}$. Then $f_p(x) \approx \phi(x) = \tilde{\phi}_2 \circ \phi_1(x)$ for any $x \in [0, 1] \setminus \mathcal{H}$. In fact, for the convenience of indexing, we construct ϕ_1 and $\tilde{\phi}_2$ such that $\phi_1(x) = \theta$ for $x \in Q_\theta$ and $f_p(x_\theta) = \tilde{\phi}_2(\theta)$ for $\theta \in \{0, 1, \dots, N^2L^2 - 1\}$. See Figure 3 (a) and (b) for an example of ϕ_1 . The simplification reduces the difficulty of controlling the approximation error for all points in $[0, 1] \setminus \mathcal{H}$ to a simpler task of controlling the regression error at N^2L^2 points. The most difficult part of the proof of Theorem 1.1 is to construct such two ReLU FNNs ϕ_1 and $\tilde{\phi}_2$ using $\mathcal{O}(NL)$ neurons.

Similarly in the d -dimensional case, $[0, 1]^d$ is divided into a union of important regions $\{Q_\theta\}$ and a “don’t-care” region \mathcal{H} , where each Q_θ is associated with a representative $\mathbf{x}_\theta \in Q_\theta$ such that $f(\mathbf{x}_\theta) = f_p(\mathbf{x}_\theta)$ for each index vector $\theta \in \{0, 1, \dots, J - 1\}^d$, where $J = \mathcal{O}(N^{2/d}L^{2/d})$ is the partition number per dimension (see Figure 3 (c) for an example when $d = 2$). We design a ReLU FNN $\phi = \tilde{\phi}_2 \circ \phi_1$ to approximate f_p , where ϕ_1 and $\tilde{\phi}_2$ are two ReLU FNNs such that $\phi_1(\mathbf{x}) = \theta$ for $\mathbf{x} \in Q_\theta$ and $f_p(\mathbf{x}_\theta) = \tilde{\phi}_2(\theta)$ for $\theta \in \{0, 1, \dots, J - 1\}^d$. Then $\phi(\mathbf{x}) \approx f_p(\mathbf{x})$ for $\mathbf{x} \in [0, 1]^d \setminus \mathcal{H}$. Here, $\phi_1(\mathbf{x})$ can be constructed via $\phi_1(\mathbf{x}) = [\phi_1(x_1), \phi_1(x_2), \dots, \phi_1(x_d)]^T$, where each one-dimensional ReLU FNN ϕ_1 follows the same construction as in the one-dimensional case.

However, constructing a d -dimensional ReLU FNN $\tilde{\phi}_2$ for the d -dimensional regression problem above is not straightforward, which motivates us to reduce the d -dimensional regression problem to a one-dimensional regression problem via an injective projection map ϕ_2 in a form of ReLU FNN such that it maps $\theta \in \{0, 1, \dots, J - 1\}^d$ to \mathbb{N} . Then we construct a ReLU FNN ϕ_3 such that $\tilde{\phi}_2(\theta) = \phi_3 \circ \phi_2(\theta)$ for $\theta \in \{0, 1, \dots, J - 1\}^d$. The

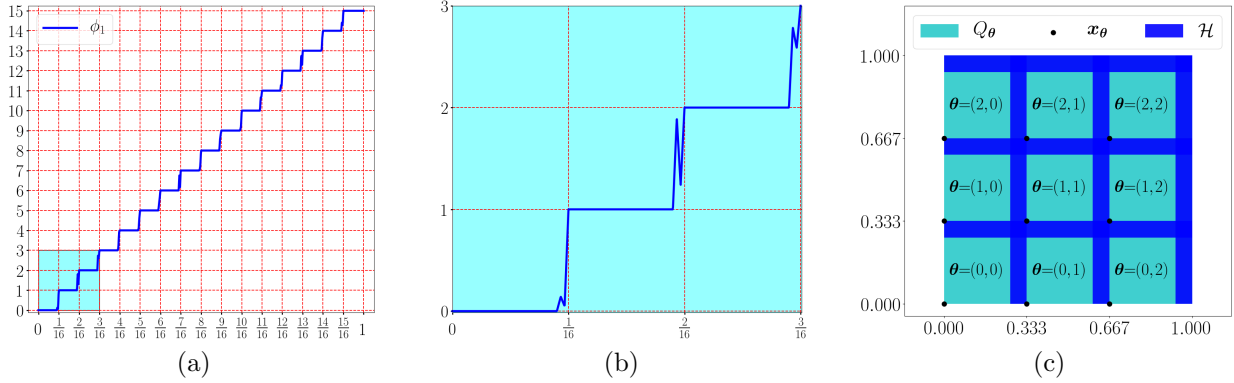


FIG. 3. (a) An illustration of ϕ_1 when $d = 1$, $N = 2$, and $L = 2$. (b) A zoomed-in image of the cyan region in (a). In the “don’t-care” region \mathcal{H} , e.g., around $\frac{1}{16}$ and $\frac{2}{16}$, ϕ_1 may oscillate and it is difficult to control its behavior. (c) An illustration of Q_θ , \mathbf{x}_θ , and the “don’t-care” region \mathcal{H} when $d = 2$ and the partition number per dimension $J = 3$.

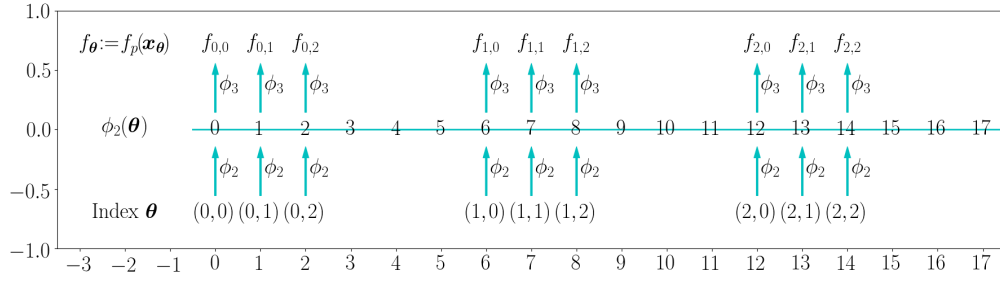


FIG. 4. An illustration of ϕ_2 and ϕ_3 when $d = 2$ and the partition number per dimension $J = 3$. The indices $\{\theta\}$ are corresponding to those in Figure 3 (c). For example, given $\theta = (0, 0)$ and \mathbf{x} in Q_θ in Figure 3 (c), ϕ_1 maps \mathbf{x} to θ , which is mapped to 0 by ϕ_2 as in this figure. Next, 0 will be mapped to $f_\theta := f_p(\mathbf{x}_\theta)$ by ϕ_3 finally.

construction of ϕ_3 is the same as that of $\tilde{\phi}_2$ in the one-dimensional case since they are ReLU FNNs mapping $\mathcal{O}(N^2 L^2)$ integers to given function values. Hence, we can also use the notation $\phi_3 \circ \phi_2(\theta)$ when $d = 1$ via setting ϕ_2 as an identity map. The computation flow of ϕ_1 , ϕ_2 , and ϕ_3 is illustrated in Figure 4, and the overall structure of the ReLU FNN ϕ constructed in a form of $\phi_3 \circ \phi_2 \circ \phi_1$ is visualized in Figure 5.

Finally, we discuss how to construct ϕ_1 , ϕ_2 , and ϕ_3 using deep ReLU FNNs with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ using several propositions as we shall prove in Section 2.5 later. The construction of ϕ_1 is based on Proposition 2.1 and 2.2 below. Directly constructing a deep FNN achieving an approximation goal is challenging, while designing a shallow FNN with one hidden layer for the same approximation task is simpler. This motivates the introduction of Proposition 2.1 and 2.2 to show how to construct a deep FNN to approximate a shallow FNN. As shown in Figure 6, suppose a shallow FNN (shown in Figure 6 (left)) is able to fit given samples in Proposition 2.1, the application of Proposition 2.1 leads to an FNN with two hidden layers (shown in Figure 6 (middle)) approximating this shallow FNN, while the application of Proposition 2.2 results in a deep FNN (shown in Figure 6 (right)) approximating the two-hidden-layer FNN. Hence, the resulting deep FNN can also fit the same given samples as the shallow FNN does. This idea as visualized in Figure 6 will be repeatedly applied in the proof of Theorem 1.1.

PROPOSITION 2.1. *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 $x_0 < x_1 < \dots < x_{N_1(N_2+1)}$ and $y_i \geq 0$ for $i = 0, 1, \dots, N_1(N_2 + 1)$, there exists $\phi \in \text{NN}(\#input = 1; \text{widthvec} = [2N_1, 2N_2 + 1])$ satisfying the following conditions.*

1. $\phi(x_i) = y_i$ for $i = 0, 1, \dots, N_1(N_2 + 1)$;

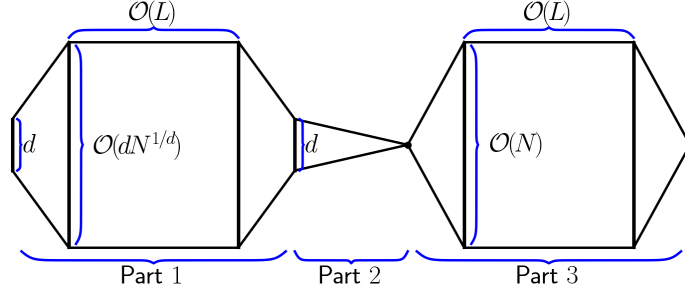


FIG. 5. The structure of the target FNN $\phi = \phi_3 \circ \phi_2 \circ \phi_1$. The first part of ϕ is ϕ_1 , an FNN with width $\mathcal{O}(dN^{1/d})$ and depth $\mathcal{O}(L)$ constructed via Proposition 2.1 and 2.2; the second part of ϕ is ϕ_2 , an FNN with an input layer, an output layer, and no hidden layer; the third part of ϕ is ϕ_3 , an FNN with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ constructed via Proposition 2.3.

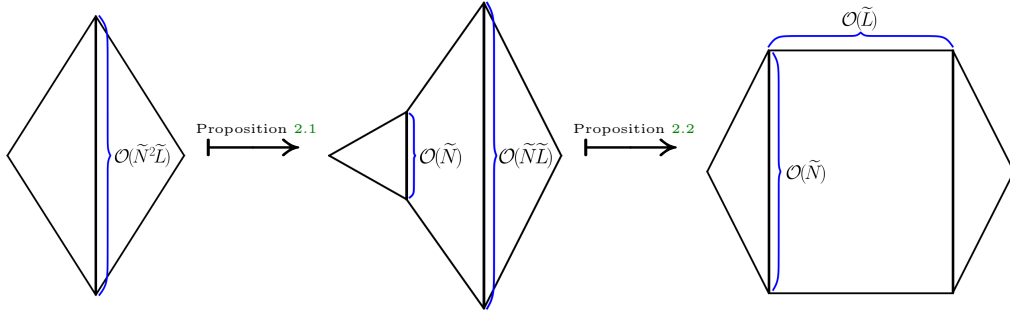


FIG. 6. Given $\mathcal{O}(\tilde{N}^2\tilde{L})$ samples, it is simple to construct a ReLU FNN with one hidden layer of size $\mathcal{O}(\tilde{N}^2\tilde{L})$ fitting these samples. The application of Proposition 2.1 with $N_1 = \mathcal{O}(\tilde{N})$ and $N_2 = \mathcal{O}(\tilde{N}\tilde{L})$ gives a ReLU FNN with two hidden layers fitting the same samples as well. Finally, the application of Proposition 2.2 with $N = \tilde{N}$ and $L = \tilde{L}$ shows the existence of a deep FNN equivalent to the two-hidden-layer FNN.

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

PROPOSITION 2.2. Given any $N, L \in \mathbb{N}^+$, for arbitrary $\phi_1 \in \text{NN}(\#\text{input} = 1; \text{widthvec} = [N, NL])$, there exists $\phi_2 \in \text{NN}(\#\text{input} = 1; \text{maxwidth} \leq 2N + 4; \#\text{layer} \leq L + 2)$ such that $\phi_1(x) = \phi_2(x)$ for any $x \in \mathbb{R}$.

Afterwards, we will use the notation $\text{NN}(\#\text{input} = 1; \text{widthvec} = [N, NL]) \subseteq \text{NN}(\#\text{input} = 1; \text{maxwidth} \leq 2N + 4; \#\text{layer} \leq L + 2)$ in the sense of Proposition 2.1.

The construction of ϕ_2 is inspired by the J -ary representation, i.e., $\phi_2(\theta) = 2J^d \left(\frac{\theta_d}{2J^d} + \sum_{i=1}^{d-1} \frac{\theta_i}{J^i} \right)$ for $\theta \in \{0, 1, \dots, J-1\}^d$, which can be implemented easily via a one-layer ReLU FNN as we shall see in the proof of Theorem 1.1 later.

The construction of ϕ_3 is a direct result of Proposition 2.3 below, the proof of which relies on both the idea illustrated in Figure 6 and the bit extraction technique in [5].

PROPOSITION 2.3. Given any $\varepsilon > 0$, arbitrary $N, L, K \in \mathbb{N}^+$ with $K \leq N^2L^2$, a sample set $\{y_k \geq 0 : k = 0, 1, \dots, K-1\}$ with $|y_k - y_{k-1}| \leq \varepsilon$, for $k = 1, 2, \dots, K-1$, there exists $\phi \in \text{NN}(\#\text{input} = 1; \text{maxwidth} \leq 12N + 14; \#\text{layer} \leq 4L + 10)$ such that

1. $|\phi(k) - y_k| \leq \varepsilon$, for $k = 0, 1, \dots, K-1$;
2. $0 \leq \phi(t) \leq \max\{y_k : k = 0, 1, \dots, K-1\}$, for any $t \in \mathbb{R}$.

Due to the highly nonlinear nature of function compositions, though it is possible to construct ϕ as a deep ReLU FNN approximating f_p uniformly well in $[0, 1]^d \setminus \mathcal{H}$ with the desired approximation rate using the above ideas, the configuration of ϕ might be complicated in \mathcal{H} and we can only bound the L^∞ -norm of ϕ in \mathcal{H} with a constant independent of N and L . In fact, the behavior of $\phi(\mathbf{x})$ when $\mathbf{x} \in \mathcal{H}$ may oscillate, e.g., due to

the oscillation of ϕ_1 as shown in Figure 3 (b). Fortunately, we would like to emphasize that our construction of such a ϕ is valid for \mathcal{H} with an arbitrarily small measure. Hence, it is enough to care about the L^∞ -norm of ϕ in the “don’t-care” region \mathcal{H} , since the approximation is essentially valid in $[0, 1]^d$ in the L^∞ sense and completely valid in the L^p sense for $p \in [1, \infty)$.

With the above propositions ready, let us prove Theorem 1.1 in the case of $d = 1$ in Section 2.3 and the case of $d \geq 2$ in Section 2.4. We further assume that $\omega_f(r) > 0$ for any $r > 0$, excluding a simple case when f is a constant function.

2.3. Proof of Theorem 1.1 for The Case $d = 1$. The key idea for the proof of Theorem 1.1 is to construct an almost piecewise constant ReLU FNN to approximate f within a budget of $\mathcal{O}(NL)$ neurons to achieve the desired approximation power. Note that $|f(x) - f(0)| \leq \omega_f(1)$ for $x \in [0, 1]$. Set $\bar{f} = f - f(0) + \omega_f(1)$, then $0 \leq \bar{f}(x) \leq 2\omega_f(1)$ for $x \in [0, 1]$. Let $M = N^2L$ and $\delta > 0$ be a small number such that

$$(2.1) \quad ML\delta = N^2L^2\delta \leq \min\{1/2, \eta\}.$$

The proof can be divided into four steps as follows:

1. Divide $[0, 1]$ into a union of sub-intervals $\{Q_\theta : \theta = 0, 1, \dots, ML - 1\}$ and a “don’t-care” region \mathcal{H} ;
2. Construct a sub-FNN ϕ_1 to output θ given $x \in Q_\theta$, i.e., $\phi_1(x) = \theta$;
3. Construct a sub-FNN ϕ_2 mapping the index θ approximately to $\bar{f}(\frac{\theta}{ML}) = f(\frac{\theta}{ML}) - f(0) + \omega_f(1)$;
4. Construct the final target FNN $\phi = \phi_2 \circ \phi_1$.

Step 1: Divide $[0, 1]$ into a union of sub-intervals $\{Q_\theta\}_\theta$ and a “don’t-care” region \mathcal{H} .

Define

$$(2.2) \quad \mathcal{H} := \cup_{\theta=1}^{ML} [\frac{\theta}{ML} - \delta, \frac{\theta}{ML}],$$

and

$$Q_\theta := [\frac{\theta}{ML}, \frac{\theta+1}{ML} - \delta], \quad \text{for } \theta = 0, 1, \dots, ML - 1.$$

Apparently, $\mu(\mathcal{H}) \leq ML\delta$ and $[0, 1] = (\cup_{\theta=0}^{ML-1} Q_\theta) \cup \mathcal{H}$. See Figure 2 for an illustration of $\{Q_\theta\}_\theta$ and the “don’t-care” region \mathcal{H} .

Step 2: Construct a sub-FNN ϕ_1 mapping $x \in Q_\theta$ to θ .

There are $ML = N^2L^2$ intervals $\{Q_\theta\}$. It is difficult to directly construct a ReLU FNN with at most $\mathcal{O}(NL)$ neurons mapping $x \in Q_\theta$ to θ for $\theta = 0, 1, \dots, N^2L^2 - 1$. However, there exists a ReLU FNN with at most $\mathcal{O}(NL)$ neurons mapping $x \in Q_\theta$ to θ for N^2L θ 's by Proposition 2.1. Note that the one-dimensional index θ is equivalent to a two dimensional index pair (m, ℓ) via $\theta = mL + \ell$. Hence, we can apply Proposition 2.1 to construct FNNs ψ_1 and ψ_2 with at most $\mathcal{O}(NL)$ neurons such that $\psi_1(x) = m$ and $\psi_2(x) = \ell$, and construct the desired FNN ϕ_1 via the linear combination of ψ_1 and ψ_2 as follows.

For the sample set $\{(\frac{m}{M}, m) : m = 0, 1, \dots, M\} \cup \{(\frac{m+1}{M} - \delta, m) : m = 0, 1, \dots, M - 1\} \subseteq \mathbb{R}^2$, whose cardinality is $2M + 1 = N \cdot ((2NL - 1) + 1) + 1$, by Proposition 2.1, there exist $\psi_1 \in \text{NN}(\#input = 1; \text{widthvec} = [2N, 2(2NL - 1) + 1]) = \text{NN}(\#input = 1; \text{widthvec} = [2N, 4NL - 1])$ such that

- $\psi_1(1) = M$ and $\psi_1(\frac{m}{M}) = \psi_1(\frac{m+1}{M} - \delta) = m$ for $m = 0, 1, \dots, M - 1$;
- ψ_1 is linear on each interval $[\frac{m}{M}, \frac{m+1}{M} - \delta]$ for $m = 0, 1, \dots, M - 1$.

Then

$$(2.3) \quad \psi_1(x) = m, \quad \text{if } x \in Q_\theta \text{ and } \theta = mL + \ell, \quad \text{for } m = 0, 1, \dots, M - 1, \ell = 0, 1, \dots, L - 1.$$

For the sample set $\{(\frac{\ell}{ML}, \ell) : \ell = 0, 1, \dots, L\} \cup \{(\frac{\ell+1}{ML} - \delta, \ell) : \ell = 0, 1, \dots, L - 1\} \subseteq \mathbb{R}^2$, whose cardinality is $2L + 1 = 1 \cdot ((2L - 1) + 1) + 1$, by Proposition 2.1, there exists $\psi_2 \in \text{NN}(\#input = 1; \text{widthvec} = [2, 2(2L - 1) + 1]) = \text{NN}(\#input = 1; \text{widthvec} = [2, 4L - 1])$ such that

- $\psi_2(\frac{1}{M}) = L$ and $\psi_2(\frac{\ell}{ML}) = \psi_2(\frac{\ell+1}{ML} - \delta) = \ell$ for $\ell = 0, 1, \dots, L-1$;
- ψ_2 is linear on each interval $[\frac{\ell}{ML}, \frac{\ell+1}{ML} - \delta]$ for $\ell = 0, 1, \dots, L-1$.

Then

$$(2.4) \quad \psi_2\left(x - \frac{1}{M}\psi_1(x)\right) = \psi_2\left(x - \frac{m}{M}\right) = \ell, \quad \text{if } x \in Q_\theta \text{ and } \theta = mL + \ell, \quad \text{for } m = 0, 1, \dots, M-1, \ell = 0, 1, \dots, L-1.$$

Define

$$\phi_1(x) := L\psi_1(x) + \psi_2\left(x - \frac{1}{M}\psi_1(x)\right), \quad \text{for any } x \in \mathbb{R}.$$

By (2.3) and (2.4), if $x \in Q_\theta$ and $\theta = mL + \ell$ for $m = 0, 1, \dots, M-1$, $\ell = 0, 1, \dots, L-1$, we have

$$\phi_1(x) = L\psi_1(x) + \psi_2\left(x - \frac{1}{M}\psi_1(x)\right) = Lm + \psi_2\left(x - \frac{m}{M}\right) = Lm + \ell = \theta.$$

By Proposition 2.2, $\psi_1 \in \text{NN}(\#\text{input} = 1; \text{widthvec} = [2N, 4NL - 1]) \subseteq \text{NN}(\#\text{input}; \text{maxwidth} \leq 8N + 4; \#\text{layer} \leq L + 2)$ and $\psi_2 \in \text{NN}(\#\text{input} = 1; \text{widthvec} = [2, 4L - 1]) \subseteq \text{NN}(\#\text{input}; \text{maxwidth} \leq 12; \#\text{layer} \leq L + 2)$, then $\phi_1 \in \text{NN}(\#\text{input} = 1; \text{maxwidth} \leq 4N + 10; \#\text{layer} \leq 2L + 4)$.

