

Optimal Approximation Error in Width and Depth for ReLU Networks*

Zuowei Shen[†] Haizhao Yang[‡] Shijun Zhang[§]

Abstract

This paper concentrates on the approximation power of deep feed-forward neural networks in terms of width and depth. It is proved by construction that ReLU networks with width $\mathcal{O}(\max\{d\lfloor N^{1/d} \rfloor, N+2\})$ and depth $\mathcal{O}(L)$ can approximate a Hölder continuous function on $[0, 1]^d$ with a approximation rate $\mathcal{O}(\lambda\sqrt{d}(N^2L^2\log_3(N+2))^{-\alpha/d})$, where $\alpha \in (0, 1]$ and $\lambda > 0$ are Hölder order and constant, respectively. Such a rate is optimal up to a constant in terms of width and depth separately, while existing results are only nearly optimal without the log factor in the approximation rate. More generally, for an arbitrary continuous function f on $[0, 1]^d$, the approximation rate becomes $\mathcal{O}(\sqrt{d}\omega_f((N^2L^2\log_3(N+2))^{-1/d}))$, where $\omega_f(\cdot)$ is the modulus of continuity. We also extend our analysis to any continuous function f on a bounded set. Particularly, if ReLU networks with depth 31 and width $\mathcal{O}(N)$ is used to approximate one-dimensional Lipschitz continuous functions on $[0, 1]$, the approximation rate in terms of the total number of parameters, W , becomes $\mathcal{O}(\frac{1}{W\ln W})$, which has not been discovered in the literature.

Key words. Deep ReLU Networks; Hölder Continuity; Optimal Approximation Theory; Bit Extraction; VC-dimension.

1 Introduction

Over the past few decades, the expressiveness of neural networks has been widely studied from many points of view, e.g. in terms of combinatorics [17], topology [4], Vapnik-Chervonenkis (VC) dimension [3, 7, 20], fat-shattering dimension [1, 11], information theory [19], classical approximation theory [2, 5, 8, 12, 14, 21, 21, 22, 23, 24, 26, 27], optimization [9, 10, 18], etc. The error analysis of neural networks consists of three parts: the approximation error, the optimization error, and the generalization error. This paper focuses on the approximation error for ReLU networks.

*Submitted to the editors DATE.

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

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

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

The approximation errors of feed-forward neural networks with various activation functions have been studied for different types of functions, e.g., smooth functions [6, 13, 14, 15, 25], piecewise smooth functions [19], band-limited functions [16], continuous functions [22, 23, 24, 26]. In [22], it was shown that a ReLU network with width $C_1(d) \cdot N$ and depth $C_2(d) \cdot L$ can attain an approximation error $C_3(d) \cdot \omega_f(N^{-2/d}L^{-2/d})$ to approximate a continuous function f on $[0, 1]^d$, where $C_1(d)$, $C_2(d)$, and $C_3(d)$ are three constants in d with explicit formulas to specify their values, and $\omega_f(\cdot)$ is the modulus of continuity of $f \in C([0, 1]^d)$ defined via

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

Such an approximation rate is optimal in terms of N and L up to a logarithmic term and the corresponding optimal approximation theory is still open. To address this open problem, we provide a constructive proof in this paper to show that ReLU networks of width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ can approximate an arbitrary continuous function f on $[0, 1]^d$ with an optimal approximation error $\mathcal{O}(\omega_f(N^2 L^2 \ln N)^{-\alpha/d})$ in terms of N and L . As shown by our main result, Theorem 1.1 below, the approximation rate obtained here admits explicit formulas to specify its prefactors when $\omega_f(\cdot)$ is known.

Theorem 1.1. *Given a continuous function $f \in C([0, 1]^d)$, for any $N \in \mathbb{N}^+$, $L \in \mathbb{N}^+$, and $p \in [1, \infty]$, there exists a function ϕ implemented by a ReLU network with width $C_1 \max\{d\lfloor N^{1/d} \rfloor, N + 2\}$ and depth $11L + C_2$ such that*

$$\|f - \phi\|_{L^p([0, 1]^d)} \leq 131\sqrt{d}\omega_f\left((N^2 L^2 \log_3(N + 2))^{-1/d}\right),$$

where $C_1 = 16$ and $C_2 = 18$ if $p \in [1, \infty)$; $C_1 = 3^{d+3}$ and $C_2 = 18 + 2d$ if $p = \infty$.

Note that $3^{d+3} \max\{d\lfloor N^{1/d} \rfloor, N + 2\} \leq 3^{d+3} \max\{dN, 3N\} \leq 3^{d+4}dN$. Given any $\tilde{N}, \tilde{L} \in \mathbb{N}^+$ with $\tilde{N} \geq 3^{d+4}d$ and $\tilde{L} \geq 29 + 2d$, there exist $N, L \in \mathbb{N}^+$ such that

$$3^{d+4}dN \leq \tilde{N} < 3^{d+4}d(N + 1) \quad \text{and} \quad 11L + 18 + 2d \leq \tilde{L} < 11(L + 1) + 18 + 2d.$$

It follows that

$$N \geq \frac{N + 1}{3} > \frac{\tilde{N}}{3^{d+5}d} \quad \text{and} \quad L \geq \frac{L + 1}{2} > \frac{1}{2} \cdot \frac{\tilde{L} - 18 - 2d}{11} = \frac{\tilde{L} - 18 - 2d}{22}.$$

Then we have an immediate corollary from Theorem 1.1.

Corollary 1.2. *Given a continuous function $f \in C([0, 1]^d)$, for any $\tilde{N}, \tilde{L} \in \mathbb{N}^+$ with $\tilde{N} \geq 3^{d+4}d$ and $\tilde{L} \geq 29 + 2d$, there exists a function ϕ implemented by a ReLU network with width \tilde{N} and depth \tilde{L} such that*

$$\|f - \phi\|_{L^\infty([0, 1]^d)} \leq 131\sqrt{d}\omega_f\left(\left(\left(\frac{\tilde{N}}{3^{d+5}d}\right)^2 \left(\frac{\tilde{L} - 18 - 2d}{22}\right)^2 \log_3\left(\frac{\tilde{N}}{3^{d+5}d} + 2\right)\right)^{-1/d}\right).$$

As a special case of Theorem 1.1 for explicit error characterization, let us take Hölder continuous functions as an example. Let $\text{Hölder}([0, 1]^d, \alpha, \lambda)$ denote the space of Hölder continuous functions on $[0, 1]^d$ of order $\alpha \in (0, 1]$ with a Hölder constant $\lambda > 0$. We have an immediate corollary of Theorem 1.1 as follows.

64 **Corollary 1.3.** *Given a Hölder continuous function $f \in \text{Hölder}([0, 1]^d, \alpha, \lambda)$, for any*
65 *$N \in \mathbb{N}^+$, $L \in \mathbb{N}^+$, and $p \in [1, \infty]$, there exists a function ϕ implemented by a ReLU*
66 *network with width $C_1 \max\{d\lfloor N^{1/d} \rfloor, N + 2\}$ and depth $11L + C_2$ such that*

$$67 \quad \|f - \phi\|_{L^p([0, 1]^d)} \leq 131\lambda\sqrt{d}(N^2L^2\log_3(N + 2))^{-\alpha/d},$$

68 *where $C_1 = 16$ and $C_2 = 18$ if $p \in [1, \infty)$; $C_1 = 3^{d+3}$ and $C_2 = 18 + 2d$ if $p = \infty$.*

69 To better illustrate the importance of our theory, we summarize our key contribu-
70 tions as follows.

71 (1) Upper bound: We provide a quantitative and non-asymptotic approximation rate
72 $131\sqrt{d}\omega_f((N^2L^2\log_3(N + 2))^{-1/d})$ in terms of width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ for any
73 $f \in C([0, 1]^d)$ in Theorem 1.1.

74 (1.1) This approximation error analysis can be extended to $f \in C(E)$ for any $E \subseteq$
75 $[-R, R]^d$ with $R > 0$ as we shall see later in Theorem 2.5.

76 (1.2) In the case of one-dimensional Lipschitz continuous functions on $[0, 1]$, the
77 approximation rate in Theorem 1.1 becomes $\mathcal{O}(\frac{1}{W \ln W})$ for ReLU networks
78 with 31 hidden layers and $\mathcal{O}(W)$ parameters. To the best of our knowledge,
79 the approximation rate $\mathcal{O}(\frac{1}{W \ln W})$ is better than existing known results for
80 approximating Lipschitz continuous functions on $[0, 1]$.