Step 3: Construct a sub-FNN ϕ_2 mapping the index θ approximately to $\bar{f}(\frac{\theta}{ML})$.

Define

$$y_\theta := \bar{f}\left(\frac{\theta}{ML}\right) \geq 0, \quad \text{for } \theta = 0, 1, \dots, ML - 1,$$

then $|y_\theta - y_{\theta-1}| \leq \omega_f(\frac{1}{ML})$. By Proposition 2.3, there exists $\phi_2 \in \text{NN}(\#\text{input} = 2; \text{maxwidth} \leq 12N + 14; \#\text{layer} \leq 4L + 10)$ such that

$$(2.5) \quad |\phi_2(\theta) - y_\theta| \leq \omega_f\left(\frac{1}{ML}\right), \quad \text{for } \theta = 0, 1, \dots, ML - 1,$$

and

$$(2.6) \quad 0 \leq \phi_2(t) \leq \max\{y_\theta : \theta = 0, 1, \dots, ML - 1\} \leq 2\omega_f(1), \quad \text{for any } t \in \mathbb{R}.$$

Step 4: Construct the final target FNN ϕ using ϕ_1 and ϕ_2 .

Define

$$\bar{\phi} := \phi_2 \circ \phi_1 \quad \text{and} \quad \phi := \bar{\phi} + f(0) - \omega_f(1).$$

Then $\phi = \bar{\phi} + f(0) - \omega_f(1) \in \text{NN}(\#\text{input} = 1; \text{maxwidth} \leq 12N + 14; \#\text{layer} \leq 6L + 14)$ as desired. It suffices to estimate the approximation rate. By Equation (2.6), it holds that

$$(2.7) \quad 0 \leq \bar{\phi}(x) = \phi_2 \circ \phi_1(x) \leq 2\omega_f(1), \quad \text{for any } x \in [0, 1].$$

By Equation (2.3), (2.4), and (2.5), for $x \in Q_\theta$, $\theta = 0, 1, \dots, ML - 1$, we have

$$|\bar{f}(x) - \bar{\phi}(x)| = |\bar{f}(x) - \phi_2 \circ \phi_1(x)| \leq |\bar{f}(x) - y_\theta| + |y_\theta - \phi_2(\theta)| = |\bar{f}(x) - \bar{f}(\frac{\theta}{ML})| + \omega_f(\frac{1}{ML}) \leq 2\omega_f(\frac{1}{ML}).$$

It follows that

$$(2.8) \quad |\phi(x) - f(x)| = |\bar{f}(x) - \bar{\phi}(x)| \leq 2\omega_f(\frac{1}{ML}) = 2\omega_f(N^{-2}L^{-2}), \quad \text{for any } x \in [0, 1] \setminus \mathcal{H}.$$

By (2.1), (2.2), and (2.7), we get $\mu(\mathcal{H}) \leq ML\delta \leq \eta$ and $\|\phi\|_{L^\infty([0,1])} \leq |f(0)| + \omega_f(1)$. So we finish the proof for the case $d = 1$.

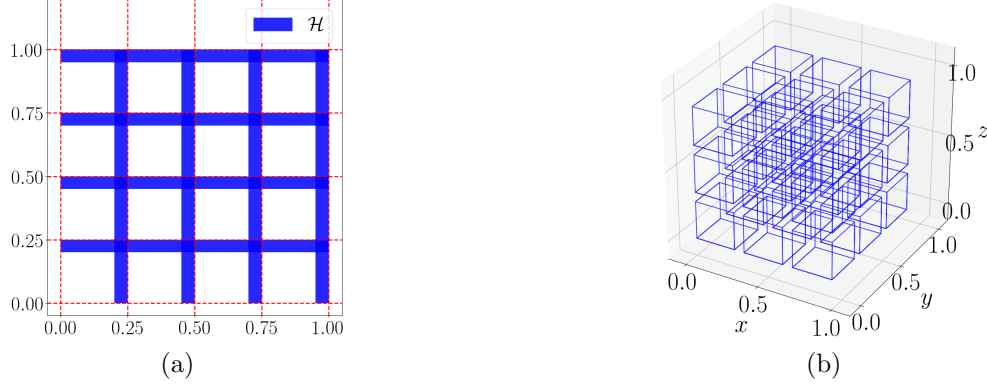


FIG. 7. (a) An illustration of the “don’t-care” region \mathcal{H} for $J = 4$ and $d = 2$. (b) An illustration of Q_{θ} sub-cubes depicted in blue for $\theta \in \{0, 1, \dots, J-1\}^d$, where $J = 3$ and $d = 3$.

2.4. Proof of Theorem 1.1 for The Case $d \geq 2$. We also need to construct an almost piecewise constant ReLU FNN with $\mathcal{O}(NL)$ neurons to approximate f in the d -dimensional case. The main difficulty is to design a ReLU FNN to reduce the d -dimensional approximation problem to a one-dimensional approximation problem as introduced in Section 2.2. It is clear that $|f(\mathbf{x}) - f(0)| \leq \omega_f(\sqrt{d})$ for any $\mathbf{x} \in [0, 1]^d$. Define $\bar{f} = f - f(0) + \omega_f(\sqrt{d})$, then $0 \leq \bar{f}(\mathbf{x}) \leq 2\omega_f(\sqrt{d})$ for any $\mathbf{x} \in [0, 1]^d$. Let $M = N^2L$, $J = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor$, and $\delta = \delta(N, L, d, f) > 0$ be a sufficiently small number satisfying

$$(2.9) \quad Jd\delta = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor d\delta \leq \min\{1, \eta\}.$$

The proof can be divided into six steps as follows:

1. Divide $[0, 1]^d$ into a union of sub-cubes $\{Q_{\theta}\}_{\theta \in \{0, 1, \dots, J-1\}^d}$ and a “don’t-care” region \mathcal{H} ;
2. Construct a sub-FNN ϕ_1 mapping \mathbf{x} to the d -dimensional index θ if $\mathbf{x} \in Q_{\theta}$;
3. Construct a sub-FNN ϕ_2 mapping the index set $\{0, 1, \dots, J-1\}^d$ to an auxiliary set $\mathcal{A}_1 \subseteq \{\frac{k}{2J^d} : k = 0, 1, \dots, 2J^d\}$;
4. Find a continuous piecewise linear function g such that \bar{f} and $g \circ \phi_2 \circ \phi_1$ have the same value at the elements of $\{0, 1, \dots, J-1\}^d/J$, i.e. $\bar{f} \approx g \circ \phi_2 \circ \phi_1$;
5. Construct a sub-FNN ϕ_3 mapping \mathcal{A}_1 to $\{f(\mathbf{x}) : \mathbf{x} \in \{0, 1, \dots, J-1\}^d/J\}$ such that $g \approx \phi_3$ on \mathcal{A}_1 ;
6. Construct the final target FNN $\phi = \phi_3 \circ \phi_2 \circ \phi_1 + f(0) - \omega_f(\sqrt{d}) \approx \bar{f} + f(0) - \omega_f(\sqrt{d}) = f$.

Step 1: Divide $[0, 1]^d$ into a union of sub-cubes $\{Q_{\theta}\}_{\theta \in \{0, 1, \dots, J-1\}^d}$ and a “don’t-care” region \mathcal{H} .

Define

$$\mathcal{H} = \cup_{j=1}^d \left\{ \mathbf{x} = [x_1, x_2, \dots, x_d]^T \in [0, 1]^d : x_j \in \cup_{i=1}^J \left[\frac{i}{J} - \delta, \frac{i}{J} \right] \right\},$$

which is the so called “don’t-care” region and it separates the d -dimensional cube into J^d “important” sub-cubes. To index these d -dimensional smaller sub-cubes, define

$$Q_{\theta} = \{ \mathbf{x} = [x_1, x_2, \dots, x_d]^T \in [0, 1]^d : x_i \in \left[\frac{\theta_i}{J}, \frac{\theta_i+1}{J} - \delta \right], i = 1, 2, \dots, d \}$$

for each d -dimensional index vector $\theta = [\theta_1, \theta_2, \dots, \theta_d]^T \in \{0, 1, \dots, J-1\}^d$. Apparently, $\mu(\mathcal{H}) \leq Jd\delta$, $[0, 1]^d = \cup_{\theta \in \{0, 1, \dots, J-1\}^d} Q_{\theta} \cup \mathcal{H}$, and the intersection of \mathcal{H} and Q_{θ} belongs to the boundary of Q_{θ} (see Figure 7 for illustrations).

Step 2: Construct a sub-FNN ϕ_1 mapping \mathbf{x} to the d -dimensional index θ if $\mathbf{x} \in Q_{\theta}$.

For the sample set $\{(\frac{j}{J}, j) : j = 0, 1, \dots, J\} \cup \{(\frac{j+1}{J} - \delta, j) : j = 0, 1, \dots, J-1\} \subseteq \mathbb{R}^2$, whose cardinality is $2J+1 = \lfloor N^{1/d} \rfloor \cdot ((2\lfloor N^{1/d} \rfloor \lfloor L^{2/d} \rfloor - 1) + 1) + 1$, by Proposition 2.1, there exists $\phi_1 \in \text{NN}(\#input = 1; \text{widthvec} = [2\lfloor N^{1/d} \rfloor, 2(2\lfloor N^{1/d} \rfloor \lfloor L^{2/d} \rfloor - 1) + 1]) = \text{NN}(\#input = 1; \text{widthvec} = [2\lfloor N^{1/d} \rfloor, 4\lfloor N^{1/d} \rfloor \lfloor L^{2/d} \rfloor - 1])$ such that

- $\phi_1(1) = J$, and $\phi_1(\frac{j}{J}) = \phi_1(\frac{j+1}{J} - \delta) = j$ for $j = 0, 1, \dots, J-1$;
- ϕ_1 is linear on each interval $[\frac{j}{J}, \frac{j+1}{J} - \delta]$ for $j = 0, 1, \dots, J-1$.

Then $\boldsymbol{\theta} = [\phi_1(x_1), \phi_1(x_2), \dots, \phi_1(x_d)]^T$ for any $\mathbf{x} = [x_1, x_2, \dots, x_d]^T \in Q_{\boldsymbol{\theta}}$. Define

$$\phi_1(\mathbf{x}) := [\phi_1(x_1), \phi_1(x_2), \dots, \phi_1(x_d)]^T, \quad \text{for any } \mathbf{x} \in \mathbb{R}^d.$$

Hence, $\phi_1(\mathbf{x}) = \boldsymbol{\theta}$ if $\mathbf{x} \in Q_{\boldsymbol{\theta}}$ for $\boldsymbol{\theta} \in \{0, 1, \dots, J-1\}^d$.

Step 3: Construct a sub-FNN ϕ_2 mapping the index set $\{0, 1, \dots, J-1\}^d$ to an auxiliary set $\mathcal{A}_1 \subseteq \{\frac{k}{2J^d} : k = 0, 1, \dots, 2J^d\}$.

Inspired by the binary representation, we define

$$(2.10) \quad \phi_2(\boldsymbol{\theta}) = \phi_2(\theta_1, \theta_2, \dots, \theta_d) := \frac{\theta_d}{2J^d} + \sum_{i=1}^{d-1} \frac{\theta_i}{J^i},$$

then ϕ_2 is an FNN mapping the index set $\{0, 1, \dots, J-1\}^d$ to

$$\mathcal{A}_1 := \left\{ \frac{\theta_d}{2J^d} + \sum_{i=1}^{d-1} \frac{\theta_i}{J^i} : \boldsymbol{\theta} \in \{0, 1, \dots, J-1\}^d \right\} \subseteq \left\{ \frac{k}{2J^d} : k = 0, 1, \dots, 2J^d \right\}.$$

Step 4: Construct a continuous piecewise linear function g such that \bar{f} and $g \circ \phi_2 \circ \phi_1$ have the same value at the elements of $\{0, 1, \dots, J-1\}^d/J$.

Define another auxiliary set

$$\mathcal{A}_2 := \left\{ \frac{1}{2J^{d-1}} + \frac{\theta_d}{2J^d} + \sum_{i=1}^{d-1} \frac{\theta_i}{J^i} : \boldsymbol{\theta} \in \{0, 1, \dots, J-1\}^d \right\},$$

then $\mathcal{A}_1 \cup \mathcal{A}_2 \cup \{1\} = \{\frac{k}{2J^d} : k = 0, 1, \dots, 2J^d\}$.

Let $g : [0, 1] \rightarrow \mathbb{R}$ be a continuous piecewise linear function with break points $\{\frac{k}{2J^d} : k = 0, 1, \dots, 2J^d\}$ and the values of g at these break points satisfy the following properties:

- At the break point 1, let $g(1) = f(\frac{[J-1]^d}{J})$, where $[J-1]^d = [J-1, J-1, \dots, J-1]^T \in \{0, 1, \dots, J-1\}^d$;
- The values of g at the break points in \mathcal{A}_1 are set as

$$(2.11) \quad g\left(\frac{\theta_d}{2J^d} + \sum_{i=1}^{d-1} \frac{\theta_i}{J^i}\right) = \bar{f}\left(\frac{\boldsymbol{\theta}}{J}\right), \quad \text{for any } \boldsymbol{\theta} \in \{0, 1, \dots, J-1\}^d;$$

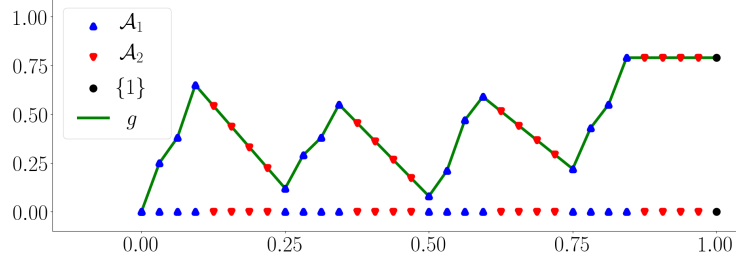
- The values of g at the break points in \mathcal{A}_2 are assigned to reduce the variation of g . To this end, for each $\boldsymbol{\theta} \in \{0, 1, \dots, J-1\}^d$, let g be linear on $[\xi_1^{\boldsymbol{\theta}}, \xi_3^{\boldsymbol{\theta}}]$, where $\xi_1^{\boldsymbol{\theta}} = \frac{J-1}{2J^d} + \sum_{i=1}^{d-1} \frac{\theta_i}{J^i} \in \mathcal{A}_1$ and $\xi_3^{\boldsymbol{\theta}} = \frac{1}{J^{d-1}} + \sum_{i=1}^{d-1} \frac{\theta_i}{J^i} \in \mathcal{A}_1 \cup \{1\}$, which requires that $g(\xi_2^{\boldsymbol{\theta}}) = \frac{(\xi_2^{\boldsymbol{\theta}} - \xi_1^{\boldsymbol{\theta}})g(\xi_3^{\boldsymbol{\theta}}) + (\xi_3^{\boldsymbol{\theta}} - \xi_2^{\boldsymbol{\theta}})g(\xi_1^{\boldsymbol{\theta}})}{\xi_3^{\boldsymbol{\theta}} - \xi_1^{\boldsymbol{\theta}}}$, where $\xi_2^{\boldsymbol{\theta}} = \frac{1}{2J^{d-1}} + \frac{\theta_d}{2J^d} + \sum_{i=1}^{d-1} \frac{\theta_i}{J^i} \in \mathcal{A}_2$.

Apparently, such a function g exists (see Figure 8 for an example) and satisfies

$$\left| g\left(\frac{k}{2J^d}\right) - g\left(\frac{k-1}{2J^d}\right) \right| \leq \max \left\{ \omega_f\left(\frac{1}{J}\right), \omega_f(\sqrt{d})/J \right\} \leq \omega_f\left(\frac{\sqrt{d}}{J}\right), \quad \text{for } k = 1, 2, \dots, 2J^d,$$

and

$$0 \leq g\left(\frac{k}{2J^d}\right) \leq 2\omega_f(\sqrt{d}), \quad \text{for } k = 0, 1, \dots, 2J^d.$$

FIG. 8. An illustration of \mathcal{A}_1 , \mathcal{A}_2 , $\{1\}$, and g for $d = 2$ and $J = 4$.

Step 5: Construct a sub-FNN ϕ_3 mapping \mathcal{A}_1 to $\{\bar{f}(x) : x \in \{0, 1, \dots, J-1\}^d/J\}$.

Since $2J^d = 2(\lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor)^d \leq 2(N^2 L^2) \leq N^2 \tilde{L}^2$, where $\tilde{L} = 2L$, by Proposition 2.3, there exists $\psi \in \text{NN}(\#input = 1; \maxwidth \leq 12N + 14; \#layer \leq 4\tilde{L} + 10) = \text{NN}(\#input = 1; \maxwidth \leq 12N + 14; \#layer \leq 8L + 10)$ such that

$$|\psi(k) - g(\frac{k}{2J^d})| \leq \omega_f(\frac{\sqrt{d}}{J}), \quad \text{for } k = 0, 1, \dots, 2J^d - 1,$$

and

$$0 \leq \psi(t) \leq \max\{g(\frac{k}{2J^d}) : k = 0, 1, \dots, 2J^d - 1\} \leq 2\omega_f(\sqrt{d}), \quad \text{for any } t \in \mathbb{R}.$$

Define $\phi_3(t) := \psi(2J^d t)$ for any $t \in \mathbb{R}$. It follows that

$$(2.12) \quad |\phi_3(\frac{k}{2J^d}) - g(\frac{k}{2J^d})| = |\psi(k) - g(\frac{k}{2J^d})| \leq \omega_f(\frac{\sqrt{d}}{J}), \quad \text{for } k = 0, 1, \dots, 2J^d - 1,$$

and

$$(2.13) \quad 0 \leq \phi_3(t) = \psi(2J^d t) \leq 2\omega_f(\sqrt{d}), \quad \text{for any } t \in \mathbb{R}.$$

Step 6: Construct the final target FNN $\phi = \phi_3 \circ \phi_2 \circ \phi_1 + f(0) - \omega_f(\sqrt{d})$.

Define $\phi := \bar{\phi} + f(0) - \omega_f(\sqrt{d})$, where $\bar{\phi} := \phi_3 \circ \phi_2 \circ \phi_1$. By Proposition 2.2, ϕ_1 constructed in Step 2 is in $\text{NN}(\#input = 1; \text{widthvec} = [2\lfloor N^{1/d} \rfloor, 4\lfloor N^{1/d} \rfloor \lfloor L^{2/d} \rfloor - 1]) \subseteq \text{NN}(\#input = 1; \maxwidth \leq 8\lfloor N^{1/d} \rfloor + 4; \#layer \leq \lfloor L^{2/d} \rfloor + 2)$. Hence, for ϕ_1 constructed in Step 2 and ϕ_2 constructed in Step 3, we have $\phi_2 \circ \phi_1 \in \text{NN}(\#input = d; \maxwidth \leq 8d\lfloor N^{1/d} \rfloor + 4d; \#layer \leq \lfloor L^{2/d} \rfloor + 2)$. It follows that $\phi = \phi_3 \circ \phi_2 \circ \phi_1 + f(0) - \omega_f(\sqrt{d}) \in \text{NN}(\#input = d; \maxwidth \leq \max\{8d\lfloor N^{1/d} \rfloor + 4d, 12N + 14\}; \#layer \leq 9L + 12)$ as desired.

Now let us estimate the approximation error. By (2.11), (2.12), and the definition of ϕ_1 , for any $\mathbf{x} = [x_1, x_2, \dots, x_d]^T \in Q_\theta$, we have

$$|\bar{f}(\mathbf{x}) - \bar{\phi}(\mathbf{x})| = |\bar{f}(\mathbf{x}) - \phi_3(\frac{1}{2J^d} \phi_1(x_d) + \sum_{i=1}^{d-1} \frac{1}{J^i} \phi_1(x_i))| = |\bar{f}(\mathbf{x}) - \phi_3(\frac{\theta_d}{2J^d} + \sum_{i=1}^{d-1} \frac{\theta_i}{J^i})|,$$

which is bounded by

$$|\bar{f}(\mathbf{x}) - g(\frac{\theta_d}{2J^d} + \sum_{i=1}^{d-1} \frac{\theta_i}{J^i})| + |g(\frac{\theta_d}{2J^d} + \sum_{i=1}^{d-1} \frac{\theta_i}{J^i}) - \phi_3(\frac{\theta_d}{2J^d} + \sum_{i=1}^{d-1} \frac{\theta_i}{J^i})| \leq |\bar{f}(\mathbf{x}) - \bar{f}(\frac{\theta}{J})| + \omega_f(\frac{\sqrt{d}}{J}) = 2\omega_f(\frac{\sqrt{d}}{J}).$$

In sum,

$$(2.14) \quad |f(\mathbf{x}) - \phi(\mathbf{x})| = |\bar{f}(\mathbf{x}) - \bar{\phi}(\mathbf{x})| \leq 2\omega_f(\frac{\sqrt{d}}{J}) \leq 2\omega_f(8\sqrt{d}N^{-2/d}L^{-2/d}), \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \mathcal{H},$$

where the last inequality comes from the fact $J = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor \geq \frac{N^{2/d} L^{2/d}}{8}$ for any $N, L \in \mathbb{N}^+$. By (2.9), $\mu(\mathcal{H}) \leq Jd\delta \leq \eta$. It remains to show the upper bound of ϕ . By (2.13) and the definition of $\bar{\phi}$, it holds that $0 \leq \bar{\phi}(\mathbf{x}) \leq 2\omega_f(\sqrt{d})$, for any $\mathbf{x} \in \mathbb{R}^d$. Together with $\phi = \bar{\phi} + f(0) - \omega_f(\sqrt{d})$, it follows that $\|\phi\|_{L^\infty([0,1]^d)} \leq |f(0)| + \omega_f(\sqrt{d})$. Thus, we finish the proof for the case $d \geq 2$.

2.5. Proofs of Propositions in Section 2.2. This section proves Propositions 2.1 to 2.3 in Section 2.2. In fact, Proposition 2.1 is a part of Lemma 2.2 in [57]. For the purpose of being self-contained, we present it as follows.

LEMMA 2.4 (Lemma 2.2 of [57]). *For any $m, n \in \mathbb{N}^+$, given any $m(n+1) + 1$ samples $(x_i, y_i) \in \mathbb{R}^2$ with $x_0 < x_1 < x_2 < \dots < x_{m(n+1)}$ and $y_i \geq 0$ for $i = 0, 1, \dots, m(n+1)$, there exists $\phi \in \text{NN}(\#input = 1; \text{widthvec} = [2m, 2n+1])$ satisfying the following conditions.*

1. $\phi(x_i) = y_i$ for $i = 0, 1, \dots, m(n+1)$;
2. ϕ is linear on each interval $[x_{i-1}, x_i]$ for $i \in \{(n+1)j : j = 1, 2, \dots, m\}$;
3. $\sup_{x \in [x_0, x_{m(n+1)}]} |\phi(x)| \leq 3 \max_{i \in \{0, 1, \dots, m(n+1)\}} y_i \prod_{k=1}^n \left(1 + \frac{\max\{x_{j(n+1)+n} - x_{j(n+1)+k-1} : j=0, 1, \dots, m-1\}}{\min\{x_{j(n+1)+k} - x_{j(n+1)+k-1} : j=0, 1, \dots, m-1\}} \right)$.

The key idea to prove Proposition 2.2 is to re-assemble $\mathcal{O}(L)$ sub-FNNs with width $\mathcal{O}(N)$ and depth $\mathcal{O}(1)$ in the shallower FNN in the middle of Figure 6 to form a deeper one with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ on the right of Figure 6.

Proof of Proposition 2.2. For any $\phi \in \text{NN}(\#input = 1; \text{widthvec} = [N, NL])$, the architecture of ϕ can be described as

$$x \xrightarrow{\mathbf{W}_1, \mathbf{b}_1} \mathbf{g} \xrightarrow{\mathbf{W}_2, \mathbf{b}_2} \mathbf{h} \xrightarrow{\mathbf{W}_3, \mathbf{b}_3} \phi(x),$$

where \mathbf{g} and \mathbf{h} are the output of the first hidden layer and the second hidden layer, respectively. Note that

$$(2.15) \quad \mathbf{g} = \sigma(\mathbf{W}_1 x + \mathbf{b}_1), \quad \mathbf{h} = \sigma(\mathbf{W}_2 \mathbf{g} + \mathbf{b}_2), \quad \text{and} \quad \phi(x) = \mathbf{W}_3 \mathbf{h} + \mathbf{b}_3.$$

We can evenly divide \mathbf{h} , \mathbf{b}_2 , \mathbf{W}_2 , and \mathbf{W}_3 into L parts as follows:

$$\begin{bmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \\ \vdots \\ \mathbf{h}_L \end{bmatrix}, \quad \begin{bmatrix} \mathbf{b}_{2,1} \\ \mathbf{b}_{2,2} \\ \vdots \\ \mathbf{b}_{2,L} \end{bmatrix}, \quad \begin{bmatrix} \mathbf{W}_{2,1} \\ \mathbf{W}_{2,2} \\ \vdots \\ \mathbf{W}_{2,L} \end{bmatrix},$$

and $[\mathbf{W}_{3,1}, \mathbf{W}_{3,2}, \dots, \mathbf{W}_{3,L}]$, where $\mathbf{h}_\ell, \mathbf{b}_{2,\ell} \in \mathbb{R}^N$, $\mathbf{W}_{2,\ell} \in \mathbb{R}^{N \times N}$, and $\mathbf{W}_{3,\ell} \in \mathbb{R}^{1 \times N}$ for $\ell = 1, 2, \dots, L$. Define $y_\ell := \mathbf{W}_{3,\ell} \mathbf{h}_\ell$ and $s_\ell := \sum_{j=1}^\ell y_j$ for $\ell = 1, 2, \dots, L$. Note that $\mathbf{g} \geq 0$ since it is the output of the ReLU function. Hence, it is easy to check that the desired deep FNN can be constructed as follows:

$$x \rightarrow \mathbf{g} \rightarrow \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{g} \end{bmatrix} \rightarrow \begin{bmatrix} \sigma(y_1) \\ \sigma(-y_1) \\ \mathbf{h}_2 \\ \mathbf{g} \end{bmatrix} \rightarrow \begin{bmatrix} \sigma(s_1) \\ \sigma(-s_1) \\ \sigma(y_2) \\ \sigma(-y_2) \\ \mathbf{h}_3 \\ \mathbf{g} \end{bmatrix} \rightarrow \dots \rightarrow \begin{bmatrix} \sigma(s_{L-2}) \\ \sigma(-s_{L-2}) \\ \sigma(y_{L-1}) \\ \sigma(-y_{L-1}) \\ \mathbf{h}_L \\ \mathbf{g} \end{bmatrix} \rightarrow \begin{bmatrix} \sigma(s_{L-1}) \\ \sigma(-s_{L-1}) \\ \sigma(y_L) \\ \sigma(-y_L) \end{bmatrix} \rightarrow \phi(x),$$

where “ \rightarrow ” represents the composition of a ReLU activation function and an appropriate linear transform with weights and bias from the transforms in (2.15) up to the change of their signs. It is clear that the FNN with above architecture has width $2N + 4$ and depth $L + 2$. So, we finish the proof. \square

The proof of Proposition 2.3 is based on the bit extraction technique in [5, 25]. In fact, we modify this technique to extract the sum of many bits rather than one bit and this modification can be summarized in Lemma 2.5 and 2.6 below.

LEMMA 2.5. *For any $L \in \mathbb{N}^+$, there exists $\phi \in \text{NN}(\#input = 2; \text{maxwidth} \leq 7; \#layer \leq 2L + 1)$ such that for any $x = \text{bin } 0.x_1x_2\cdots x_L$, we have $\phi(x, \ell) = \sum_{j=1}^{\ell} x_j$ for $\ell = 1, 2, \dots, L$.*

Proof of Lemma 2.5. For any $x = \text{bin } 0.x_1x_2\cdots x_L$, we define $\xi_j := \text{bin } 0.x_jx_{j+1}\cdots x_L$ for $j = 1, 2, \dots, L$, and

$$\mathcal{T}(t) := \begin{cases} 1, & t \geq 0, \\ 0, & t < 0. \end{cases}$$

Then $x_j = \mathcal{T}(\xi_j - 1/2)$ for $j = 1, 2, \dots, L$, and $\xi_{j+1} = 2\xi_j - x_j$ for $j = 1, 2, \dots, L-1$.

Next, ReLU FNNs will be constructed to represent $\mathcal{T}(\cdot)$ and identity maps in the above discussion. In fact, $\mathcal{T}(t) = \sigma(t/\delta + 1) - \sigma(t/\delta)$ for any $t \notin (-\delta, 0)$. Hence, if we set $\delta = 1/2 - \sum_{j=2}^L 2^{-j} = 2^{-L}$, then

$$(2.16) \quad x_j = \mathcal{T}(\xi_j - 1/2) = \sigma((\xi_j - 1/2)/\delta + 1) - \sigma((\xi_j - 1/2)/\delta)$$

for $j = 1, 2, \dots, L$. Note that $2\xi_j - x_j$ is positive for $j = 1, 2, \dots, L-1$. Hence,

$$(2.17) \quad \xi_{j+1} = 2\xi_j - x_j = \sigma(2\xi_j - x_j), \quad \text{for } j = 1, 2, \dots, L-1.$$

Now let us establish a formula to represent $\sum_{j=1}^{\ell} x_j$ in a form of a ReLU FNN as follows:

$$\sum_{j=1}^{\ell} x_j = \sum_{j=1}^L \mathcal{T}(\ell - j)x_j = \sum_{j=1}^L (\sigma(\ell - j + 1) - \sigma(\ell - j))x_j, \quad \text{for } \ell = 1, 2, \dots, L,$$

where the last equality comes from the fact $\mathcal{T}(n) = \sigma(n + 1) - \sigma(n)$ for any integer n . The fact that $t_1 t_2 = \sigma(t_1 + t_2 - 1)$ for any $t_1, t_2 \in \{0, 1\}$ implies

$$(2.18) \quad \sum_{j=1}^{\ell} x_j = \sum_{j=1}^L (\sigma(\ell - j + 1) - \sigma(\ell - j))x_j = \sum_{j=1}^L \sigma(\sigma(\ell - j + 1) - \sigma(\ell - j) + x_j - 1).$$

Define

$$(2.19) \quad a_j := \sigma(\sigma(\ell - j + 1) - \sigma(\ell - j) + x_j - 1), \quad \text{and} \quad s_j := \sum_{i=1}^j a_i, \quad \text{for } j = 1, 2, \dots, L.$$

By the definitions of a_j , s_j , x_j , and ξ_j in Equation (2.16), (2.17), (2.18), and (2.19), it is easy to construct a ReLU FNN in $\text{NN}(\#input = 2; \text{maxwidth} \leq 7; \#layer \leq 2L + 1)$ outputting $\sum_{j=1}^{\ell} x_j = \sum_{j=1}^L a_j = s_L$ given the input (x, ℓ) with $x = \text{bin } 0.x_1x_2\cdots x_L$ and $\ell \in \{1, 2, \dots, L\}$. The skeleton of the desired FNN is shown in Figure 9. Hence, we finish the proof. \square

Next, we introduce Lemma 2.6 as an advanced version of Lemma 2.5.

LEMMA 2.6. *For any $N, L \in \mathbb{N}^+$, any $\theta_{m,\ell} \in \{0, 1\}$ for $m = 0, 1, \dots, M-1$, $\ell = 0, 1, \dots, L-1$, where $M = N^2L$, there exists $\phi \in \text{NN}(\#input = 2; \text{maxwidth} \leq 4N + 5; \#layer \leq 3L + 4)$ such that $\phi(m, \ell) = \sum_{j=0}^{\ell} \theta_{m,j}$, for $m = 0, 1, \dots, M-1$, $\ell = 0, 1, \dots, L-1$.*

Proof of Lemma 2.6. Define $y_m := \text{bin } 0.\theta_{m,0}\theta_{m,1}\cdots\theta_{m,L-1}$ for $m = 0, 1, \dots, M-1$. By Proposition 2.1, for the sample set $\{(m, y_m) : m = 0, 1, \dots, M\}$, whose cardinality is $M + 1 = N((NL - 1) + 1) + 1$, there exists $\phi_1 \in \text{NN}(\#input = 1; \text{widthvec} = [2N, 2(NL - 1) + 1]) = \text{NN}(\#input = 1; \text{widthvec} = [2N, 2NL - 1])$ such that $\phi_1(m) = y_m$, for $m = 0, 1, \dots, M-1$. By Lemma 2.5, there exists $\phi_2 \in \text{NN}(\#input = 2; \text{maxwidth} \leq 7; \#layer \leq 2L + 1)$ such that $\phi_2(x, \ell) = \sum_{j=1}^{\ell} x_j$ for $\ell = 1, 2, \dots, L$, if $x = \text{bin } 0.x_1x_2\cdots x_L$.

Note that $\phi_2(\phi_1(m), \ell + 1) = \phi_2(y_m, \ell + 1) = \sum_{j=0}^{\ell} \theta_{m,j}$, for $m = 0, 1, \dots, M-1$, $\ell = 0, 1, \dots, L-1$. Define

$$\phi(t_1, t_2) := \phi_2(\sigma(\phi_1(t_1)), \sigma(t_2 + 1)), \quad \text{for any } t_1, t_2 \in \mathbb{R}.$$

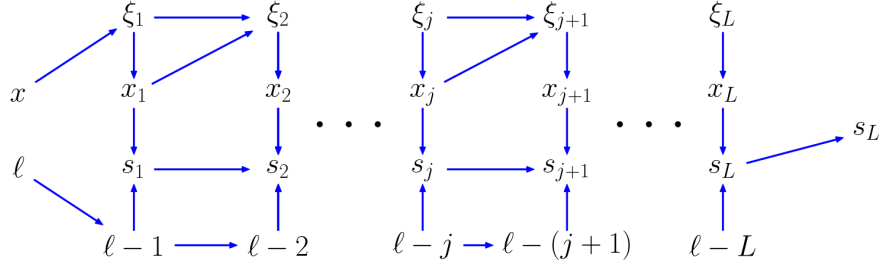


FIG. 9. A skeleton of the target ReLU FNN to illustrate how to output $\sum_{j=1}^L x_j = \sum_{j=1}^L a_j = s_L$ given the input (x, ℓ) with $x = \text{bin } 0.x_1x_2\cdots x_L$ and $\ell \in \{1, 2, \dots, L\}$, following the iterative definitions of a_j , s_j , x_j , and ξ_j in Equation (2.16), (2.17), (2.18), and (2.19).

Then $\phi(m, \ell) = \phi_2(\sigma(\phi_1(m)), \sigma(\ell + 1)) = \phi_2(\phi_1(m), \ell + 1) = \phi_2(y_m, \ell + 1) = \sum_{j=0}^{\ell} \theta_{m,j}$, for $m = 0, 1, \dots, M - 1$, $\ell = 0, 1, \dots, L - 1$. And the architecture of ϕ is

$$\begin{bmatrix} t_1 \\ t_2 \end{bmatrix} \rightarrow \begin{bmatrix} \sigma(\phi_1(t_1)) \\ \sigma(t_2 + 1) \end{bmatrix} \rightarrow \phi_2(\sigma(\phi_1(t_1)), \sigma(t_2 + 1)) = \phi(t_1, t_2).$$

By Proposition 2.2, $\phi_1 \in \text{NN}(\#input = 1; \text{widthvec} = [2N, 2NL - 1]) \subseteq \text{NN}(\#input = 1; \text{maxwidth} \leq 4N + 4; \#layer \leq L + 2)$ implies $\phi \in \text{NN}(\#input = 2; \text{maxwidth} \leq 4N + 5; \#layer \leq 3L + 4)$ as we desire. Therefore, we have finished the proof. \square

Next, we apply Lemma 2.6 to prove Lemma 2.7 below, which is a key intermediate conclusion to prove Proposition 2.3.

LEMMA 2.7. For any $\varepsilon > 0$, $L, N \in \mathbb{N}^+$, any sample set $\{y_{m,\ell} \geq 0 : m = 0, 1, \dots, M - 1, \ell = 0, 1, \dots, L - 1\}$ with $|y_{m,\ell} - y_{m,\ell-1}| \leq \varepsilon$, for $m = 0, 1, \dots, M - 1$, $\ell = 1, 2, \dots, L - 1$, where $M = N^2L$, there exists $\phi \in \text{NN}(\#input = 2; \text{maxwidth} \leq 12N + 14; \#layer \leq 3L + 7)$ such that

1. $|\phi(m, \ell) - y_{m,\ell}| \leq \varepsilon$, for $m = 0, 1, \dots, M - 1$, $\ell = 0, 1, \dots, L - 1$;
2. $0 \leq \phi(t_1, t_2) \leq \max\{y_{m,\ell} : m = 0, 1, \dots, M - 1, \ell = 0, 1, \dots, L - 1\}$, for any $t_1, t_2 \in \mathbb{R}$.

Proof of Lemma 2.7. Define

$$a_{m,\ell} := \lfloor y_{m,\ell} / \varepsilon \rfloor, \quad \text{for } m = 0, 1, \dots, M - 1, \ell = 0, 1, \dots, L - 1.$$

We will construct a sub-FNN mapping the index (m, ℓ) to $a_{m,\ell}\varepsilon$ for $m = 0, 1, \dots, M - 1$, $\ell = 0, 1, \dots, L - 1$. Set $b_{m,0} := 0$ and $b_{m,\ell} := a_{m,\ell} - a_{m,\ell-1}$ for $m = 0, 1, \dots, M - 1$, $\ell = 1, \dots, L - 1$. Since $|y_{m,\ell} - y_{m,\ell-1}| \leq \varepsilon$ for all m and ℓ , we have $b_{m,\ell} \in \{-1, 0, 1\}$. Hence, there exist $c_{m,\ell}$ and $d_{m,\ell} \in \{0, 1\}$ such that $b_{m,\ell} = c_{m,\ell} - d_{m,\ell}$. It follows that

$$a_{m,\ell} = a_{m,0} + \sum_{j=1}^{\ell} (a_{m,j} - a_{m,j-1}) = a_{m,0} + \sum_{j=1}^{\ell} b_{m,j} = a_{m,0} + \sum_{j=0}^{\ell} b_{m,j} = a_{m,0} + \sum_{j=0}^{\ell} c_{m,j} - \sum_{j=0}^{\ell} d_{m,j}.$$

For the sample set $\{(m, a_{m,0}) : m = 0, 1, \dots, M - 1\} \cup \{(M, 0)\}$, whose cardinality is $M + 1 = N \cdot ((NL - 1) + 1) + 1$, by Proposition 2.1, there exists $\phi_1 \in \text{NN}(\#input = 1; \text{widthvec} = [2N, 2(NL - 1) + 1]) = \text{NN}(\#input = 1; \text{widthvec} = [2N, 2NL - 1])$ such that $\phi_1(m) = a_{m,0}$ for $m = 0, 1, \dots, M - 1$. By Lemma 2.6, there exist $\phi_2, \phi_3 \in \text{NN}(\#input = 2; \text{maxwidth} \leq 4N + 5; \#layer \leq 3L + 4)$ such that $\phi_2(m, \ell) = \sum_{j=0}^{\ell} c_{m,j}$ and $\phi_3(m, \ell) = \sum_{j=0}^{\ell} d_{m,j}$ for

$m = 0, 1, \dots, M-1$, $\ell = 0, 1, \dots, L-1$. Hence, it holds that

$$a_{m,\ell} = a_{m,0} + \sum_{j=0}^{\ell} c_{m,j} - \sum_{j=0}^{\ell} d_{m,j} = \phi_1(m) + \phi_2(m, \ell) - \phi_3(m, \ell),$$

for $m = 0, 1, \dots, M-1$, $\ell = 0, 1, \dots, L-1$. So, define

$$\bar{\phi}(t_1, t_2) := \left(\sigma(\phi_1(t_1)) + \sigma(\phi_2(t_1, t_2)) - \sigma(\phi_3(t_1, t_2)) \right) \varepsilon, \quad \text{for any } t_1, t_2 \in \mathbb{R}.$$

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

$$\begin{aligned} \bar{\phi}(m, \ell) &= \left(\sigma(\phi_1(m)) + \sigma(\phi_2(m, \ell)) - \sigma(\phi_3(m, \ell)) \right) \varepsilon \\ &= (\phi_1(m) + \phi_2(m, \ell) - \phi_3(m, \ell)) \varepsilon \\ &= \left(a_{m,0} + \sum_{j=0}^{\ell} c_{m,j} - \sum_{j=0}^{\ell} d_{m,j} \right) \varepsilon \\ &= a_{m,\ell} \varepsilon. \end{aligned}$$

And the architecture of $\bar{\phi}$ is

$$\begin{bmatrix} t_1 \\ t_2 \end{bmatrix} \rightarrow \begin{bmatrix} \sigma(\phi_1(t_1)) \\ \sigma(\phi_2(t_1, t_2)) \\ \sigma(\phi_3(t_1, t_2)) \end{bmatrix} \rightarrow \left(\sigma(\phi_1(t_1)) + \sigma(\phi_2(t_1, t_2)) - \sigma(\phi_3(t_1, t_2)) \right) \varepsilon = \bar{\phi}(t_1, t_2).$$

By Proposition 2.2, $\phi_1 \in \text{NN}(\#input = 1; \text{widthvec} = [2N, 2NL - 1]) \subseteq \text{NN}(\#input = 1; \text{maxwidth} \leq 4N + 4; \#layer \leq L + 2)$, which means $\bar{\phi} \in \text{NN}(\#input = 2; \text{maxwidth} \leq 12N + 14; \#layer \leq 3L + 5)$. Since $\bar{\phi}$ may not be bounded appropriately, we define

$$\phi(t_1, t_2) := \min\{\sigma(\bar{\phi}(t_1, t_2)), y_{\max}\}, \quad \text{for any } t_1, t_2 \in \mathbb{R},$$

where

$$y_{\max} := \max\{y_{m,\ell} : m = 0, 1, \dots, M-1, \ell = 0, 1, \dots, L-1\}.$$

It is clear that

$$0 \leq \phi(t_1, t_2) \leq y_{\max} = \max\{y_{m,\ell} : m = 0, 1, \dots, M-1, \ell = 0, 1, \dots, L-1\}, \quad \text{for any } t_1, t_2 \in \mathbb{R},$$

and that

$$\phi(m, \ell) = \min\{\sigma(\bar{\phi}(m, \ell)), y_{\max}\} = \min\{a_{m,\ell} \varepsilon, y_{\max}\} = a_{m,\ell} \varepsilon,$$

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

It follows that

$$|\phi(m, \ell) - y_{m,\ell}| = |a_{m,\ell} \varepsilon - y_{m,\ell}| = \left| \lfloor y_{m,\ell} / \varepsilon \rfloor \varepsilon - y_{m,\ell} \right| \leq \varepsilon,$$

for $m = 0, 1, \dots, M-1$, $\ell = 0, 1, \dots, L-1$. It remains to establish the architecture of ϕ . Since $\min\{t_1, t_2\} = \frac{t_1 + t_2 - |t_1 - t_2|}{2} = \frac{\sigma(t_1 + t_2) - \sigma(-t_1 - t_2) - \sigma(t_1 - t_2) - \sigma(t_2 - t_1)}{2}$ for any $t_1, t_2 \in \mathbb{R}$, the architecture of ϕ is

$$\begin{bmatrix} t_1 \\ t_2 \end{bmatrix} \rightarrow \sigma(\bar{\phi}(t_1, t_2)) \rightarrow \begin{bmatrix} \sigma(\sigma(\bar{\phi}(t_1, t_2)) + y_{\max}) \\ \sigma(-\sigma(\bar{\phi}(t_1, t_2)) - y_{\max}) \\ \sigma(\sigma(\bar{\phi}(t_1, t_2)) - y_{\max}) \\ \sigma(-\sigma(\bar{\phi}(t_1, t_2)) + y_{\max}) \end{bmatrix} \rightarrow \min\{\sigma(\bar{\phi}(t_1, t_2)), y_{\max}\} = \phi(t_1, t_2).$$

Therefore, $\phi \in \text{NN}(\#input = 2; \text{maxwidth} \leq 12N + 14; \#layer \leq 3L + 7)$. Hence, we finish the proof. \square

Finally, we apply Lemma 2.7 to prove Proposition 2.3.