81 (2) Lower bound: Through the VC-dimension bounds of ReLU networks given in [7], we
82 show that the approximation rate $131\lambda\sqrt{d}(N^2L^2\log_3(N + 2))^{-\alpha/d}$ in terms of width
83 $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ for Hölder($[0, 1]^d, \alpha, \lambda$) is optimal as follows.

84 (2.1) When the width is fixed, both the approximation upper and lower bounds take
85 the form of $CL^{-2\alpha/d}$ for a positive constant C .

86 (2.2) When the depth is fixed, both the approximation upper and lower bounds take
87 the form of $C(N^2 \ln N)^{-\alpha/d}$ for a positive constant C .

88 Compared to the nearly optimal rate $\mathcal{O}(N^{-2/d}L^{-2/d})$ for Hölder continuous functions
89 in our previous paper [22], this paper achieves the optimal rate $131\lambda\sqrt{d}(N^2L^2\log_3(N +$
90 $2))^{-\alpha/d}$ using innovative construction. For example, we propose a novel modification of
91 the bit extraction technique in [3] and construct new ReLU networks to approximate
92 step functions more efficiently than [22].

93 Remark that the full error analysis of deep learning consists of three parts: approxi-
94 mation, generalization, and optimization errors. Our constructive analysis here provides
95 an upper bound of the approximation error that characterizes the discrepancy between
96 the target function and the best approximator generated by networks. Instead of deriv-
97 ing an approximator in the form of a network to attain the approximation error bound
98 when the whole target function f is given, which is the goal of traditional approxima-
99 tion research, deep learning usually aims to identify a network from finite samples of

f such that this network can also approximate f well on unseen samples, i.e., reducing the optimization error on given samples and the generalization error on unseen samples simultaneously. The approximation error analysis in this paper is a key part of the generalization error analysis of deep learning.

The rest of this paper is organized as follows. In Section 2, we prove Theorem 1.1 by assuming Theorem 2.1 is true, show the optimality of Theorem 1.1, and extend our analysis to continuous functions defined on any bounded set. Next, Theorem 2.1 is proved in Section 3 based on Proposition 3.1 and 3.2, the proofs of which can be found in Section 4. Finally, Section 5 concludes this paper with a short discussion.

2 Theoretical analysis

In this section, we first prove Theorem 1.1 and discuss its optimality by assuming Theorem 2.1 is true. Next, we extend our analysis to general continuous functions defined on any bounded set in \mathbb{R}^d . Notations throughout this paper is summarized in Section 2.1.

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, $\mathbf{v} = [v_1, \dots, v_d]^T = \begin{bmatrix} v_1 \\ \vdots \\ v_d \end{bmatrix} \in \mathbb{R}^d$ is a column vector with $\mathbf{v}(i) = v_i$ being the i -th element. Besides, “[” and “]” are used to partition matrices (vectors) into blocks, e.g., $\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}$.

- For any $p \in [1, \infty)$, the p -norm (or ℓ^p -norm) of a vector $\mathbf{x} = [x_1, x_2, \dots, x_d]^T \in \mathbb{R}^d$ is defined by

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

- For any $x \in \mathbb{R}$, let $\lfloor x \rfloor := \max\{n : n \leq x, n \in \mathbb{Z}\}$ and $\lceil x \rceil := \min\{n : n \geq x, n \in \mathbb{Z}\}$.
- Assume $\mathbf{n} \in \mathbb{N}^d$, then $f(\mathbf{n}) = \mathcal{O}(g(\mathbf{n}))$ means that there exists positive C independent of \mathbf{n} , f , and g such that $f(\mathbf{n}) \leq Cg(\mathbf{n})$ when all entries of \mathbf{n} go to $+\infty$.
- For any $\theta \in [0, 1)$, suppose its binary representation is $\theta = \sum_{\ell=1}^{\infty} \theta_{\ell} 2^{-\ell}$ with $\theta_{\ell} \in \{0, 1\}$, we introduce a special notation $\text{bin}0.\theta_1\theta_2\cdots\theta_L$ to denote the L -term binary representation of θ , i.e., $\text{bin}0.\theta_1\theta_2\cdots\theta_L := \sum_{\ell=1}^L \theta_{\ell} 2^{-\ell}$.
- Let $\mu(\cdot)$ denote the Lebesgue measure.
- Let 1_S be the characteristic function on a set S , i.e., 1_S is equal to 1 on S and 0 outside S .
- Let $|S|$ denote the size of a set S , i.e., the number of all elements in S .