Proof of Proposition 2.3. Let $M = N^2L$, then we may assume $K = ML$ since we can assume $y_{K-1} = y_K = y_{K+1} = \dots = y_{ML-1}$ if $K < ML$. For the sample set $\{(mL, m) : m = 0, 1, \dots, M\} \cup \{(mL + L - 1, m) : m = 0, 1, \dots, M - 1\}$, whose cardinality is $2M + 1 = N \cdot ((2NL - 1) + 1) + 1$, by Proposition 2.1 there exist $\phi_1 \in \text{NN}(\#input = 1; \text{widthvec} = [2N, 2(2NL - 1) + 1]) = \text{NN}(\#input = 1; \text{widthvec} = [2N, 4NL - 1])$ such that

- $\phi_1(ML) = M$ and $\phi_1(mL) = \phi_1(mL + L - 1) = m$ for $m = 0, 1, \dots, M - 1$;
- ϕ_1 is linear on each interval $[mL, mL + L - 1]$ for $m = 0, 1, \dots, M - 1$.

Note that any number k in $\{0, 1, \dots, K - 1\}$ can be indexed as $k = mL + \ell$ for $m = 0, 1, \dots, M - 1$, $\ell = 0, 1, \dots, L - 1$. Hence, we can denote $y_k = y_{mL + \ell}$ as $y_{m, \ell}$ for $m = 0, 1, \dots, M - 1$, $\ell = 0, 1, \dots, L - 1$, which gives

$$(2.20) \quad \begin{cases} \phi_1(k) = m, \\ k - L\phi_1(k) = \ell. \end{cases}$$

Then by Lemma 2.7, there exists $\phi_2 \in \text{NN}(\#input = 2; \text{maxwidth} \leq 12N + 14; \#layer \leq 3L + 7)$ such that

$$(2.21) \quad |\phi_2(m, \ell) - y_{m, \ell}| \leq \varepsilon, \quad \text{for } m = 0, 1, \dots, M - 1, \ell = 0, 1, \dots, L - 1,$$

and

$$(2.22) \quad 0 \leq \phi_2(t_1, t_2) \leq \max\{y_{m, \ell} : m = 0, 1, \dots, M - 1, \ell = 0, 1, \dots, L - 1\}, \quad \text{for any } t_1, t_2 \in \mathbb{R}.$$

Define $\phi(t) := \phi_2(\sigma(\phi_1(t)), \sigma(t) - L\sigma(\phi_1(t)))$ for any $t \in \mathbb{R}$. By Equation (2.20) and (2.21), it holds that $|\phi(k) - y_k| = |\phi_2(\sigma(\phi_1(k)), \sigma(k) - L\sigma(\phi_1(k))) - y_k| = |\phi_2(m, \ell) - y_{m, \ell}| \leq \varepsilon$ for $k \in \{0, 1, \dots, K - 1\}$. By Equation (2.22), for any $t \in \mathbb{R}$,

$$\begin{aligned} 0 \leq \phi(t) &= \phi_2(\sigma(\phi_1(t)), \sigma(t) - L\sigma(\phi_1(t))) \\ &\leq \max\{y_{m, \ell} : m = 0, 1, \dots, M - 1, \ell = 0, 1, \dots, L - 1\} \\ &= \max\{y_k : k = 0, 1, \dots, K - 1\}. \end{aligned}$$

And the architecture of ϕ is

$$t \rightarrow \begin{bmatrix} \sigma(\phi_1(t)) \\ \sigma(t) \end{bmatrix} \rightarrow \begin{bmatrix} \sigma(\phi_1(t)) \\ \sigma(\sigma(t) - L\sigma(\phi_1(t))) \\ \sigma(L\sigma(\phi_1(t)) - \sigma(t)) \end{bmatrix} \rightarrow \phi_2(\sigma(\phi_1(t)), \sigma(t) - L\sigma(\phi_1(t))) = \phi(t).$$

By Proposition 2.2, $\phi_1 \in \text{NN}(\#input = 1; \text{widthvec} = [2N, 4NL - 1]) \subseteq \text{NN}(\#input = 1; \text{maxwidth} \leq 8N + 4; \#layer \leq L + 2)$, which means $\phi \in \text{NN}(\#input = 1; \text{maxwidth} \leq 12N + 14; \#layer \leq 4L + 10)$ as expected. So we finish the proof. \square

We would like to remark that the key idea in the proof of Proposition 2.3 is the bit extraction technique in Lemma 2.5, which allows us to store L bits in a binary number $\text{bin} 0.x_1x_2 \dots x_L$ and extract each bit x_i . The extraction operator can be efficiently carried out via a deep ReLU neural network demonstrating the power of depth.

3. Optimality of the Constructive Approximation. This section will show that the approximation rate in Theorem 1.1 is tight and there is no room to improve for the function class $\text{Lip}(\nu, \alpha, d)$. First, we show that the approximation rate in Theorem 1.1 for the function class $\text{Lip}(\nu, \alpha, d)$ cannot be improved to $\mathcal{O}(\omega_f(N^{-(2/d+\rho)}L^{-(2/d+\rho)}))$ in the $L^\infty([0, 1]^d)$ -norm for any $\rho > 0$. This means that the approximation rate in Theorem 1.1 is “nearly” tight, though the norm in Theorem 1.1 is different to the $L^\infty([0, 1]^d)$ -norm. Second, to finally prove the tightness using the same norm, we introduce a special norm denoted as $L^{p, \infty}([0, 1]^d)$ -norm for $p \in [1, \infty)$ later. We will show that Theorem 1.1 is valid in the $L^{p, \infty}([0, 1]^d)$ -norm and the approximation rate cannot be improved to $\mathcal{O}(\omega_f(N^{-(2/d+\rho)}L^{-(2/d+\rho)}))$ in the $L^{p, \infty}([0, 1]^d)$ -norm as well.

3.1. Unachievable Approximation Rate $\mathcal{O}(\omega_f(N^{-(2/d+\rho)}L^{-(2/d+\rho)}))$. Now we show that the approximation rate in Theorem 1.1 is “nearly” tight for the function class $\text{Lip}(\nu, \alpha, d)$ by presenting an asymptotic unachievable approximation rate $\mathcal{O}(\omega_f(N^{-(2/d+\rho)}L^{-(2/d+\rho)}))$ for any $\rho > 0$ in Theorem 3.1 below.

THEOREM 3.1. *Given any $\rho > 0$ and $C > 0$, there exists $f \in \text{Lip}(\nu, \alpha, d)$ such that, for any $J_0 > 0$, there exist $N, L \in \mathbb{N}$ with $NL \geq J_0$ satisfying*

$$\inf_{\phi \in \text{NN}(\#\text{input}=d; \text{maxwidth} \leq N; \#\text{layer} \leq L)} \|\phi - f\|_{L^\infty([0,1]^d)} \geq C\nu N^{-(2\alpha/d+\rho)} L^{-(2\alpha/d+\rho)}.$$

In fact, we can show a stronger result than Theorem 3.1. Under the same conditions as in Theorem 3.1, for any $\mathcal{H} \in [0,1]^d$ with $\mu(\mathcal{H}) \leq 2^{-(d+K^d+1)}K^{-d}$, where $K = \lfloor (NL)^{2/d+\rho/(2\alpha)} \rfloor$, it can be proved that

$$(3.1) \quad \inf_{\phi \in \text{NN}(\#\text{input}=d; \text{maxwidth} \leq N; \#\text{layer} \leq L)} \|\phi - f\|_{L^\infty([0,1]^d \setminus \mathcal{H})} \geq C\nu N^{-(2\alpha/d+\rho)} L^{-(2\alpha/d+\rho)}.$$

We will prove (3.1) by contradiction, then Theorem 3.1 holds as a consequence. The result of (3.1) will be used later in Section 3.2. Assuming Equation (3.1) is false, we have the following claim.

CLAIM 3.2. *There exist $\rho > 0$ and $C > 0$ such that given any $f \in \text{Lip}(\nu, \alpha, d)$, there exists $J_0 = J_0(\rho, C, f) > 0$ such that, for any $N, L \in \mathbb{N}$ with $NL \geq J_0$, there exist $\phi \in \text{NN}(\#\text{input} = d; \text{maxwidth} \leq N; \#\text{layer} \leq L)$ and $\mathcal{H} \in [0,1]^d$ with $\mu(\mathcal{H}) \leq 2^{-(d+K^d+1)}K^{-d}$, where $K = \lfloor (NL)^{2/d+\rho/(2\alpha)} \rfloor$, satisfying*

$$\|f - \phi\|_{L^\infty([0,1]^d \setminus \mathcal{H})} \leq C\nu N^{-(2\alpha/d+\rho)} L^{-(2\alpha/d+\rho)}.$$

Now let us disprove this claim to show Theorem 3.1 and Equation (3.1) are true.

Disproof of Claim 3.2. Without the loss of generality, we assume $\nu = 1$; in the case of $\nu \neq 1$, the proof is similar. We will disprove Claim 3.2 using the VC dimension. Recall that the VC dimension of a class of functions is defined as the cardinality of the largest set of points that this class of functions can shatter. Denote the VC dimension of a function set \mathcal{F} by $\text{VCDim}(\mathcal{F})$. By [25], there exists $C_1 > 0$ such that

$$(3.2) \quad \text{VCDim}(\text{NN}(\#\text{input} = d; \#\text{parameter} \leq W; \#\text{layer} \leq L)) \leq C_1 W L \ln W.$$

Hence,

$$(3.3) \quad \begin{aligned} & \text{VCDim}(\text{NN}(\#\text{input} = d; \text{maxwidth} \leq N; \#\text{layer} \leq L)) \\ & \leq \text{VCDim}(\text{NN}(\#\text{input} = d; \#\text{parameter} \leq (LN + d + 2)(N + 1); \#\text{layer} \leq L)) \\ & \leq C_1 (LN + d + 2)(N + 1) L \ln((LN + d + 2)(N + 1)) \\ & := b_u(N, L). \end{aligned}$$

Then we will use Claim 3.2 to estimate a lower bound of

$$(3.4) \quad \text{VCDim}(\text{NN}(\#\text{input} = d; \text{maxwidth} \leq N; \#\text{layer} \leq L)),$$

and this lower bound is asymptotically larger than $b_u(N, L)$, which leads to a contradiction.

More precisely, we will construct $\{f_\beta : \beta \in \mathcal{B}\} \subseteq \text{Lip}(1, \alpha, d)$, which can shatter $b_\ell(N, L) := K^d$ points, where \mathcal{B} is a set defined later. Then by Claim 3.2, there exists $\{\phi_\beta : \beta \in \mathcal{B}\}$ such that this set can shatter $b_\ell(N, L)$ points. Finally, $b_\ell(N, L) = K^d = \lfloor (NL)^{2/d+\rho/(2\alpha)} \rfloor^d$ is asymptotically larger than $b_u(N, L) = C_1 (LN + d + 2)(N + 1) L \ln((LN + d + 2)(N + 1))$, which leads to a contradiction. More details can be found below.

Step 1: Construct $\{f_\beta : \beta \in \mathcal{B}\} \subseteq \text{Lip}(1, \alpha, d)$ that scatters $b_\ell(N, L)$ points.

Let $Q(\mathbf{x}_0, r) \subseteq [0,1]^d$ be a cube, whose center and sidelength are \mathbf{x}_0 and r , respectively. Then we define a function ζ_Q on $[0,1]^d$ corresponding to $Q = Q(\mathbf{x}_0, r) \subseteq [0,1]^d$ such that:

- $\zeta_Q(\mathbf{x}_0) = (r/2)^\alpha/2$;
- $\zeta_Q(\mathbf{x}) = 0$ for any $\mathbf{x} \notin Q \setminus \partial Q$, where ∂Q is the boundary of Q ;
- ζ_Q is linear on the line that connects \mathbf{x}_0 and \mathbf{x} for any $\mathbf{x} \in \partial Q$.

Divide $[0, 1]^d$ into K^d non-overlapping sub-cubes $\{Q_\theta\}_\theta$ as follows:

$$Q_\theta := \{\mathbf{x} = [x_1, x_2, \dots, x_d]^T \in [0, 1]^d : x_i \in [\frac{\theta_i-1}{K}, \frac{\theta_i}{K}], \ i = 1, 2, \dots, d\},$$

for any index vector $\theta = [\theta_1, \theta_2, \dots, \theta_d]^T \in \{1, 2, \dots, K\}^d$. Define

$$\mathcal{B} := \{\beta : \beta \text{ is a map from } \{1, 2, \dots, K\}^d \text{ to } \{-1, 1\}\}.$$

For each $\beta \in \mathcal{B}$, we define

$$f_\beta(\mathbf{x}) := \sum_{\theta \in \{1, 2, \dots, K\}^d} \beta(\theta) \zeta_{Q_\theta}(\mathbf{x}),$$

where $\zeta_{Q_\theta}(\mathbf{x})$ is the associated function introduced just above. It is easy to check that $f_\beta \in \text{Lip}(1, \alpha, d)$ and $\{f_\beta : \beta \in \mathcal{B}\}$ can shatter $b_\ell(N, L) = K^d$ points.

Step 2: Construct $\{\phi_\beta : \beta \in \mathcal{B}\}$ that scatters $b_\ell(N, L)$ points.

By Claim 3.2, there exist $\rho > 0$ and $C_2 > 0$ such that for any $f_\beta \in \{f_\beta : \beta \in \mathcal{B}\}$ there exists $J_\beta > 0$ such that for all $N, L \in \mathbb{N}$ with $NL \geq J_\beta$, there exist $\phi_\beta \in \text{NN}(\#\text{input} = 1; \text{maxwidth} \leq N; \#\text{layer} \leq L)$ and \mathcal{H}_β with $\mu(\mathcal{H}_\beta) \leq 2^{-(d+K^d+1)} K^{-d}$ such that

$$|f_\beta(\mathbf{x}) - \phi_\beta(\mathbf{x})| \leq C_2(NL)^{-\alpha(2/d+\rho/\alpha)}, \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \mathcal{H}_\beta.$$

Set $\mathcal{H} = \cup_{\beta \in \mathcal{B}} \mathcal{H}_\beta$ and $J_1 = \max_{\beta \in \mathcal{B}} J_\beta$. Then it holds that

$$(3.5) \quad \mu(\mathcal{H}) \leq 2^{K^d} 2^{-(d+K^d+1)} K^{-d} = (2K)^{-d}/2.$$

It follows that for all $\beta \in \mathcal{B}$ and $N, L \in \mathbb{N}$ with $NL \geq J_1$, we have

$$(3.6) \quad |f_\beta(\mathbf{x}) - \phi_\beta(\mathbf{x})| \leq C_2(NL)^{-\alpha(2/d+\rho/\alpha)}, \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \mathcal{H}.$$

For each index vector $\theta \in \{1, 2, \dots, K\}^d$ and any $\mathbf{x} \in \frac{1}{2}Q_\theta$, where $\frac{1}{2}Q_\theta$ denotes the cube whose sidelength is half of that of Q_θ sharing the same center of Q_θ , since Q_θ has a sidelength $\frac{1}{K}$, we have

$$(3.7) \quad |f_\beta(\mathbf{x})| = |g_{Q_\theta}(\mathbf{x})| \geq |g_{Q_\theta}(\mathbf{x}_{Q_\theta})|/2 = \left(\frac{1}{2K}\right)^\alpha/4,$$

where \mathbf{x}_{Q_θ} is the center of Q_θ . For fixed d, α , and ρ , there exists $J_2 > 0$ large enough such that, for any $N, L \in \mathbb{N}$ with $NL \geq J_2$, we have

$$(3.8) \quad \frac{1}{2^{2+\alpha}} \lfloor (NL)^{2/d+\rho/(2\alpha)} \rfloor^{-\alpha} > C_2(NL)^{-\alpha(2/d+\rho/\alpha)}.$$

By (3.5), for any $\theta \in \{1, 2, \dots, K\}^d$, we have

$$\mu(\mathcal{H}) \leq (2K)^{-d}/2 < (2K)^{-d} = \mu(\tfrac{1}{2}Q_\theta),$$

which means $\frac{1}{2}Q_\theta \cap ([0, 1]^d \setminus \mathcal{H})$ is not empty. Therefore, there exists $\mathbf{x}_\theta \in \frac{1}{2}Q_\theta \cap ([0, 1]^d \setminus \mathcal{H})$ for each $\theta \in \{1, 2, \dots, K\}^d$ such that

$$|f_\beta(\mathbf{x}_\theta)| \geq \frac{1}{2^{2+\alpha}} \lfloor (NL)^{2/d+\rho/(2\alpha)} \rfloor^{-\alpha} > C_2(NL)^{-\alpha(2/d+\rho/\alpha)} \geq |f_\beta(\mathbf{x}_\theta) - \phi_\beta(\mathbf{x}_\theta)|,$$

where the first, the second, and the last inequalities come from (3.7), (3.8), and (3.6), respectively. In other words, for any $\beta \in \mathcal{B}$ and $\boldsymbol{\theta} \in \{1, 2, \dots, K\}^d$, $f_\beta(\mathbf{x}_\boldsymbol{\theta})$ and $\phi_\beta(\mathbf{x}_\boldsymbol{\theta})$ have the same sign. Then $\{\phi_\beta : \beta \in \mathcal{B}\}$ shatters $\{\mathbf{x}_\boldsymbol{\theta} : \boldsymbol{\theta} \in \{1, 2, \dots, K\}^d\}$ since $\{f_\beta : \beta \in \mathcal{B}\}$ shatters $\{\mathbf{x}_\boldsymbol{\theta} : \boldsymbol{\theta} \in \{1, 2, \dots, K\}^d\}$ as discussed in Step 1. Hence,

$$(3.9) \quad \text{VCDim}(\{\phi_\beta : \beta \in \mathcal{B}\}) \geq K^d = b_\ell(N, L).$$

Step 3: Contradiction.

By (3.3) and (3.9), for any $N, L \in \mathbb{N}$ with $NL \geq J_0 = \max\{J_1, J_2\}$, we have

$$b_\ell(N, L) \leq \text{VCDim}(\{\phi_\beta : \beta \in \mathcal{B}\}) \leq \text{VCDim}(\text{NN}(\#\text{input} = d; \text{maxwidth} \leq N; \#\text{layer} \leq L)) \leq b_u(N, L),$$

implying that

$$[(NL)^{2/d+\rho/(2\alpha)}]^d \leq C_1(LN + d + 2)(N + 1)L \ln((LN + d + 2)(N + 1)),$$

which is a contradiction for sufficiently large $N, L \in \mathbb{N}$. So we finish the proof. \square

By Theorem 3.1, for any $\rho > 0$, the approximation rate cannot be $\mathcal{O}(N^{-(2\alpha/d+\rho)}L^{-(2/\alpha+\rho)})$, if we use FNNs in $\text{NN}(\#\text{input} = d; \text{maxwidth} \leq N; \#\text{layer} \leq L)$ to approximate functions in $\text{Lip}(\nu, \alpha, d)$. By a similar argument, we can show that the approximation rate cannot be $\mathcal{O}(N^{-2\alpha/d}L^{-(2/\alpha+\rho)})$ nor $\mathcal{O}(N^{-(2\alpha/d+\rho)}L^{-2\alpha/d})$. Hence, the approximation rate in Theorem 1.1 is nearly tight.

3.2. Tightness in the $L^{p,\infty}$ -Norm. In Section 2, the quantitative approximation rate was obtained using the L^∞ -norm on $[0, 1]^d \setminus \mathcal{H}$ or the L^p -norm on $[0, 1]^d$ for $p \in [1, \infty)$ as a consequence, while the L^∞ -norm on $[0, 1]^d$ was used to prove the unachievable approximation rate in Section 3.1. In order to make our results tight in the same norm, we introduce a new norm, the $L^{1,\infty}$ -norm below. Similarly, we can define $L^{p,\infty}$ -norm for $p \in [1, \infty)$. The main intuition of this $L^{1,\infty}$ -norm is to rely on a small “don’t-care” region, where the approximation error can be large, while uniformly controlling the approximation error in other “important” regions.

DEFINITION 3.3. *Shrinking function.* We define the shrinking function as

$$\lambda(x) := \frac{2^{-(2^x)} - 2^{-(2^{x+1})}}{x}, \quad \text{for any } x \in [1, \infty).$$

It is clear that the shrinking function $\lambda(\cdot)$ is strictly decreasing in x on $[1, \infty)$ and satisfies

$$(3.10) \quad \sum_{j=J}^{\infty} j\lambda(j) \leq 2^{-(2^J)}, \quad \text{for any } J \in \mathbb{N}^+.$$

Actually, the shrinking function can be simply defined as $\lambda(x) := \frac{2^{-x}}{x}$ for any $x \in [1, \infty)$, but the proof will be more complicated. Let $\lambda_{\text{inv}}(\cdot)$ be the inverse function of $\lambda(\cdot)$, i.e., $\lambda(\lambda_{\text{inv}}(x)) = x$ and $\lambda_{\text{inv}}(\lambda(x)) = x$. Now we can define the “don’t-care” region by using the shrinking function $\lambda(\cdot)$.

DEFINITION 3.4. *Shrinking region.* Let $\mathcal{S}(\ell) = \cup_{j=1}^{\ell} [j/\ell - \lambda(\ell), j/\ell]$. Then the shrinking region corresponding to integers k and d is defined as

$$(3.11) \quad \Omega(k, d) := \cup_{\ell=k}^{\infty} \left(\cup_{i=1}^d \{ \mathbf{x} = [x_1, \dots, x_d]^T \in \mathbb{R}^d : x_i \in \mathcal{S}(\ell) \} \right).$$

The shrinking region $\Omega(k, d)$ is the so-called “don’t-care” region in this paper. Furthermore, $\Omega(k, d)$ shrinks gradually to a set of rational points in $[0, 1]^d$ as $k \rightarrow \infty$, and the shrinking speed is determined by the shrinking function $\lambda(\cdot)$. More precisely,

$$(3.12) \quad \mu(\Omega(k, d)) \leq \sum_{j=k}^{\infty} dj\lambda(j) \leq d2^{-(2^k)}, \quad \text{for any } k, d \in \mathbb{N}^+.$$

Using the “don’t-care” region, we are ready to define a new norm denoted as $L^{1,\infty}$.

DEFINITION 3.5. *For any $f \in L^1([0, 1]^d) \cap L^\infty([0, 1]^d)$, we define*

$$\|f\|_{L^{1,\infty}([0,1]^d)} := \sup \left\{ \frac{1}{\ln k} \|f\|_{L^\infty([0,1]^d \setminus \Omega(k,d))} : k \geq 3, k \in \mathbb{N} \right\} + \|f\|_{L^1([0,1]^d)}.$$

It follows that $\mu(\Omega(k,d))$ decays quickly, so $\|f\|_{L^\infty([0,1]^d \setminus \Omega(k,d))}$ approaches to $\|f\|_{L^\infty([0,1]^d)}$ as $k \rightarrow \infty$. To define a norm weaker than the L^∞ -norm, we need the factor $\frac{1}{\ln k}$ in Definition 3.5, which can be replaced by other functions going to 0 as $k \rightarrow +\infty$. It is easy to verify that $\|\cdot\|_{L^{1,\infty}([0,1]^d)}$ ($\|\cdot\|_{1,\infty}$ for short) is a norm satisfying

$$\|f\|_{L^1([0,1]^d)} \leq \|f\|_{1,\infty} \leq 2\|f\|_{L^\infty([0,1]^d)}, \quad \text{for any } f \in L^1([0,1]^d) \cap L^\infty([0,1]^d).$$

In fact, to make the $L^{1,\infty}$ -norm well-defined, the shrinking function $\lambda(\cdot)$ can be any strictly decreasing function defined on $[1, \infty)$ with $\lim_{x \rightarrow \infty} \lambda(x) = 0$.

Theorem 3.6 below provides a quantitative approximation rate in the $L^{1,\infty}$ -norm, which is asymptotically tight by Theorem 3.7.

THEOREM 3.6. *Let f be a continuous function on $[0, 1]^d$. For any $L, N \in \mathbb{N}^+$, there exists a ReLU FNN ϕ with width $\max\{8d\lfloor N^{1/d} \rfloor + 4d, 12N + 14\}$ and depth $9L + 12$ such that*

$$\|f - \phi\|_{1,\infty} \leq 5\omega_f(8\sqrt{d}N^{-2/d}L^{-2/d}).$$

Proof of Theorem 3.6. We divide the proof into two cases: the case $d = 1$ and the case $d \geq 2$.

Case 1: The case $d = 1$.

All the settings are the same as in Section 2.3 except the conditions on δ . The reader is referred to Section 2.3 for the definitions of notations. In Section 2.3, we require (2.1), but now we let $\delta = \delta(N, L, f) > 0$ be a sufficiently small number satisfying Equation (3.13), (3.14), and (3.15). It remains to estimate the approximation error in the $L^{1,\infty}$ -norm. By Equation (2.7) and (2.8), we get

$$\begin{aligned} \|f - \phi\|_{L^1([0,1])} &= \sum_{m=0}^{M-1} \sum_{\ell=0}^{L-1} \int_{Q_{m,\ell}} |f(x) - \phi(x)| dx + \int_{\mathcal{H}} |\bar{f}(x) - \bar{\phi}(x)| dx \\ &\leq \sum_{m=0}^{M-1} \sum_{\ell=0}^{L-1} \frac{1}{ML} 2\omega_f\left(\frac{1}{ML}\right) + 4ML\delta\omega_f(1) \\ &\leq 3\omega_f\left(\frac{1}{ML}\right), \end{aligned}$$

where the last inequality comes from

$$(3.13) \quad 4ML\delta\omega_f(1) \leq \omega_f\left(\frac{1}{ML}\right),$$

when δ is chosen to be small enough. In fact, we also require that δ satisfies

$$(3.14) \quad \frac{4\omega_f(1)}{\ln \lambda_{\text{inv}}(\delta)} \leq 2\omega_f\left(\frac{1}{ML}\right),$$

and

$$(3.15) \quad \frac{\lambda_{\text{inv}}(\delta)}{ML} \in \mathbb{N}^+,$$

where $\lambda(\cdot)$ and $\lambda_{\text{inv}}(\cdot)$ were introduced in Definition 3.3. By (3.14), we have

$$\sup \left\{ \frac{\|f - \phi\|_{L^\infty([0,1] \setminus \Omega(k,d))}}{\ln k} : k \geq \lambda_{\text{inv}}(\delta), k \in \mathbb{N} \right\} \leq \sup \left\{ \frac{4\omega_f(1)}{\ln k} : k \geq \lambda_{\text{inv}}(\delta), k \in \mathbb{N} \right\} \leq \frac{4\omega_f(1)}{\ln \lambda_{\text{inv}}(\delta)} \leq 2\omega_f\left(\frac{1}{ML}\right).$$

By (3.15) and the definition of $\Omega(k, d)$ in (3.11), if $k \leq \lambda_{\text{inv}}(\delta)$, we have

$$\Omega(k, 1) \supseteq \Omega(\lambda_{\text{inv}}(\delta), 1) \supseteq \cup_{j=1}^{\lambda_{\text{inv}}(\delta)} \left[\frac{j}{\lambda_{\text{inv}}(\delta)} - \lambda(\lambda_{\text{inv}}(\delta)), \frac{j}{\lambda_{\text{inv}}(\delta)} \right] \supseteq \cup_{j=1}^{ML} \left[\frac{j}{ML} - \delta, \frac{j}{ML} \right] = \mathcal{H},$$

where the last equality comes from (2.2). Together with (2.8), it holds that

$$\sup \left\{ \frac{1}{\ln k} \|f - \phi\|_{L^\infty([0,1] \setminus \Omega(k,1))} : 3 \leq k \leq \lambda_{\text{inv}}(\delta), k \in \mathbb{N} \right\} \leq \frac{1}{\ln 3} \|f - \phi\|_{L^\infty([0,1] \setminus \mathcal{H})} \leq 2\omega_f\left(\frac{1}{ML}\right).$$

Then we get

$$\sup \left\{ \frac{1}{\ln k} \|f - \phi\|_{L^\infty([0,1] \setminus \Omega(k,1))} : k \geq 3, k \in \mathbb{N} \right\} \leq 2\omega_f\left(\frac{1}{ML}\right).$$

Therefore, we have

$$\|f - \phi\|_{1,\infty} = \sup \left\{ \frac{1}{\ln k} \|f - \phi\|_{L^\infty([0,1] \setminus \Omega(k,1))} : k \geq 3, k \in \mathbb{N} \right\} + \|f - \phi\|_{L^1([0,1])} \leq 5\omega_f\left(\frac{1}{ML}\right) = 5\omega_f(N^{-2}L^{-2}).$$

So we finish the proof for the case $d = 1$.

Case 2: The case $d \geq 2$.

All the settings are the same as in Section 2.4 except the conditions on δ . The reader is referred to Section 2.4 for the definitions of notations. In Section 2.4, we require (2.9), but now we let $\delta = \delta(N, L, d, f) > 0$ be a sufficiently small number satisfying Equation (3.16), (3.17), and (3.18). It remains to estimate the approximation error in the $L^{1,\infty}$ -norm. By Equation (2.13) and (2.14), we have

$$\begin{aligned} \|f - \phi\|_{L^1([0,1]^d)} &= \int_{\mathcal{H}} |\bar{f} - \bar{\phi}| d\mathbf{x} + \int_{[0,1]^d \setminus \mathcal{H}} |f - \phi| d\mathbf{x} \\ &\leq \mu(\mathcal{H}) (2\omega_f(\sqrt{d}) + 2\omega_f(\sqrt{d})) + \sum_{\boldsymbol{\theta} \in \{0,1,\dots,n-1\}^d} \int_{Q_{\boldsymbol{\theta}}} |f - \phi| d\mathbf{x} \\ &\leq 4J\delta d\omega_f(\sqrt{d}) + \sum_{\boldsymbol{\theta} \in \{0,1,\dots,n-1\}^d} 2\omega_f\left(\frac{\sqrt{d}}{J}\right) \mu(Q_{\boldsymbol{\theta}}) \\ &\leq 3\omega_f\left(\frac{\sqrt{d}}{J}\right), \end{aligned}$$

where the last inequality comes from

$$(3.16) \quad 4J\delta d\omega_f(\sqrt{d}) \leq \omega_f\left(\frac{\sqrt{d}}{J}\right).$$

In fact, we also require

$$(3.17) \quad \frac{4\omega_f(\sqrt{d})}{\ln \lambda_{\text{inv}}(\delta)} \leq 2\omega_f\left(\frac{\sqrt{d}}{J}\right)$$

and

$$(3.18) \quad \frac{\lambda_{\text{inv}}(\delta)}{J} \in \mathbb{N}^+.$$

By Equation (3.17), we have

$$\sup \left\{ \frac{\|f - \phi\|_{L^\infty([0,1]^d \setminus \Omega(k,d))}}{\ln k} : k \geq \lambda_{\text{inv}}(\delta), k \in \mathbb{N} \right\} \leq \sup \left\{ \frac{4\omega_f(\sqrt{d})}{\ln k} : k \geq \lambda_{\text{inv}}(\delta), k \in \mathbb{N} \right\} \leq \frac{4\omega_f(\sqrt{d})}{\ln \lambda_{\text{inv}}(\delta)} \leq 2\omega_f\left(\frac{\sqrt{d}}{J}\right).$$

By (3.18) and the definition of $\Omega(k, d)$, if $k \leq \lambda_{\text{inv}}(\delta)$, we have

$$\Omega(k, d) \supseteq \Omega(\lambda_{\text{inv}}(\delta), d) \supseteq \cup_{i=1}^d \left\{ \mathbf{x} \in [0,1]^d : x_i \in \cup_{j=1}^J \left[\frac{j}{J} - \delta, \frac{j}{J} \right] \right\} = \mathcal{H}.$$

Together with (2.14), it holds that

$$\sup \left\{ \frac{1}{\ln k} \|f - \phi\|_{L^\infty([0,1]^d \setminus \Omega(k,d))} : 3 \leq k \leq \lambda_{\text{inv}}(\delta), k \in \mathbb{N} \right\} \leq \frac{1}{\ln 3} \|f - \phi\|_{L^\infty([0,1]^d \setminus \mathcal{H})} \leq 2\omega_f\left(\frac{\sqrt{d}}{J}\right).$$

Furthermore,

$$\sup \left\{ \frac{1}{\ln k} \|f - \phi\|_{L^\infty([0,1]^d \setminus \Omega(k,d))} : k \geq 3, k \in \mathbb{N} \right\} \leq 2\omega_f\left(\frac{\sqrt{d}}{J}\right).$$

Therefore,

$$\|f - \phi\|_{1,\infty} = \sup \left\{ \frac{\|f - \phi\|_{L^\infty([0,1]^d \setminus \Omega(k,d))}}{\ln k} : k \geq 3, k \in \mathbb{N} \right\} + \|f - \phi\|_{L^1([0,1])} \leq 5\omega_f\left(\frac{\sqrt{d}}{J}\right) \leq 5\omega_f(8\sqrt{d}N^{-2/d}L^{-2/d}),$$

where the last inequality comes from the fact $J = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor \geq \frac{N^{2/d}L^{2/d}}{8}$ for any $N, L \in \mathbb{N}^+$. So we finish the proof for the case $d \geq 2$. \square

Note that $\omega_f(r) \leq \nu r^\alpha$ for any $r \geq 0$ if $f \in \text{Lip}(\nu, \alpha, d)$. Hence, we have the quantitative approximation rate $5(8\sqrt{d})^\alpha \nu N^{-2\alpha/d} L^{-2\alpha/d}$ in the $L^{1,\infty}$ -norm if we use ReLU FNNs with width $\max\{8d\lfloor N^{1/d} \rfloor + 4d, 12N + 14\}$ and depth $9L + 12$ to approximate functions in $\text{Lip}(\nu, \alpha, d)$ on $[0, 1]^d$.

The following theorem estimates an asymptotic and unachivable approximation rate in the $L^{1,\infty}$ -norm, showing that the quantitative approximation rate in Theorem 3.6 is tight for the function class $\text{Lip}(\nu, \alpha, d)$.

THEOREM 3.7. *For any $\rho > 0$ and $C > 0$, there exists $f \in \text{Lip}(\nu, \alpha, d)$, for any $J_0 = J_0(\rho, C, f) > 0$, there exist $N, L \in \mathbb{N}$ with $NL \geq J_0$ such that*

$$\inf_{\phi \in \text{NN}(\#\text{input}=d; \text{maxwidth} \leq N; \#\text{layer} \leq L)} \|\phi - f\|_{1,\infty} \geq C\nu N^{-(2\alpha/d+\rho)} L^{-(2\alpha/d+\rho)}.$$

Proof of Theorem 3.7. The proof is mainly based on Theorem 3.1 and Equation (3.1). By Equation (3.12), let $\tilde{\rho} = \rho/2$ and $\tilde{J}_0 \geq J_0 + 3$ be a large number satisfying

$$\frac{1}{\ln(NL)} (NL)^{-(2\alpha/d+\rho/2)} \geq (NL)^{-(2\alpha/d+\rho)}, \quad \text{for any } N, L \text{ with } NL \geq \tilde{J}_0,$$

and

$$\mu(\Omega(NL, d)) \leq d2^{-(2^{NL})} \leq 2^{-(d+K^d)} K^{-d}, \quad \text{for any } N, L \text{ with } NL \geq \tilde{J}_0, \text{ where } K = \lfloor (NL)^{2/d+\tilde{\rho}/(2\alpha)} \rfloor.$$

Set $\mathcal{H} = \Omega(NL, d)$ and denote $\text{NN}(\#\text{input} = d; \text{maxwidth} \leq N; \#\text{layer} \leq L)$ by \mathcal{N} for short. Then by Theorem 3.1, there exist $N, L \in \mathbb{N}$ with $NL \geq \tilde{J}_0 \geq J_0 + 3$ such that

$$\inf_{\phi \in \mathcal{N}} \|\phi - f\|_{L^\infty([0,1]^d \setminus \Omega(NL,d))} \geq C\nu(NL)^{-(2\alpha/d+\tilde{\rho})} = C\nu(NL)^{-(2\alpha/d+\rho/2)}.$$

It follows that

$$\begin{aligned} \inf_{\phi \in \mathcal{N}} \|\phi - f\|_{1,\infty} &\geq \inf_{\phi \in \mathcal{N}} \sup \left\{ \frac{1}{\ln k} \|\phi - f\|_{L^\infty([0,1]^d \setminus \Omega(k,d))} : k \geq 3, k \in \mathbb{N} \right\} \\ &\geq \inf_{\phi \in \mathcal{N}} \frac{1}{\ln(NL)} \|\phi - f\|_{L^\infty([0,1]^d \setminus \Omega(NL,d))} \\ &\geq C\nu(NL)^{-(2\alpha/d+\rho)}. \end{aligned}$$

So, we finish the proof. \square

4. Neural Networks in Practice. This section is concerned with neural networks in practice, e.g., approximating functions defined on irregular domains or domains with a low-dimensional structure, and neural network computation in parallel computing.

4.1. Approximation on Irregular Domain. In this section, we consider approximating continuous functions defined on irregular domains by deep ReLU FNNs. The construction is through extending the target function to a cubic domain, applying Theorem 1.1, and finally restricting the constructed FNN back to the irregular domain.

Given any uniformly continuous and real-valued function f defined on a metric space S with a metric $d_S(\cdot, \cdot)$, we define the (optimal) modulus of continuity of f on a subset $E \subseteq S$ as

$$\omega_f^E(r) := \sup\{|f(\mathbf{x}_1) - f(\mathbf{x}_2)| : \mathbf{x}_1, \mathbf{x}_2 \in E, d_S(\mathbf{x}_1, \mathbf{x}_2) \leq r\}, \quad \text{for any } r \geq 0.$$

For the purpose of consistency and simplicity, $\omega_f(\cdot)$ is short of $\omega_f^{[0,1]^d}(\cdot)$.

First, let us present two lemmas for (approximately) extending (almost) continuous functions on E to (almost) continuous functions on S . These lemmas are similar to the well-known results for extending Lipchitz or differentiable functions in [43, 59]. We generalize these results to a broader class of functions required in the proof of Theorem 4.3.

LEMMA 4.1 (Approximate Extension of Almost-Continuous Functions). *Assume S is a metric space with a metric $d_S(\cdot, \cdot)$ and $\omega : [0, \infty) \rightarrow [0, \infty)$ is an increasing function with*

$$(4.1) \quad \omega(r_1 + r_2) \leq \omega(r_1) + \omega(r_2), \quad \text{for any } r_1, r_2 \in [0, \infty).$$

Let the real-valued function f be defined on a subset $E \subseteq S$ and satisfy

$$(4.2) \quad |f(\mathbf{x}_1) - f(\mathbf{x}_2)| \leq \omega(d_S(\mathbf{x}_1, \mathbf{x}_2) + \Delta), \quad \text{for any } \mathbf{x}_1, \mathbf{x}_2 \in E,$$

where Δ is a positive constant independent of f . Then there exists a function g defined on S such that

$$0 \leq f(\mathbf{x}) - g(\mathbf{x}) \leq \omega(\Delta), \quad \text{for any } \mathbf{x} \in E$$

and

$$|g(\mathbf{x}_1) - g(\mathbf{x}_2)| \leq \omega(d_S(\mathbf{x}_1, \mathbf{x}_2)), \quad \text{for any } \mathbf{x}_1, \mathbf{x}_2 \in S.$$

In Lemma 4.1, g is an approximate extension of f defined on E to a new domain S with an approximation error $\omega(\Delta)$. In a special case when $\Delta = 0$ and $\omega(0) = 0$, g is an exact extension of f .