- The set difference of two sets A and B is denoted by $A \setminus B := \{x : x \in A, x \notin B\}$.
- Given any $K \in \mathbb{N}^+$ and $\delta \in (0, \frac{1}{K})$, define a trifling region $\Omega([0, 1]^d, K, \delta)$ of $[0, 1]^d$ as

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

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

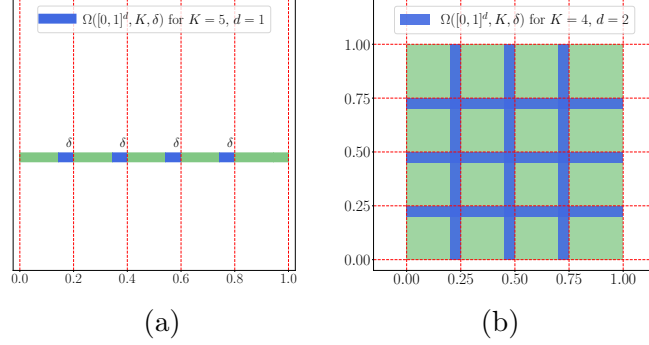


Figure 1: Two examples of trifling regions. (a) $K = 5, d = 1$. (b) $K = 4, d = 2$.

- Let $\text{Hölder}([0, 1]^d, \alpha, \lambda)$ denote the space of Hölder continuous functions on $[0, 1]^d$ of order $\alpha \in (0, 1]$ with a Hölder constant $\lambda > 0$.
- For a continuous piecewise linear function $f(x)$, the x values where the slope changes are typically called **breakpoints**.
- Let $\text{CPwL}(\mathbb{R}, n)$ denote the space that consists of all continuous piecewise linear functions with at most n breakpoints on \mathbb{R} .
- Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ denote the rectified linear unit (ReLU), i.e. $\sigma(x) = \max\{0, x\}$. With a slight abuse of notation, we define $\sigma : \mathbb{R}^d \rightarrow \mathbb{R}^d$ as $\sigma(\mathbf{x}) = \begin{bmatrix} \max\{0, x_1\} \\ \vdots \\ \max\{0, x_d\} \end{bmatrix}$ for any $\mathbf{x} = [x_1, \dots, x_d]^T \in \mathbb{R}^d$.
- We will use \mathcal{NN} to denote a function implemented by a ReLU network for short and use Python-type notations to specify a class of functions implemented by ReLU networks with several conditions, e.g., $\mathcal{NN}(c_1; c_2; \dots; c_m)$ is a set of functions implemented by ReLU networks satisfying m conditions given by $\{c_i\}_{1 \leq i \leq m}$, each of which may specify the number of inputs ($\#input$), the number of outputs ($\#output$), the total number of neurons in all hidden layers ($\#neuron$), the number of hidden layers (depth), the total number of parameters ($\#parameter$), and the width in each hidden layer (widthvec), the maximum width of all hidden layers (width), etc. For example, if $\phi \in \mathcal{NN}(\#input = 2; \text{widthvec} = [100, 100]; \#output = 1)$, then ϕ is a functions satisfies

- ϕ maps from \mathbb{R}^2 to \mathbb{R} .
- ϕ can be implemented by a ReLU network with two hidden layers and the number of nodes in each hidden layer is 100.
- For any function $\phi \in \mathcal{NN}(\#input = d; \text{widthvec} = [N_1, N_2, \dots, N_L]; \#output = 1)$, if we set $N_0 = d$ and $N_{L+1} = 1$, then the architecture of the network implementing ϕ can be briefly described as follows:

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

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

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

and

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

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

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

which has been illustrated in Figure 2.

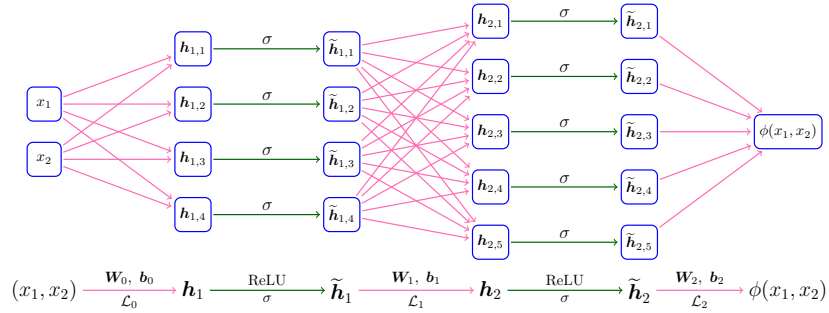


Figure 2: An example of a ReLU network with width 5 and depth 2.

- The expression “a network with width N and depth L ” means
 - The maximum width of this network for all **hidden** layers is no more than N .
 - The number of **hidden** layers of this network is no more than L .

2.2 Proof of Theorem 1.1

The key point is to construct piecewise constant functions to approximate continuous functions in the proof. However, it is impossible to construct a piecewise constant function implemented by a ReLU network due to the continuity of ReLU networks.

Thus, we introduce the trifling region $\Omega([0, 1]^d, K, \delta)$, defined in Equation (2.1), and use ReLU networks to implement piecewise constant functions outside the trifling region. To prove Theorem 1.1, we first introduce a weaker variant of Theorem 1.1, showing how to construct ReLU networks to pointwisely approximate continuous functions except for the trifling region.

Theorem 2.1. *Given a function $f \in C([0, 1]^d)$, for any $N \in \mathbb{N}^+$ and $L \in \mathbb{N}^+$, there exists a function ϕ implemented by a ReLU network with width $\max\{8d\lfloor N^{1/d} \rfloor + 3d, 16N + 30\}$ and depth $11L + 18$ such that $\|\phi\|_{L^\infty(\mathbb{R}^d)} \leq |f(\mathbf{0})| + \omega_f(\sqrt{d})$ and*

$$|f(\mathbf{x}) - \phi(\mathbf{x})| \leq 130\sqrt{d}\omega_f\left((N^2L^2\log_3(N+2))^{-1/d}\right), \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta),$$

where $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor^2 \lfloor \log_3(N+2) \rfloor^{1/d}$ and δ is an arbitrary number in $(0, \frac{1}{3K}]$.

With Theorem 2.1 that will be proved in Section 3, we can easily prove Theorem 1.1 for the case $p \in [1, \infty)$. To attain the rate in L^∞ -norm, we need to control the approximation error in the trifling region. To this end, we introduce a theorem to deal with the approximation inside the trifling region $\Omega([0, 1]^d, K, \delta)$.

Theorem 2.2 (Theorem 3.7 of [27] or Theorem 2.1 of [14]). *Given any $\varepsilon > 0$, $N, L, K \in \mathbb{N}^+$, and $\delta \in (0, \frac{1}{3K}]$, assume f is a continuous function in $C([0, 1]^d)$ and $\tilde{\phi}$ can be implemented by a ReLU network with width N and depth L . If*

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

then there exists a function ϕ implemented by a new ReLU network with width $3^d(N+4)$ and depth $L + 2d$ such that

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

Now we are ready to prove Theorem 1.1 by assuming Theorem 2.1 is true, which will be proved later in Section 3.

Proof of Theorem 1.1. We may assume f is not a constant function since it is a trivial case. Then $\omega_f(r) > 0$ for any $r > 0$. Let us first consider the case $p \in [1, \infty)$. Set $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor^2 \lfloor \log_3(N+2) \rfloor^{1/d}$ and choose a small $\delta \in (0, \frac{1}{3K}]$ such that

$$\begin{aligned} Kd\delta(2|f(\mathbf{0})| + 2\omega_f(\sqrt{d}))^p &= \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor^2 \lfloor \log_3(N+2) \rfloor^{1/d} d\delta(2|f(\mathbf{0})| + 2\omega_f(\sqrt{d}))^p \\ &\leq \left(\omega_f\left((N^2L^2\log_3(N+2))^{-1/d}\right) \right)^p. \end{aligned}$$