Proof of Lemma 4.1. Define

$$g(\mathbf{x}) := \sup_{\mathbf{z} \in E} (f(\mathbf{z}) - \omega(d_S(\mathbf{z}, \mathbf{x}) + \Delta)).$$

By Equation (4.2), we have $f(\mathbf{x}_1) - \omega(d_S(\mathbf{x}_1, \mathbf{x}_2) + \Delta) \leq f(\mathbf{x}_2)$ for any $\mathbf{x}_1, \mathbf{x}_2 \in E$. It holds that $g(\mathbf{x}) \leq f(\mathbf{x})$ for any $\mathbf{x} \in E$. Together with

$$g(\mathbf{x}) = \sup_{\mathbf{z} \in E} (f(\mathbf{z}) - \omega(d_S(\mathbf{z}, \mathbf{x}) + \Delta)) \geq f(\mathbf{x}) - \omega(d_S(\mathbf{x}, \mathbf{x}) + \Delta) = f(\mathbf{x}) - \omega(\Delta), \quad \text{for any } \mathbf{x} \in E,$$

it follows that $0 \leq f(\mathbf{x}) - g(\mathbf{x}) \leq \omega(\Delta)$ for any $\mathbf{x} \in E$. By Equation (4.1) and the fact

$$\sup_{\mathbf{z} \in E} f_1(\mathbf{z}) - \sup_{\mathbf{z} \in E} f_2(\mathbf{z}) \leq \sup_{\mathbf{z} \in E} (f_1(\mathbf{z}) - f_2(\mathbf{z})), \quad \text{for any functions } f_1, f_2,$$

we have

$$\begin{aligned}
g(\mathbf{x}_1) - g(\mathbf{x}_2) &= \sup_{\mathbf{z} \in E} (f(\mathbf{z}) - \omega(d_S(\mathbf{z}, \mathbf{x}_1))) - \sup_{\mathbf{z} \in E} (f(\mathbf{z}) - \omega(d_S(\mathbf{z}, \mathbf{x}_2))) \\
&\leq \sup_{\mathbf{z} \in E} (\omega(d_S(\mathbf{z}, \mathbf{x}_1)) - \omega(d_S(\mathbf{z}, \mathbf{x}_2))) \\
&\leq \sup_{\mathbf{z} \in E} \omega(d_S(\mathbf{z}, \mathbf{x}_1) - d_S(\mathbf{z}, \mathbf{x}_2)) \\
&\leq \sup_{\mathbf{z} \in E} \omega(d_S(\mathbf{x}_1, \mathbf{x}_2)) \\
&= \omega(d_S(\mathbf{x}_1, \mathbf{x}_2)),
\end{aligned}$$

for any $\mathbf{x}_1, \mathbf{x}_2 \in S$. Similarly, we have $g(\mathbf{x}_2) - g(\mathbf{x}_1) \leq \omega(d_S(\mathbf{x}_1, \mathbf{x}_2))$, which implies

$$|g(\mathbf{x}_1) - g(\mathbf{x}_2)| \leq \omega(d_S(\mathbf{x}_1, \mathbf{x}_2)).$$

So we finish the proof. \square

Next, we introduce a lemma below for extending continuous functions defined on $E \subseteq S$ to continuous functions defined on S preserving the modulus of continuity.

LEMMA 4.2 (Extension of Continuous Functions). *Suppose f is a uniformly continuous function defined on a subset $E \subseteq S$, where S is a metric space with a metric $d_S(\cdot, \cdot)$, then there exists a uniformly continuous function g on S such that $f(x) = g(x)$ for $x \in E$ and $\omega_f^E(r) = \omega_g^S(r)$ for any $r \geq 0$.*

Proof of Lemma 4.2. By the application of Lemma 4.1 with $\omega(r) = \omega_f^E(r)$ for $r \geq 0$ and $\Delta = 0$, we know that there exists $g : S \rightarrow \mathbb{R}$ such that

$$0 \leq f(\mathbf{x}) - g(\mathbf{x}) \leq \omega_f^E(\Delta) = 0, \quad \text{for any } \mathbf{x} \in E,$$

and

$$|g(\mathbf{x}_1) - g(\mathbf{x}_2)| \leq \omega_f^E(d_S(\mathbf{x}_1, \mathbf{x}_2)), \quad \text{for any } \mathbf{x}_1, \mathbf{x}_2 \in S.$$

The equation above and the uniform continuity of f imply that g is uniformly continuous. It also follows that

$$f(\mathbf{x}) = g(\mathbf{x}), \quad \text{for any } \mathbf{x} \in E, \quad \text{and} \quad \omega_g^S(r) \leq \omega_f^E(r), \quad \text{for any } r \geq 0,$$

since $\omega_g^S(\cdot)$ is the optimal modulus of continuity of g . Note that $\omega_f^E(\cdot)$ is the optimal modulus of continuity of f and

$$|f(\mathbf{x}_1) - f(\mathbf{x}_2)| = |g(\mathbf{x}_1) - g(\mathbf{x}_2)| \leq \omega_g^S(d_S(\mathbf{x}_1, \mathbf{x}_2)), \quad \text{for any } \mathbf{x}_1, \mathbf{x}_2 \in E.$$

Hence, $\omega_f^E(r) \leq \omega_g^S(r)$ for all $r \geq 0$, which implies $\omega_f^E(r) = \omega_g^S(r)$ since we have proved that $\omega_g^S(r) \leq \omega_f^E(r)$ for all $r \geq 0$. So we finish the proof. \square

Now we are ready to introduce and prove the main theorem of this section, which extends Theorem 1.1 to an irregular domain as follows.

THEOREM 4.3. *Let f be a uniformly continuous function defined on $E \subseteq [-R, R]^d$. For any arbitrary $L \in \mathbb{N}^+$, $N \in \mathbb{N}^+$, and $\eta > 0$, there exists a ReLU FNN ϕ with width $\max\{8d\lfloor N^{1/d} \rfloor + 4d, 12N + 14\}$ and depth $9L + 12$ such that*

$$\|f - \phi\|_{L^\infty(E \setminus \mathcal{H})} \leq 5\omega_f^E(16R\sqrt{d}N^{-2/d}L^{-2/d}), \quad \text{and} \quad \|\phi\|_{L^\infty(E)} \leq \|f\|_{L^\infty(E)} + 2\omega_f^E(2R\sqrt{d}),$$

where $\mathcal{H} \subseteq E$ with $\mu(\mathcal{H}) \leq \eta$.

As in Corollary 1.2, it is an immediate result to replace the $L^\infty(E \setminus \mathcal{H})$ -norm with the $L^p(E)$ -norm for $p \in [1, \infty)$. Now let us prove Theorem 4.3.

Proof of Theorem 4.3. By Lemma 4.2, f can be extended to \mathbb{R}^d such that

$$\omega_f^{\mathbb{R}^d}(r) = \omega_f^E(r), \quad \text{for any } r \geq 0.$$

Define

$$\tilde{f}(\mathbf{x}) := f(2R\mathbf{x} - R), \quad \text{for any } \mathbf{x} \in \mathbb{R}^d.$$

It follows that

$$(4.3) \quad \omega_{\tilde{f}}^{\mathbb{R}^d}(r) = \omega_{\tilde{f}}^{\mathbb{R}^d}(2Rr) = \omega_f^E(2Rr), \quad \text{for any } r \geq 0.$$

By Theorem 1.1, there exists a ReLU FNN $\tilde{\phi}$ with width $\max\{8d\lfloor N^{1/d} \rfloor + 4d, 12N + 14\}$ and depth $9L + 12$ such that

$$\|\tilde{f} - \tilde{\phi}\|_{L^\infty([0,1]^d \setminus \tilde{\mathcal{H}})} \leq 5\omega_{\tilde{f}}^{\mathbb{R}^d}(8\sqrt{d}N^{-2/d}L^{-2/d}), \quad \text{and} \quad \|\tilde{\phi}\|_{L^\infty([0,1]^d)} \leq |\tilde{f}(0)| + \omega_{\tilde{f}}^{\mathbb{R}^d}(\sqrt{d}),$$

for any $\tilde{\mathcal{H}} \in [0, 1]^d$ with $\mu(\tilde{\mathcal{H}}) \leq \frac{\eta}{(2R)^d}$. Define $\mathcal{H}_0 = 2R\tilde{\mathcal{H}} - R \subseteq [-R, R]^d$ and

$$\phi(\mathbf{x}) := \tilde{\phi}\left(\frac{1}{2R}\mathbf{x} + \frac{1}{2}\right), \quad \text{for any } \mathbf{x} \in \mathbb{R}^d.$$

Then for any $\mathbf{x} \in [-R, R]^d \setminus \mathcal{H}_0$, we have

$$\begin{aligned} |f(\mathbf{x}) - \phi(\mathbf{x})| &= |\tilde{f}\left(\frac{1}{2R}\mathbf{x} + \frac{1}{2}\right) - \tilde{\phi}\left(\frac{1}{2R}\mathbf{x} + \frac{1}{2}\right)| \\ &\leq \|\tilde{f} - \tilde{\phi}\|_{L^\infty([0,1]^d \setminus \tilde{\mathcal{H}})} \\ &\leq 5\omega_{\tilde{f}}^{\mathbb{R}^d}(8\sqrt{d}N^{-2/d}L^{-2/d}) \\ &= 5\omega_f^E(16R\sqrt{d}N^{-2/d}L^{-2/d}). \end{aligned}$$

Set $\mathcal{H} = E \cap \mathcal{H}_0$, then $\mu(\mathcal{H}) \leq \mu(\mathcal{H}_0) = (2R)^d \mu(\tilde{\mathcal{H}}) \leq \eta$ and $E \setminus \mathcal{H} \subseteq [-R, R]^d \setminus \mathcal{H}_0$, which implies

$$\|f - \phi\|_{L^\infty(E \setminus \mathcal{H})} \leq 5\omega_f^E(16R\sqrt{d}N^{-2/d}L^{-2/d}).$$

Besides,

$$\begin{aligned} \|\phi\|_{L^\infty(E)} &\leq \|\phi\|_{L^\infty([-R, R]^d)} = \|\tilde{\phi}\|_{L^\infty([0,1]^d)} \leq |\tilde{f}(0)| + \omega_{\tilde{f}}^{\mathbb{R}^d}(\sqrt{d}) \leq |\tilde{f}\left(\frac{1}{2R}\mathbf{x}_E + \frac{1}{2}\right)| + 2\omega_{\tilde{f}}^{\mathbb{R}^d}(\sqrt{d}) \\ &= |f(\mathbf{x}_E)| + 2\omega_f^E(2R\sqrt{d}) \\ &= \|f\|_{L^\infty(E)} + 2\omega_f^E(2R\sqrt{d}), \end{aligned}$$

where \mathbf{x}_E is a point in E . So we finish the proof. \square

4.2. Approximation in a Neighborhood of a Low-Dimensional Manifold. In this section, we study neural network approximation of functions defined in a neighborhood of a low-dimensional manifold and prove Theorem 1.3 in this setting. Let us first introduce Theorem 4.4 and Lemma 4.5 which are required to prove Theorem 1.3.

THEOREM 4.4 (Theorem 3.1 of [3]). *Let \mathcal{M} be a compact $d_{\mathcal{M}}$ -dimensional Riemannian submanifold of \mathbb{R}^d having condition number $1/\tau$, volume V , and geodesic covering regularity \mathcal{R} . Fix $\delta \in (0, 1)$ and $\gamma \in (0, 1)$. Let $\mathbf{A} = \sqrt{\frac{d}{d_\delta}}\Phi$, where $\Phi \in \mathbb{R}^{d_\delta \times d}$ is a random orthoprojector with*

$$d_\delta = \mathcal{O}\left(\frac{d_{\mathcal{M}} \ln(dV\mathcal{R}\tau^{-1}\delta^{-1}) \ln(1/\gamma)}{\delta^2}\right).$$

If $d_\delta \leq d$, then with probability at least $1 - \gamma$, the following statement holds: For every $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{M}$,

$$(1 - \delta)|\mathbf{x}_1 - \mathbf{x}_2| \leq |\mathbf{A}\mathbf{x}_1 - \mathbf{A}\mathbf{x}_2| \leq (1 + \delta)|\mathbf{x}_1 - \mathbf{x}_2|.$$

Theorem 4.4 shows the existence of a linear projector $\mathbf{A} \in \mathbb{R}^{d_\delta \times d}$ that maps a low-dimensional manifold in a high-dimensional space to a low-dimensional space nearly preserving distance. With this projection \mathbf{A} available, we can prove Theorem 1.3 via constructing a ReLU FNN defined in the low-dimensional space using Theorem 4.3 and hence the curse of dimensionality is lessened. The ideas of the proof are summarized in the following Table 1.

In Table 1 and the detailed proof later, we introduce a new notation $\mathcal{SL}(E)$ for any compact set $E \subseteq \mathbb{R}^d$ as the “smallest” element of E . Specifically, $\mathcal{SL}(E)$ is defined as the unique point in $\cap_{k=1}^d E_k$, where

$$E_k := \{\mathbf{x} \in E_{k-1} : x_k = s_k\}, \quad s_k := \inf\{x_k : [x_1, x_2, \dots, x_d]^T \in E_{k-1}\}, \quad \text{for } k = 1, 2, \dots, d,$$

and $E_0 = E$. The compactness of E ensures that $\cap_{k=1}^d E_k$ is in fact one point belonging to E . The introduction of $\mathcal{SL}(\cdot)$ uniquely formulates a low-dimensional function \tilde{f} representing a high-dimensional function f defined on \mathcal{M}_ε by

$$\tilde{f}(\mathbf{y}) := f(\mathbf{x}_\mathbf{y}), \quad \text{where } \mathbf{x}_\mathbf{y} = \mathcal{SL}(\{\mathbf{x} \in \mathcal{M}_\varepsilon : \mathbf{A}\mathbf{x} = \mathbf{y}\}), \quad \text{for any } \mathbf{y} \in \mathbf{A}(\mathcal{M}_\varepsilon) \subseteq \mathbb{R}^{d_\delta}.$$

As we shall see later, \tilde{f} can approximate f well because $\{\mathbf{x} \in \mathcal{M}_\varepsilon : \mathbf{A}\mathbf{x} = \mathbf{y}\}$ is contained in a small ball of radius $\mathcal{O}(\varepsilon)$ for any $\mathbf{y} \in \mathbf{A}(\mathcal{M}_\varepsilon)$. There are many other alternative ways to define $\mathcal{SL}(\cdot)$ as long as the definition ensures that $\mathcal{SL}(E)$ contains only one element. For example, $\mathcal{SL}(E)$ can be defined as any arbitrary point in E . For another example, $\mathbf{y} \in \mathbf{A}(\mathcal{M})$ cannot guarantee $\mathbf{x}_\mathbf{y} = \mathcal{SL}(\{\mathbf{x} \in \mathcal{M}_\varepsilon : \mathbf{A}\mathbf{x} = \mathbf{y}\}) \in \mathcal{M}$ in the current definition, but in practice we can choose $\mathcal{SL}(\{\mathbf{x} \in \mathcal{M} : \mathbf{A}\mathbf{x} = \mathbf{y}\})$ as $\mathbf{x}_\mathbf{y}$ to ensure that $\mathbf{x}_\mathbf{y} \in \mathcal{M}$, which might be beneficial for potential applications.

TABLE 1

Main steps of the proof of Theorem 1.3. Step 1: dimension reduction via the nearly isometric projection operator \mathbf{A} provided by Theorem 4.4 to obtain an “equivalent” function \tilde{f} of f in a low-dimensional domain using $\mathbf{x}_\mathbf{y} = \mathcal{SL}(\{\mathbf{x} \in \mathcal{M}_\varepsilon : \mathbf{A}\mathbf{x} = \mathbf{y}\})$. Step 2: construct a ReLU FNN $\tilde{\phi} \approx \tilde{f}$ by Theorem 4.3. Step 3: define a ReLU FNN ϕ in the original high-dimensional domain via the projection \mathbf{A} . Step 4: verify that the approximation error of $\phi \approx f$ satisfies our requirement.

$f(\mathbf{x})$ for $\mathbf{x} \in \mathcal{M}_\varepsilon \subseteq [0, 1]^d$	Step 4 \approx	$\phi(\mathbf{x}) := \tilde{\phi}(\mathbf{A}\mathbf{x})$ for $\mathbf{x} \in \mathcal{M}_\varepsilon \subseteq [0, 1]^d$
Step 1 \Downarrow $\mathbf{x}_\mathbf{y} = \mathcal{SL}(\{\mathbf{x} \in \mathcal{M}_\varepsilon : \mathbf{A}\mathbf{x} = \mathbf{y}\})$		Step 3 \Downarrow $\mathbf{y} = \mathbf{A}\mathbf{x}$
$\tilde{f}(\mathbf{y}) := f(\mathbf{x}_\mathbf{y})$ for $\mathbf{y} \in \mathbf{A}(\mathcal{M}_\varepsilon) \subseteq \mathbb{R}^{d_\delta}$	Step 2 \approx	$\tilde{\phi}(\mathbf{y})$ for $\mathbf{y} \in \mathbf{A}(\mathcal{M}_\varepsilon) \subseteq \mathbb{R}^{d_\delta}$

Note that in Step 2 in Table 1, the approximation $\tilde{f} \approx \tilde{\phi}$ is only valid in the L^∞ -norm outside a “don’t-care” region \mathcal{H} with an arbitrarily small measure $\mu(\mathcal{H})$. Hence, when we verify the approximation error in Step 4, it is important to ensure that $\mu(\{\mathbf{x} \in \mathcal{M}_\varepsilon : \mathbf{A}\mathbf{x} \in \mathcal{H}\})$ is well-controlled. Hence, Lemma 4.5 below is introduced to characterize the behavior of the Lebesgue measure under a linear transform. For any set $E \subseteq \mathbb{R}^d$ and a matrix $\mathbf{A} \in \mathbb{R}^{d' \times d}$, we use the notation $\mathbf{A}(E) := \{\mathbf{A}\mathbf{x} : \mathbf{x} \in E\}$.

LEMMA 4.5. Given $R > 0$ and $d', d \in \mathbb{N}^+$ with $d' \leq d$, let $\mathbf{A} \in \mathbb{R}^{d' \times d}$ be a row independent matrix, then

$$\mu(\{\mathbf{x} \in \mathbb{R}^d : \mathbf{A}\mathbf{x} \in \mathcal{H}\} \cap [-R, R]^d) \rightarrow 0 \quad \text{if} \quad \mu(\mathcal{H}) \rightarrow 0,$$

for any Lebesgue measurable set $\mathcal{H} \subseteq \mathbb{R}^{d'}$.

Proof of Lemma 4.5. Since $\mathbf{A} \in \mathbb{R}^{d' \times d}$ is row independent, there exists $\mathbf{A}_0 \in \mathbb{R}^{(d-d') \times d}$ such that $\tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{A} \\ \mathbf{A}_0 \end{bmatrix}$ is nonsingular, which implies

$$\mathbf{A} = \mathbf{I}_{d'} \mathbf{A} = [\mathbf{I}_{d'}, \mathbf{0}] \begin{bmatrix} \mathbf{A} \\ \mathbf{A}_0 \end{bmatrix} = [\mathbf{I}_{d'}, \mathbf{0}] \tilde{\mathbf{A}}.$$

Define

$$\tilde{\mathcal{H}} := \{\mathbf{x} \in \mathbb{R}^d : [\mathbf{I}_{d'}, \mathbf{0}] \mathbf{x} \in \mathcal{H}\} = \{\mathbf{x} = [x_1, x_2, \dots, x_d]^T \in \mathbb{R}^d : [x_1, \dots, x_{d'}]^T \in \mathcal{H}\} = \mathcal{H} \otimes \mathbb{R}^{d-d'}.$$

It follows that

$$\begin{aligned} \{\mathbf{x} : \mathbf{A}\mathbf{x} \in \mathcal{H}\} \cap [-R, R]^d &= \{\mathbf{x} : [\mathbf{I}_{d'}, \mathbf{0}] \tilde{\mathbf{A}}\mathbf{x} \in \mathcal{H}\} \cap [-R, R]^d \\ &= \tilde{\mathbf{A}}^{-1}(\{\mathbf{y} : [\mathbf{I}_{d'}, \mathbf{0}] \mathbf{y} \in \mathcal{H}\} \cap \tilde{\mathbf{A}}([-R, R]^d)) \\ &= \tilde{\mathbf{A}}^{-1}(\tilde{\mathcal{H}} \cap \tilde{\mathbf{A}}([-R, R]^d)). \end{aligned}$$

There exists $R' > 0$ such that $\tilde{\mathbf{A}}([-R, R]^d) \subseteq [-R', R']^d$, then

$$\tilde{\mathcal{H}} \cap \tilde{\mathbf{A}}([-R, R]^d) \subseteq \tilde{\mathcal{H}} \cap [-R', R']^d = (\mathcal{H} \otimes \mathbb{R}^{d-d'}) \cap [-R', R']^d \subseteq \mathcal{H} \otimes [-R', R']^{d-d'},$$

which means

$$\mu(\tilde{\mathcal{H}} \cap \tilde{\mathbf{A}}([-R, R]^d)) \rightarrow 0 \quad \text{if} \quad \mu(\mathcal{H}) \rightarrow 0.$$

Since $\tilde{\mathbf{A}}^{-1}$ is nonsingular,

$$\mu(\tilde{\mathbf{A}}^{-1}(\tilde{\mathcal{H}} \cap \tilde{\mathbf{A}}([-R, R]^d))) \rightarrow 0 \quad \text{if} \quad \mu(\tilde{\mathcal{H}} \cap \tilde{\mathbf{A}}([-R, R]^d)) \rightarrow 0.$$

Therefore, we have

$$\mu(\{\mathbf{x} : \mathbf{A}\mathbf{x} \in \mathcal{H}\} \cap [-R, R]^d) \rightarrow 0 \quad \text{if} \quad \mu(\mathcal{H}) \rightarrow 0.$$

So we finish the proof. \square

Now we are ready to prove Theorem 1.3.

Proof of Theorem 1.3. By Theorem 4.4, there exists a matrix $\mathbf{A} \in \mathbb{R}^{d_\delta \times d}$ such that

$$(4.4) \quad \mathbf{A}\mathbf{A}^T = \frac{d}{d_\delta} \mathbf{I}_{d_\delta},$$

where \mathbf{I}_{d_δ} is an identity matrix of size $d_\delta \times d_\delta$, and

$$(4.5) \quad (1 - \delta)|\mathbf{x}_1 - \mathbf{x}_2| \leq |\mathbf{A}\mathbf{x}_1 - \mathbf{A}\mathbf{x}_2| \leq (1 + \delta)|\mathbf{x}_1 - \mathbf{x}_2|, \quad \text{for any } \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{M}.$$

Given any $\mathbf{y} \in \mathbf{A}(\mathcal{M}_\varepsilon)$, then $\{\mathbf{x} \in \mathcal{M}_\varepsilon : \mathbf{A}\mathbf{x} = \mathbf{y}\}$ is a nonzero compact set. Let $\mathbf{x}_\mathbf{y} = \mathcal{SL}(\{\mathbf{x} \in \mathcal{M}_\varepsilon : \mathbf{A}\mathbf{x} = \mathbf{y}\})$, then we define \tilde{f} on $\mathbf{A}(\mathcal{M}_\varepsilon)$ as $\tilde{f}(\mathbf{y}) = f(\mathbf{x}_\mathbf{y})$.