By Theorem 2.1, there exists a function ϕ implemented by a ReLU network with width

$$\max\{8d\lfloor N^{1/d} \rfloor + 3d, 16N + 30\} \leq 16 \max\{d\lfloor N^{1/d} \rfloor, N + 2\}$$

and depth $11L + 18$ such that $\|\phi\|_{L^\infty(\mathbb{R}^d)} \leq |f(\mathbf{0})| + \omega_f(\sqrt{d})$ and

$$|f(\mathbf{x}) - \phi(\mathbf{x})| \leq 130\sqrt{d}\omega_f\left((N^2L^2\log_3(N+2))^{-1/d}\right), \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta),$$

212 It follows from $\mu(\Omega([0, 1]^d, K, \delta)) \leq Kd\delta$ and $\|f\|_{L^\infty([0, 1]^d)} \leq |f(\mathbf{0})| + \omega_f(\sqrt{d})$ that

$$\begin{aligned}
\|f - \phi\|_{L^p([0, 1]^d)}^p &= \int_{\Omega([0, 1]^d, K, \delta)} |f(\mathbf{x}) - \phi(\mathbf{x})|^p d\mathbf{x} + \int_{[0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)} |f(\mathbf{x}) - \phi(\mathbf{x})|^p d\mathbf{x} \\
&\leq Kd\delta(2|f(\mathbf{0})| + 2\omega_f(\sqrt{d}))^p + \left(130\sqrt{d}\omega_f\left((N^2L^2\log_3(N+2))^{-1/d}\right)\right)^p \\
213 \quad &\leq \left(\omega_f\left((N^2L^2\log_3(N+2))^{-1/d}\right)\right)^p + \left(130\sqrt{d}\omega_f\left((N^2L^2\log_3(N+2))^{-1/d}\right)\right)^p \\
&\leq \left(131\sqrt{d}\omega_f\left((N^2L^2\log_3(N+2))^{-1/d}\right)\right)^p.
\end{aligned}$$

214 Hence, $\|f - \phi\|_{L^p([0, 1]^d)} \leq 131\sqrt{d}\omega_f\left((N^2L^2\log_3(N+2))^{-1/d}\right)$.

215 Next, let us discuss the case $p = \infty$. Set $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor^2 \lfloor \log_3(N+2) \rfloor^{1/d}$ and
216 choose a small $\delta \in (0, \frac{1}{3K}]$ such that

$$217 \quad d \cdot \omega_f(\delta) \leq \omega_f\left((N^2L^2\log_3(N+2))^{-1/d}\right).$$

218 By Theorem 2.1, there exists a function $\tilde{\phi}$ implemented by a ReLU network with width
219 $\max\{8d\lfloor N^{1/d} \rfloor + 3d, 16N + 30\}$ and depth $11L + 18$ such that

$$220 \quad |f(\mathbf{x}) - \tilde{\phi}(\mathbf{x})| \leq 130\sqrt{d}\omega_f\left((N^2L^2\log_3(N+2))^{-1/d}\right) =: \varepsilon,$$

221 for any $\mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)$. By Theorem 2.2, there exists a function ϕ imple-
222 mented by a ReLU network with width

$$223 \quad 3^d \left(\max\{8d\lfloor N^{1/d} \rfloor + 3d, 16N + 30\} + 4 \right) \leq 3^{d+3} \max\{d\lfloor N^{1/d} \rfloor, N + 2\}$$

224 and depth $11L + 18 + 2d$ such that

$$225 \quad |f(\mathbf{x}) - \phi(\mathbf{x})| \leq \varepsilon + d \cdot \omega_f(\delta) \leq 131\sqrt{d}\omega_f\left((N^2L^2\log_3(N+2))^{-1/d}\right), \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

226 So we finish the proof. \square

227 2.3 Optimality

228 This section will show that the approximation rates in Theorem 1.1 and Corollary 1.3
229 are optimal and there is no room to improve for the function class $\text{H\"older}([0, 1]^d, \alpha, \lambda)$.
230 Therefore, the approximation rate for the whole continuous functions space in terms of
231 width and depth in Theorem 1.1 cannot be improved. A typical method to characterize
232 the optimal approximation theory of neural networks is to study the connection between
233 the approximation error and VapnikChervonenkis (VC) dimension [14, 22, 25, 26, 27]. This
234 method relies on the VC-dimension upper bound given in [7]. In this paper, we adopt
235 this method with several modifications to simplify the proof.

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

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

where $|\cdot|$ denotes the size of a set. This equation means, given any $\theta_i \in \{0, 1\}$ for $i = 1, 2, \dots, m$, there exists $h \in H$ such that $h(\mathbf{x}_i) = \theta_i$.