For any $\mathbf{y}_1, \mathbf{y}_2 \in \mathbf{A}(\mathcal{M}_\varepsilon)$, let $\mathbf{x}_i = \mathcal{SL}(\{\mathbf{x} \in \mathcal{M}_\varepsilon : \mathbf{A}\mathbf{x} = \mathbf{y}_i\})$, then $\mathbf{x}_i \in \mathcal{M}_\varepsilon$ for $i = 1, 2$. By the definition of \mathcal{M}_ε , there exist $\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2 \in \mathcal{M}$ such that $|\tilde{\mathbf{x}}_i - \mathbf{x}_i| \leq \varepsilon$ for $i = 1, 2$. It follows that

$$|\tilde{f}(\mathbf{y}_1) - \tilde{f}(\mathbf{y}_2)| = |f(\mathbf{x}_1) - f(\mathbf{x}_2)| \leq \omega_f(|\mathbf{x}_1 - \mathbf{x}_2|) \leq \omega_f(|\tilde{\mathbf{x}}_1 - \tilde{\mathbf{x}}_2| + 2\varepsilon) \leq \omega_f\left(\frac{1}{1-\delta}|\mathbf{A}\tilde{\mathbf{x}}_1 - \mathbf{A}\tilde{\mathbf{x}}_2| + 2\varepsilon\right),$$

where the last inequality comes from (4.5). By the triangular inequality, we have

$$\begin{aligned} |\tilde{f}(\mathbf{y}_1) - \tilde{f}(\mathbf{y}_2)| &\leq \omega_f\left(\frac{1}{1-\delta}|\mathbf{A}\mathbf{x}_1 - \mathbf{A}\mathbf{x}_2| + \frac{1}{1-\delta}|\mathbf{A}\mathbf{x}_1 - \mathbf{A}\tilde{\mathbf{x}}_1| + \frac{1}{1-\delta}|\mathbf{A}\mathbf{x}_2 - \mathbf{A}\tilde{\mathbf{x}}_2| + 2\varepsilon\right) \\ &\leq \omega_f\left(\frac{1}{1-\delta}|\mathbf{A}\mathbf{x}_1 - \mathbf{A}\mathbf{x}_2| + \frac{2\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}} + 2\varepsilon\right) \\ &\leq \omega_f\left(\frac{1}{1-\delta}|\mathbf{y}_1 - \mathbf{y}_2| + \frac{2\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}} + 2\varepsilon\right). \end{aligned}$$

Set $\omega(r) = \omega_f\left(\frac{1}{1-\delta}r\right)$ for any $r \geq 0$ and $\Delta = 2\varepsilon\sqrt{\frac{d}{d_\delta}} + 2\varepsilon(1-\delta)$, then

$$|\tilde{f}(\mathbf{y}_1) - \tilde{f}(\mathbf{y}_2)| \leq \omega(|\mathbf{y}_1 - \mathbf{y}_2| + \Delta), \quad \text{for any } \mathbf{y}_1, \mathbf{y}_2 \in \mathbf{A}(\mathcal{M}_\varepsilon) \subseteq \mathbb{R}^{d_\delta}.$$

By Lemma 4.1, there exists \tilde{g} defined on \mathbb{R}^{d_δ} such that

$$(4.6) \quad |\tilde{g}(\mathbf{y}) - \tilde{f}(\mathbf{y})| \leq \omega(\Delta) = \omega_f\left(\frac{2\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}} + 2\varepsilon\right), \quad \text{for any } \mathbf{y} \in \mathbf{A}(\mathcal{M}_\varepsilon),$$

and

$$|\tilde{g}(\mathbf{y}_1) - \tilde{g}(\mathbf{y}_2)| \leq \omega(|\mathbf{y}_1 - \mathbf{y}_2|) = \omega_f\left(\frac{1}{1-\delta}|\mathbf{y}_1 - \mathbf{y}_2|\right), \quad \text{for any } \mathbf{y}_1, \mathbf{y}_2 \in \mathbb{R}^{d_\delta}.$$

It follows that

$$(4.7) \quad \omega_{\tilde{g}}^{\mathbb{R}^{d_\delta}}(r) \leq \omega_f\left(\frac{r}{1-\delta}\right), \quad \text{for any } r \geq 0.$$

By Equation (4.4) and the definition of \mathcal{M}_ε in Equation (1.2), it is easy to check that

$$\mathbf{A}(\mathcal{M}_\varepsilon) \subseteq \mathbf{A}([0, 1]^d) \subseteq [-\sqrt{\frac{d}{d_\delta}}, \sqrt{\frac{d}{d_\delta}}]^{d_\delta}.$$

By the application of Theorem 4.3 with $E = [-\sqrt{\frac{d}{d_\delta}}, \sqrt{\frac{d}{d_\delta}}]^{d_\delta}$, there exists a ReLU FNN $\tilde{\phi}$ with width $\max\{8d_\delta\lfloor N^{1/d_\delta} \rfloor + 4d_\delta, 12N + 14\}$ and depth $9L + 12$ such that

$$(4.8) \quad \|\tilde{g} - \tilde{\phi}\|_{L^\infty(E \setminus \tilde{\mathcal{H}})} \leq 5\omega_{\tilde{g}}^E\left(\frac{16d}{\sqrt{d_\delta}}N^{-2/d_\delta}L^{-2/d_\delta}\right) \quad \text{and} \quad \|\tilde{\phi}\|_{L^\infty(E)} \leq \|\tilde{g}\|_{L^\infty(E)} + 2\omega_{\tilde{g}}^E\left(\frac{2d}{\sqrt{d_\delta}}\right),$$

where $\tilde{\mathcal{H}} \subseteq E$ with $\mu(\tilde{\mathcal{H}}) \leq \eta$, and $\eta > 0$ is a small number to be determined later.

Define $\phi := \tilde{\phi} \circ \mathbf{A}$, i.e., $\phi(\mathbf{x}) := \tilde{\phi}(\mathbf{A}\mathbf{x})$ for any $\mathbf{x} \in \mathbb{R}^d$. Then ϕ is also a ReLU FNN with width $\max\{8d_\delta\lfloor N^{1/d_\delta} \rfloor + 4d_\delta, 12N + 14\}$ and depth $9L + 12$. Let $\mathcal{H} := \{\mathbf{x} : \mathbf{A}\mathbf{x} \in \tilde{\mathcal{H}}\} \cap \mathcal{M}_\varepsilon$. By Lemma 4.5,

$$(4.9) \quad \eta \rightarrow 0 \implies \mu(\tilde{\mathcal{H}}) \rightarrow 0 \implies \mu(\mathcal{H}) \rightarrow 0 \implies \mu_\varrho(\mathcal{H}) \rightarrow 0.$$

Therefore,

$$\begin{aligned} \|f - \phi\|_{L^p([0, 1]^d, \mu_\varrho)} &= \left(\int_{\mathcal{M}_\varepsilon} |f(\mathbf{x}) - \phi(\mathbf{x})|^p d\mu_\varrho(\mathbf{x}) \right)^{1/p} \\ &\leq \left(\int_{\mathcal{M}_\varepsilon \setminus \mathcal{H}} |f(\mathbf{x}) - \phi(\mathbf{x})|^p d\mu_\varrho(\mathbf{x}) \right)^{1/p} + \left(\int_{\mathcal{H}} |f(\mathbf{x}) - \phi(\mathbf{x})|^p d\mu_\varrho(\mathbf{x}) \right)^{1/p} := \mathcal{S}_1 + \mathcal{S}_2. \end{aligned}$$

First, let us estimate \mathcal{S}_1 . For any $\mathbf{x} \in \mathcal{M}_\varepsilon \setminus \mathcal{H}$, set $\mathbf{y} = \mathbf{A}\mathbf{x}$ and $\mathbf{x}_y = \mathcal{SL}(\{\mathbf{z} \in \mathbb{R}^d : \mathbf{A}\mathbf{z} = \mathbf{y}\})$, there exist $\tilde{\mathbf{x}}, \tilde{\mathbf{x}}_y \in \mathcal{M}$ such that $|\tilde{\mathbf{x}} - \mathbf{x}| \leq \varepsilon$ and $|\tilde{\mathbf{x}}_y - \mathbf{x}_y| \leq \varepsilon$. It follows that

$$\begin{aligned} |\mathbf{x} - \mathbf{x}_y| &\leq |\tilde{\mathbf{x}} - \tilde{\mathbf{x}}_y| + 2\varepsilon \leq \frac{1}{1-\delta}|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{A}\tilde{\mathbf{x}}_y| + 2\varepsilon \\ &\leq \frac{1}{1-\delta}(|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{A}\mathbf{x}| + |\mathbf{A}\mathbf{x} - \mathbf{A}\mathbf{x}_y| + |\mathbf{A}\mathbf{x}_y - \mathbf{A}\tilde{\mathbf{x}}_y|) + 2\varepsilon \\ (4.10) \quad &= \frac{1}{1-\delta}(|\mathbf{A}\tilde{\mathbf{x}} - \mathbf{A}\mathbf{x}| + |\mathbf{A}\mathbf{x}_y - \mathbf{A}\tilde{\mathbf{x}}_y|) + 2\varepsilon \\ &\leq \frac{2\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}} + 2\varepsilon. \end{aligned}$$

In fact, the above equation implies that $\{\mathbf{x} \in \mathcal{M}_\varepsilon : \mathbf{A}\mathbf{x} = \mathbf{y}\}$ is contained in a small ball of radius $\mathcal{O}(\varepsilon)$ for any $\mathbf{y} \in \mathbf{A}(\mathcal{M}_\varepsilon)$ as we mentioned previously.

Together with Equation (4.6), we have

$$\begin{aligned}
|f(\mathbf{x}) - \phi(\mathbf{x})| &\leq |f(\mathbf{x}) - f(\mathbf{x}_\mathbf{y})| + |f(\mathbf{x}_\mathbf{y}) - \phi(\mathbf{x})| \\
&\leq \omega_f\left(\frac{2\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}} + 2\varepsilon\right) + |\tilde{f}(\mathbf{y}) - \tilde{\phi}(\mathbf{y})| \\
&\leq \omega_f\left(\frac{2\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}} + 2\varepsilon\right) + |\tilde{f}(\mathbf{y}) - \tilde{g}(\mathbf{y})| + |\tilde{g}(\mathbf{y}) - \tilde{\phi}(\mathbf{y})| \\
&\leq \omega_f\left(\frac{2\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}} + 2\varepsilon\right) + \omega_f\left(\frac{2\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}} + 2\varepsilon\right) + 5\omega_g^E\left(\frac{16d}{\sqrt{d_\delta}}N^{-2/d_\delta}L^{-2/d_\delta}\right) \\
&\leq 2\omega_f\left(\frac{2\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}} + 2\varepsilon\right) + 5\omega_f\left(\frac{16d}{(1-\delta)\sqrt{d_\delta}}N^{-2/d_\delta}L^{-2/d_\delta}\right).
\end{aligned}$$

It follows that

$$\mathcal{S}_1 \leq 2\omega_f\left(\frac{2\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}} + 2\varepsilon\right) + 5\omega_f\left(\frac{16d}{(1-\delta)\sqrt{d_\delta}}N^{-2/d_\delta}L^{-2/d_\delta}\right).$$

Then we focus on \mathcal{S}_2 . By Equation (4.8) and the definition of $\phi = \tilde{\phi} \circ \mathbf{A}$,

$$\begin{aligned}
\|f - \phi\|_{L^\infty(\mathcal{M}_\varepsilon)} &\leq \|f\|_{L^\infty(\mathcal{M}_\varepsilon)} + \|\phi\|_{L^\infty(\mathcal{M}_\varepsilon)} \\
&= \|f\|_{L^\infty(\mathcal{M}_\varepsilon)} + \|\tilde{\phi}\|_{L^\infty(\mathbf{A}(\mathcal{M}_\varepsilon))} \\
&\leq \|f\|_{L^\infty(\mathcal{M}_\varepsilon)} + \|\tilde{g}\|_{L^\infty(E)} + 2\omega_g^E\left(\frac{2d}{\sqrt{d_\delta}}\right) \\
&\leq \|f\|_{L^\infty(\mathcal{M}_\varepsilon)} + \|\tilde{g}\|_{L^\infty(E)} + 2\omega_f\left(\frac{2d}{(1-\delta)\sqrt{d_\delta}}\right).
\end{aligned}$$

Hence,

$$\mathcal{S}_2 \leq \left(\|f\|_{L^\infty(\mathcal{M}_\varepsilon)} + \|\tilde{g}\|_{L^\infty(E)} + 2\omega_f\left(\frac{2d}{(1-\delta)\sqrt{d_\delta}}\right) \right)^p \mu_\varrho(\mathcal{H})^{1/p}.$$

Together with Equation (4.9), we know

$$\eta \rightarrow 0 \implies \mu_\varrho(\mathcal{H}) \rightarrow 0 \implies \mathcal{S}_2 \rightarrow 0.$$

So we can choose a sufficiently small $\eta > 0$ such that

$$\mathcal{S}_2 \leq \omega_f\left(\frac{2\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}} + 2\varepsilon\right).$$

Therefore,

$$\begin{aligned}
\|f - \phi\|_{L^p([0,1]^d, \mu_\varrho)} &= \mathcal{S}_1 + \mathcal{S}_2 \leq 3\omega_f\left(\frac{2\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}} + 2\varepsilon\right) + 5\omega_f\left(\frac{16d}{(1-\delta)\sqrt{d_\delta}}N^{-2/d_\delta}L^{-2/d_\delta}\right) \\
&\leq 3\omega_f\left(\frac{4\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}}\right) + 5\omega_f\left(\frac{16d}{(1-\delta)\sqrt{d_\delta}}N^{-2/d_\delta}L^{-2/d_\delta}\right).
\end{aligned}$$

Hence, we have finished the proof of this theorem. \square

It is worth emphasizing that the approximation error $\mathcal{O}\left(\omega_f(\mathcal{O}(\varepsilon)) + \omega_f(\mathcal{O}(N^{-2/d_\delta}L^{-2/d_\delta}))\right)$ in Theorem 1.3 is equal to $\mathcal{O}\left(\omega_f(\mathcal{O}(N^{-2/d_\delta}L^{-2/d_\delta}))\right)$ when $\varepsilon \leq N^{-2/d_\delta}L^{-2/d_\delta}$.

The application of Theorem 4.4 and the proof of Theorem 1.3 in fact inspire an efficient two-step algorithm for high-dimensional learning problems: in the first step, high-dimensional data are projected to a low-dimensional space via a random projection; in the second step, a deep learning algorithm is applied to learn from the low-dimensional data. By Theorem 4.4 and 1.3, the deep learning algorithm in the low-dimensional space can still provide good results with a high probability.

4.3. Optimal ReLU FNN Structure in Parallel Computing. In this section, we show how to select the best ReLU FNN to approximate functions in $\text{Lip}(\nu, \alpha, d)$ on a d -dimensional cube, if the approximation error ε and the number of parallel computing cores (processors) p are given. We choose the best ReLU FNN by minimizing the time complexity in each training iteration. The analysis in this section is valid up to a constant prefactor.

Assume $\phi_{\theta} \in \text{NN}(\# \text{input} = d; \text{widthvec} = [N]^L)$, $N, L \in \mathbb{N}^+$, where θ is the vector including all parameters of ϕ_{θ} . By the basic knowledge of parallel computing (see [33] for more details), we have the following Table 2.

TABLE 2
Time complexity of one training iteration for an FNN of width N and depth L .

Number of cores p	Time Complexity	
	Evaluating $\phi_{\theta}(\mathbf{x})$	Evaluating $\frac{\partial \phi_{\theta}(\mathbf{x})}{\partial \theta}$
$p \in [1, N]$	$\mathcal{O}(N^2 L/p)$	$\mathcal{O}(N^2 L/p)$
$p \in (N, N^2]$	$\mathcal{O}(L(N^2/p + \ln \frac{p}{N}))$	$\mathcal{O}(L(N^2/p + \ln \frac{p}{N}))$
$p \in (N^2, \infty)$	$\mathcal{O}(L \ln N)$	$\mathcal{O}(L \ln N)$

For the sake of simplicity, we assume that the training batch size is $\mathcal{O}(1)$. Denote the time complexity of each training iteration as $T(n, L)$, then

$$T(N, L) = \begin{cases} \mathcal{O}(N^2 L/p), & p \in [1, N], \\ \mathcal{O}(L(N^2/p + \ln \frac{p}{N})), & p \in (N, N^2], \\ \mathcal{O}(L \ln N), & p \in (N^2, \infty). \end{cases}$$

Theorem 3.6 and 3.7 imply that the approximation error ε is essentially $\mathcal{O}((NL)^{-2\alpha/d})$. Hence, we can get the optimal size of ReLU FNNs via the optimization problem below:

$$(4.11) \quad \begin{aligned} (N_{\text{opt}}, L_{\text{opt}}) &= \arg \min_{N, L} T(N, L) \\ \text{subject to } &\begin{cases} \varepsilon = \mathcal{O}((NL)^{-2\alpha/d}), \\ N, L, p \in \mathbb{N}^+. \end{cases} \end{aligned}$$

To simplify the discussion, we have the following assumptions:

- Dropping the notation $\mathcal{O}(\cdot)$ sometimes while assuming asymptotic analysis with the abuse of notations.
- N , L , and p are allowed to be real numbers.
- Replacing $\mathcal{O}((NL)^{-(2\alpha/d+\rho)}) \leq \varepsilon \leq \mathcal{O}((NL)^{-2\alpha/d})$ with $\varepsilon = (NL)^{-2\alpha/d}$ since $\rho > 0$ is arbitrary.

With $\varepsilon = (NL)^{-2\alpha/d}$, we have

$$(4.12) \quad \begin{aligned} \bar{T}(N, L) &:= \begin{cases} N^2 L/p, & p \in [1, N], \\ L(N^2/p + \ln \frac{p}{N}), & p \in (N, N^2], \\ L(1 + \ln N), & p \in [N^2, \infty), \end{cases} \\ &= \begin{cases} N\varepsilon^{-d/(2\alpha)}/p, & N \in [p, \infty), \\ N\varepsilon^{-d/(2\alpha)}/p + \frac{1}{N}\varepsilon^{-d/(2\alpha)} \ln \frac{p}{N}, & N \in [\sqrt{p}, p], \\ \frac{1+\ln N}{N}\varepsilon^{-d/(2\alpha)}, & N \in [1, \sqrt{p}). \end{cases} \end{aligned}$$

Then we get $T(N, L) = \mathcal{O}(\bar{T}(N, L))$. Therefore, the optimization problem in (4.11) can be simplified to

$$(4.13) \quad \begin{aligned} (N_{\text{opt}}, L_{\text{opt}}) &= \arg \min_{N, L} \bar{T}(N, L) \\ \text{subject to } &\begin{cases} \varepsilon = (NL)^{-2\alpha/d}, \\ N, L, p \in [1, \infty). \end{cases} \end{aligned}$$

By (4.12), $\bar{T}(N, L)$ is independent of L on the condition that $\varepsilon = (NL)^{-2\alpha/d}$. Therefore, we may denote $\bar{T}(N, L)$ by $\bar{T}(N)$. Now we consider two cases: the case $p = \mathcal{O}(1)$ and the case $p \gg \mathcal{O}(1)$.

Case 1: The case $p = \mathcal{O}(1)$.

It is clear that $\bar{T}(N)$ is increasing in N when $N \in [p, \infty)$ by (4.12). Together with $p = \mathcal{O}(1)$, then $\mathcal{O}(\sqrt{p}) = \mathcal{O}(p) = \mathcal{O}(1)$. Therefore, $N_{\text{opt}} = \mathcal{O}(1)$ and $L_{\text{opt}} = \mathcal{O}(\varepsilon^{-d/(2\alpha)})$. Note that we regard d as a constant ($\mathcal{O}(1)$) in above analysis, N_{opt} should be $\mathcal{O}(d)$ in fact.

Case 2: The case $p \gg \mathcal{O}(1)$.

Since $\varepsilon = (NL)^{-2\alpha/d}$, we have $N \leq \varepsilon^{-d/(2\alpha)}$. We only need to consider the monotonicity of $\bar{T}(N)$ on $[1, \varepsilon^{-d/(2\alpha)}]$. Together with (4.12), this case can be divided into two sub-cases: the sub-case $\sqrt{p} \leq \varepsilon^{-d/(2\alpha)}$ and the sub-case $\sqrt{p} > \varepsilon^{-d/(2\alpha)}$.

Case 2.1: The sub-case $\sqrt{p} > \varepsilon^{-d/(2\alpha)}$.

$\sqrt{p} > \varepsilon^{-d/(2\alpha)}$ implies $[1, \varepsilon^{-d/(2\alpha)}] \subseteq [1, \sqrt{q}]$. Hence, $\bar{T}(N)$ is decreasing in N on $[1, \varepsilon^{-d/(2\alpha)}]$. It follows that $N_{\text{opt}} = \mathcal{O}(\varepsilon^{-d/(2\alpha)})$ and that $L_{\text{opt}} = \mathcal{O}(1)$.

Case 2.2: The sub-case $\sqrt{p} \leq \varepsilon^{-d/(2\alpha)}$.

For this sub-case, N_{opt} and N_{opt} are hard to estimate. However, we can give a rough range of N_{opt} . Since $\bar{T}(N)$ is decreasing in N on $[1, \sqrt{p}]$ and increasing in N on $[p, \infty)$, the minimum of $\bar{T}(N)$ is achieved on $[\sqrt{p}, p]$. Hence, $N_{\text{opt}} \in [\mathcal{O}(\sqrt{p}), \mathcal{O}(p)] \cap [\mathcal{O}(\sqrt{p}), \mathcal{O}(\varepsilon^{-d/(2\alpha)})]$ and $L_{\text{opt}} = \mathcal{O}(\varepsilon^{-d/(2\alpha)}/N_{\text{opt}})$.

5. Conclusion and future work. This paper aims at a quantitative and optimal approximation rate of ReLU FNNs in terms of both width and depth simultaneously to approximate continuous functions. It was shown that ReLU FNNs with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ can approximate an arbitrary continuous function on a d -dimensional cube with an approximation rate $5\omega_f(8\sqrt{d}N^{-2/d}L^{-2/d})$. In particular, when f is a Hölder continuous function of order α with a Lipschitz constant ν , the approximation rate is $5(8\sqrt{d})^\alpha \nu N^{-2\alpha/d} L^{-2\alpha/d}$ and it is asymptotically tight. We also extended our analysis to the case when the domain of f is irregular and showed the same approximation rate. In practical applications, it is usually believed that real data are sampled from an ε -neighborhood of a $d_{\mathcal{M}}$ -dimensional smooth manifold $\mathcal{M} \subseteq [0, 1]^d$ with $d_{\mathcal{M}} \ll d$. In the case of an essentially low-dimensional domain, we show an approximation rate $3\omega_f\left(\frac{4\varepsilon}{1-\delta}\sqrt{\frac{d}{d_\delta}}\right) + 5\omega_f\left(\frac{16d}{(1-\delta)\sqrt{d_\delta}}N^{-2/d_\delta}L^{-2/d_\delta}\right)$ for ReLU FNNs to approximate f in the ε -neighborhood, $d_\delta = \mathcal{O}\left(d_{\mathcal{M}}\frac{\ln(d/\delta)}{\delta^2}\right)$ for any given $\delta \in (0, 1)$.

Besides, we studied how to select the best ReLU FNN to approximate continuous function in parallel computing. In particular, ReLU FNNs with depth $\mathcal{O}(1)$ are the best choices if the number of parallel computing cores p is sufficiently large. ReLU FNNs with width $\mathcal{O}(d)$ are best choices if $p = \mathcal{O}(1)$. The width of best ReLU FNNs is between $\mathcal{O}(\sqrt{p})$ and $\mathcal{O}(p)$ if p is moderate.

We would like to remark that our analysis was based on the fully connected feed-forward neural networks and the ReLU activation function. It would be very interesting to generalize our conclusions to neural networks with other types of architectures (e.g., convolutional neural networks) and activation functions (e.g., tanh and sigmoid functions). Another important direction is to extend our results to the L^∞ -norm on the whole domain without the “don’t-care” region. Besides, if identity maps are allowed in the construction of neural networks as in the residual networks [26], the size of FNNs in our construction can be further optimized. Finally, the proposed analysis could be generalized to other function spaces with explicit formulas to characterize the approximation error. These will be left as future work.

Acknowledgments. H. Yang thanks the support of the start-up grant by the Department of Mathematics at the National University of Singapore, the support of the Ministry of Education in Singapore under the grant MOE2018-T2-2-147.

REFERENCES

- [1] O. ABDEL-HAMID, A. MOHAMED, H. JIANG, L. DENG, G. PENN, AND D. YU, *Convolutional neural networks for speech recognition*, IEEE/ACM Transactions on Audio, Speech, and Language Processing, 22 (2014), pp. 1533–1545, <https://doi.org/10.1109/TASLP.2014.2339736>.
- [2] M. ANTHONY AND P. L. BARTLETT, *Neural Network Learning: Theoretical Foundations*, Cambridge University Press, New York, NY, USA, 1st ed., 2009.
- [3] R. G. BARANIUK AND M. B. WAKIN, *Random projections of smooth manifolds*, Foundations of Computational Mathematics, 9 (2009), pp. 51–77, <https://doi.org/10.1007/s10208-007-9011-z>, <https://doi.org/10.1007/s10208-007-9011-z>.
- [4] A. R. BARRON, *Universal approximation bounds for superpositions of a sigmoidal function*, IEEE Transactions on Information Theory, 39 (1993), pp. 930–945, <https://doi.org/10.1109/18.256500>.
- [5] P. BARTLETT, V. MAIOROV, AND R. MEIR, *Almost linear VC dimension bounds for piecewise polynomial networks*, Neural Computation, 10 (1998), pp. 217–3.
- [6] M. BIANCHINI AND F. SCARSELLI, *On the complexity of neural network classifiers: A comparison between shallow and deep architectures*, IEEE Transactions on Neural Networks and Learning Systems, 25 (2014), pp. 1553–1565, <https://doi.org/10.1109/TNNLS.2013.2293637>.
- [7] E. K. BLUM AND L. K. LI, *Approximation theory and feedforward networks*, Neural Networks, 4 (1991), pp. 511 – 515, [https://doi.org/https://doi.org/10.1016/0893-6080\(91\)90047-9](https://doi.org/https://doi.org/10.1016/0893-6080(91)90047-9), <http://www.sciencedirect.com/science/article/pii/0893608091900479>.
- [8] D. S. BROOMHEAD AND D. LOWE, *Multivariable Functional Interpolation and Adaptive Networks*, Complex Systems 2, (1988), pp. 321–355.
- [9] J. CAI, D. LI, J. SUN, AND K. WANG, *Enhanced expressive power and fast training of neural networks by random projections*, CoRR, abs/1811.09054 (2018), <http://arxiv.org/abs/1811.09054>, <https://arxiv.org/abs/1811.09054>.
- [10] S. CHEN AND D. DONOHO, *Basis pursuit*, in Proceedings of 1994 28th Asilomar Conference on Signals, Systems and Computers, vol. 1, Oct 1994, pp. 41–44 vol.1, <https://doi.org/10.1109/ACSSC.1994.471413>.
- [11] C. K. CHUI, S.-B. LIN, AND D.-X. ZHOU, *Construction of neural networks for realization of localized deep learning*.
- [12] D. C. CIREŞAN, U. MEIER, J. MASCI, L. M. GAMBARELLA, AND J. SCHMIDHUBER, *Flexible, high performance convolutional neural networks for image classification*, in Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence - Volume Volume Two, IJCAI’11, AAAI Press, 2011, pp. 1237–1242, <https://doi.org/10.5591/978-1-57735-516-8/IJCAI11-210>, <http://dx.doi.org/10.5591/978-1-57735-516-8/IJCAI11-210>.
- [13] D. COSTARELLI AND A. R. SAMBUCINI, *Saturation classes for max-product neural network operators activated by sigmoidal functions*, Results in Mathematics, 72 (2017), pp. 1555 – 1569, <https://doi.org/10.1007/s00025-017-0692-6>.
- [14] D. COSTARELLI AND G. VINTI, *Convergence for a family of neural network operators in orlicz spaces*, Mathematische Nachrichten, 290 (2017), pp. 226–235, <https://onlinelibrary.wiley.com/doi/abs/10.1002/mana.201600006>.
- [15] D. COSTARELLI AND G. VINTI, *Approximation results in orlicz spaces for sequences of kantorovich max-product neural network operators*, Results in Mathematics, 73 (2018), pp. 1 – 15, <https://doi.org/10.1007/s00025-018-0799-4>.
- [16] G. CYBENKO, *Approximation by superpositions of a sigmoidal function*, MCSS, 2 (1989), pp. 303–314.
- [17] I. DAUBECHIES, R. DEVORE, S. FOUCART, B. HANIN, AND G. PETROVA, *Nonlinear approximation and (deep) relu networks*, 2019.
- [18] R. DEVORE AND A. RON, *Approximation using scattered shifts of a multivariate function*, Transactions of the American Mathematical Society, 362 (2010), pp. 6205–6229, <http://www.jstor.org/stable/40997201>.
- [19] R. A. DEVORE, *Nonlinear approximation*, Acta Numerica, 7 (1998), p. 51–150, <https://doi.org/10.1017/S0962492900002816>.
- [20] W. E, J. HAN, AND A. JENTZEN, *Deep learning-based numerical methods for high-dimensional parabolic partial differential equations and backward stochastic differential equations*, Communications in Mathematics and Statistics, 5 (2017), pp. 349–380, <https://doi.org/10.1007/s40304-017-0117-6>, <https://doi.org/10.1007/s40304-017-0117-6>.
- [21] W. E AND Q. WANG, *Exponential convergence of the deep neural network approximation for analytic functions*, CoRR, abs/1807.00297 (2018), <http://arxiv.org/abs/1807.00297>, <https://arxiv.org/abs/1807.00297>.
- [22] J. HAN, A. JENTZEN, AND W. E, *Solving high-dimensional partial differential equations using deep learning*, Proceedings of the National Academy of Sciences, 115 (2018), pp. 8505–8510, <https://doi.org/10.1073/pnas.1718942115>, <https://www.pnas.org/content/115/34/8505>, <https://arxiv.org/abs/https://www.pnas.org/content/115/34/8505.full.pdf>.
- [23] T. HANGELBROEK AND A. RON, *Nonlinear approximation using gaussian kernels*, Journal of Functional Analysis, 259 (2010), pp. 203 – 219, <https://doi.org/https://doi.org/10.1016/j.jfa.2010.02.001>, <http://www.sciencedirect.com/science/article/pii/S0022123610000467>.
- [24] B. HANIN AND M. SELLKE, *Approximating continuous functions by ReLU nets of minimal width*, (2017), <https://arxiv.org/abs/1710.11278>.
- [25] N. HARVEY, C. LIAW, AND A. MEHRABIAN, *Nearly-tight VC-dimension bounds for piecewise linear neural networks*, in Proceedings of the 2017 Conference on Learning Theory, S. Kale and O. Shamir, eds., vol. 65 of Proceedings of Machine Learning Research, Amsterdam, Netherlands, 07–10 Jul 2017, PMLR, pp. 1064–1068, <http://proceedings.mlr.press/v65/harvey17a.html>.
- [26] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, in 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2016, pp. 770–778, <https://doi.org/10.1109/CVPR.2016.90>.
- [27] K. HORNIK, *Approximation capabilities of multilayer feedforward networks*, Neural Networks, 4 (1991), pp. 251

- 257, [https://doi.org/https://doi.org/10.1016/0893-6080\(91\)90009-T](https://doi.org/https://doi.org/10.1016/0893-6080(91)90009-T), <http://www.sciencedirect.com/science/article/pii/089360809190009T>.
- [28] K. HORNIK, M. STINCHCOMBE, AND H. WHITE, *Multilayer feedforward networks are universal approximators*, Neural Networks, 2 (1989), pp. 359 – 366, [https://doi.org/https://doi.org/10.1016/0893-6080\(89\)90020-8](https://doi.org/https://doi.org/10.1016/0893-6080(89)90020-8), <http://www.sciencedirect.com/science/article/pii/0893608089900208>.
- [29] K. KAWAGUCHI, *Deep learning without poor local minima*, in Advances in Neural Information Processing Systems 29, D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, eds., Curran Associates, Inc., 2016, pp. 586–594, <http://papers.nips.cc/paper/6112-deep-learning-without-poor-local-minima.pdf>.
- [30] K. KAWAGUCHI AND Y. BENGIO, *Depth with nonlinearity creates no bad local minima in resnets*, (2018), <https://arxiv.org/abs/1810.09038>, <https://arxiv.org/abs/1810.09038>.
- [31] M. J. KEARNS AND R. E. SCHAPIRE, *Efficient distribution-free learning of probabilistic concepts*, J. Comput. Syst. Sci., 48 (1994), pp. 464–497, [https://doi.org/10.1016/S0022-0000\(05\)80062-5](https://doi.org/10.1016/S0022-0000(05)80062-5), [http://dx.doi.org/10.1016/S0022-0000\(05\)80062-5](http://dx.doi.org/10.1016/S0022-0000(05)80062-5).
- [32] A. KRIZHEVSKY, I. SUTSKEVER, AND G. E. HINTON, *Imagenet classification with deep convolutional neural networks*, in Advances in Neural Information Processing Systems 25, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, eds., Curran Associates, Inc., 2012, pp. 1097–1105, <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>.
- [33] V. KUMAR, *Introduction to Parallel Computing*, Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2nd ed., 2002.
- [34] V. KŮRKOVÁ, *Kolmogorov’s theorem and multilayer neural networks*, Neural Networks, 5 (1992), pp. 501 – 506, [https://doi.org/https://doi.org/10.1016/0893-6080\(92\)90012-8](https://doi.org/https://doi.org/10.1016/0893-6080(92)90012-8), <http://www.sciencedirect.com/science/article/pii/0893608092900128>.
- [35] G. LEWICKI AND G. MARINO, *Approximation of functions of finite variation by superpositions of a sigmoidal function*, Applied Mathematics Letters, 17 (2004), pp. 1147 – 1152, <http://www.sciencedirect.com/science/article/pii/S089396590481694X>.
- [36] S. LIANG AND R. SRIKANT, *Why deep neural networks?*, CoRR, abs/1610.04161 (2016), <http://arxiv.org/abs/1610.04161>, <https://arxiv.org/abs/1610.04161>.
- [37] S. LIN, X. LIU, Y. RONG, AND Z. XU, *Almost optimal estimates for approximation and learning by radial basis function networks*, Machine Learning, 95 (2014), pp. 147–164, <https://doi.org/10.1007/s10994-013-5406-z>, <https://doi.org/10.1007/s10994-013-5406-z>.
- [38] B. LLANAS AND F. SAINZ, *Constructive approximate interpolation by neural networks*, Journal of Computational and Applied Mathematics, 188 (2006), pp. 283 – 308, <http://www.sciencedirect.com/science/article/pii/S0377042705002566>.
- [39] Z. LU, H. PU, F. WANG, Z. HU, AND L. WANG, *The expressive power of neural networks: A view from the width*, in Advances in Neural Information Processing Systems 30, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, eds., Curran Associates, Inc., 2017, pp. 6231–6239, <http://papers.nips.cc/paper/7203-the-expressive-power-of-neural-networks-a-view-from-the-width.pdf>.
- [40] Z. LU, H. PU, F. WANG, Z. HU, AND L. WANG, *The expressive power of neural networks: A view from the width*, CoRR, abs/1709.02540 (2017), <http://arxiv.org/abs/1709.02540>, <https://arxiv.org/abs/1709.02540>.
- [41] V. MAIOROV AND A. PINKUS, *Lower bounds for approximation by mlp neural networks*, Neurocomputing, 25 (1999), pp. 81 – 91, [https://doi.org/https://doi.org/10.1016/S0925-2312\(98\)00111-8](https://doi.org/https://doi.org/10.1016/S0925-2312(98)00111-8), <http://www.sciencedirect.com/science/article/pii/S0925231298001118>.
- [42] S. G. MALLAT AND Z. ZHANG, *Matching pursuits with time-frequency dictionaries*, IEEE Transactions on Signal Processing, 41 (1993), pp. 3397–3415, <https://doi.org/10.1109/78.258082>.
- [43] E. J. MCSHANE, *Extension of range of functions*, Bull. Amer. Math. Soc., 40 (1934), pp. 837–842, <https://projecteuclid.org/443/euclid.bams/1183497871>.
- [44] H. MONTANELLI AND Q. DU, *New error bounds for deep networks using sparse grids*, (2017), <https://arxiv.org/abs/1712.08688>.
- [45] H. MONTANELLI, H. YANG, AND Q. DU, *Deep relu networks overcome the curse of dimensionality for bandlimited functions*, (2019), <https://arxiv.org/abs/1903.00735v1>, <https://arxiv.org/abs/1903.00735>.
- [46] G. F. MONTUFAR, R. PASCANU, K. CHO, AND Y. BENGIO, *On the number of linear regions of deep neural networks*, in Advances in Neural Information Processing Systems 27, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, eds., Curran Associates, Inc., 2014, pp. 2924–2932, <http://papers.nips.cc/paper/5422-on-the-number-of-linear-regions-of-deep-neural-networks.pdf>.
- [47] Q. N. NGUYEN AND M. HEIN, *The loss surface of deep and wide neural networks*, CoRR, abs/1704.08045 (2017), <http://arxiv.org/abs/1704.08045>, <https://arxiv.org/abs/1704.08045>.
- [48] J. PARK AND I. W. SANDBERG, *Universal approximation using radial-basis-function networks*, Neural Computation, 3 (1991), pp. 246–257, <https://doi.org/10.1162/neco.1991.3.2.246>.
- [49] P. PETERSEN AND F. VOIGTLAENDER, *Optimal approximation of piecewise smooth functions using deep ReLU neural networks*, Neural Networks, 108 (2018), pp. 296 – 330, <https://doi.org/https://doi.org/10.1016/j.neunet.2018.08.019>, <http://www.sciencedirect.com/science/article/pii/S0893608018302454>.
- [50] P. PETRUSHEV, *Multivariate n-term rational and piecewise polynomial approximation*, Journal of Approximation Theory, 121 (2003), pp. 158 – 197, [https://doi.org/https://doi.org/10.1016/S0021-9045\(02\)00060-6](https://doi.org/https://doi.org/10.1016/S0021-9045(02)00060-6), <http://www.sciencedirect.com>.

- [com/science/article/pii/S0021904502000606](http://www.sciencedirect.com/science/article/pii/S0021904502000606).
- [51] D. ROLNICK AND M. TEGMARK, *The power of deeper networks for expressing natural functions*, CoRR, abs/1705.05502 (2017), <http://arxiv.org/abs/1705.05502>, <https://arxiv.org/abs/1705.05502>.
 - [52] I. SAFRAN AND O. SHAMIR, *Depth-width tradeoffs in approximating natural functions with neural networks*, in Proceedings of the 34th International Conference on Machine Learning, D. Precup and Y. W. Teh, eds., vol. 70 of Proceedings of Machine Learning Research, International Convention Centre, Sydney, Australia, 06–11 Aug 2017, PMLR, pp. 2979–2987, <http://proceedings.mlr.press/v70/safran17a.html>.
 - [53] A. SAKURAI, *Tight bounds for the VC-dimension of piecewise polynomial networks*, in Advances in Neural Information Processing Systems, Neural information processing systems foundation, 1999, pp. 323–329.
 - [54] D. SCHERER, A. MÜLLER, AND S. BEHNKE, *Evaluation of pooling operations in convolutional architectures for object recognition*, in Artificial Neural Networks – ICANN 2010, K. Diamantaras, W. Duch, and L. S. Iliadis, eds., Berlin, Heidelberg, 2010, Springer Berlin Heidelberg, pp. 92–101.
 - [55] J. SCHMIDT-HIEBER, *Nonparametric regression using deep neural networks with ReLU activation function*, (2017), <https://arxiv.org/abs/1708.06633>.
 - [56] U. SHAHAM, A. CLONINGER, AND R. R. COIFMAN, *Provable approximation properties for deep neural networks*, Applied and Computational Harmonic Analysis, 44 (2018), pp. 537 – 557, <https://doi.org/https://doi.org/10.1016/j.acha.2016.04.003>, <http://www.sciencedirect.com/science/article/pii/S1063520316300033>.
 - [57] Z. SHEN, H. YANG, AND S. ZHANG, *Nonlinear Approximation via Compositions*, arXiv e-prints, (2019), arXiv:1902.10170, p. arXiv:1902.10170, <https://arxiv.org/abs/1902.10170>.
 - [58] T. SUZUKI, *Adaptivity of deep reLU network for learning in besov and mixed smooth besov spaces: optimal rate and curse of dimensionality*, in International Conference on Learning Representations, 2019, <https://openreview.net/forum?id=H1ebTsActm>.
 - [59] H. WHITNEY, *Analytic extensions of differentiable functions defined in closed sets*, Transactions of the American Mathematical Society, 36 (1934), pp. 63–89, <http://www.jstor.org/stable/1989708>.
 - [60] T. F. XIE AND F. L. CAO, *The rate of approximation of gaussian radial basis neural networks in continuous function space*, Acta Mathematica Sinica, English Series, 29 (2013), pp. 295–302, <https://doi.org/10.1007/s10114-012-1369-4>, <https://doi.org/10.1007/s10114-012-1369-4>.
 - [61] D. YAROTSKY, *Error bounds for approximations with deep ReLU networks*, Neural Networks, 94 (2017), pp. 103 – 114, <https://doi.org/https://doi.org/10.1016/j.neunet.2017.07.002>, <http://www.sciencedirect.com/science/article/pii/S0893608017301545>.
 - [62] D. YAROTSKY, *Optimal approximation of continuous functions by very deep ReLU networks*, in Proceedings of the 31st Conference On Learning Theory, S. Bubeck, V. Perchet, and P. Rigollet, eds., vol. 75 of Proceedings of Machine Learning Research, PMLR, 06–09 Jul 2018, pp. 639–649, <http://proceedings.mlr.press/v75/yarotsky18a.html>.