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

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

Definition 2.3 (VC-dimension). Let H be a class of functions from \mathcal{X} to $\{0, 1\}$. The VC-dimension of H , denoted by $\text{VCDim}(H)$, is the size of the largest shattered set, namely, $\text{VCDim}(H) := \sup\{m \in \mathbb{N}^+ : \Pi_H(m) = 2^m\}$.

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

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

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

The theorem below, Theorem 2.4, reveals the connection between VC-dimension and approximation rate.

Theorem 2.4. Assume \mathcal{F} is a function set with all elements defined on $[0, 1]^d$. For any $\varepsilon \in (0, 2/9)$, if

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

then $\text{VCDim}(\mathcal{F}) \geq (9\varepsilon)^{-d/\alpha}$.

This theorem demonstrates the connection between VC-dimension of \mathcal{F} and the approximation rate using elements of \mathcal{F} to approximate functions in $\text{H\"older}([0, 1]^d, \alpha, \lambda)$. To be precise, the VC-dimension of \mathcal{F} determines an approximation rate lower bound $\text{VCDim}(\mathcal{F})^{-\alpha/d/9}$, which is the best possible approximation rate. Denote the best approximation error of functions in $\text{H\"older}([0, 1]^d, \alpha, 1)$ approximated by ReLU networks with width N and depth L as

$$\mathcal{E}_{\alpha, d}(N, L) := \sup_{f \in \text{H\"older}([0, 1]^d, \alpha, 1)} \left(\inf_{\phi \in \mathcal{NN}(\text{width} \leq N; \text{depth} \leq L)} \|\phi - f\|_{L^\infty([0, 1]^d)} \right),$$

We have two remarks listed below.

- (i) A large VC-dimension cannot guarantee a good approximation rate. For example, it is easy to verify that

$$\text{VCDim}\left(\{f : f(x) = \cos(ax), a \in \mathbb{R}\}\right) = \infty.$$

However, functions in $\{f : f(x) = \cos(ax), a \in \mathbb{R}\}$ cannot approximate Hölder continuous functions well.

- (ii) A large VC-dimension is necessary for a good approximation rate, because the best possible approximation rate is controlled by an expression of VC-dimension, as shown in Theorem 2.4. For example, Theorem 6 and 8 of [7] implies that

$$\text{VCDim}(\mathcal{NN}(\text{width} \leq N; \text{depth} \leq L)) \leq \min\left\{\mathcal{O}(N^2 L^2 \ln(NL)), \mathcal{O}(N^3 L^2)\right\},$$

deducing

$$\underbrace{C_1(\alpha, d) \left(\min\{N^2 L^2 \ln(NL), N^3 L^2\} \right)^{-\alpha/d}}_{\text{implied by Theorem 2.4}} \leq \mathcal{E}_{\alpha, d}(N, L) \leq \underbrace{C_2(\alpha, d) \left(N^2 L^2 \ln N \right)^{-\alpha/d}}_{\text{implied by Corollary 1.3}}, \quad \textcircled{1} \quad (2.3)$$

for any $N, L \in \mathbb{N}^+$ with $N \geq 2$, where $C_1(\alpha, d)$ and $C_2(\alpha, d)$ are two positive constants determined by s, d , and $C_2(s, d)$ can be explicitly expressed. When $L = L_0$ is fixed, Equation 2.3 implies

$$C_1(\alpha, d, L_0)(N^2 \ln N)^{-\alpha/d} \leq \mathcal{E}_{\alpha, d}(N, L_0) \leq C_2(\alpha, d, L_0)(N^2 \ln N)^{-\alpha/d},$$

where $C_1(\alpha, d, L_0)$ and $C_2(\alpha, d, L_0)$ are two position constant determined by α, d, L_0 . When $N = N_0$ is fixed, Equation 2.3 implies

$$C_1(\alpha, d, N_0)L^{-2\alpha/d} \leq \mathcal{E}_{\alpha, d}(N_0, L) \leq C_2(\alpha, d, N_0)L^{-2\alpha/d},$$

where $C_1(\alpha, d, N_0)$ and $C_2(\alpha, d, N_0)$ are two position constant determined by α, d, N_0 .

Finally, let us present the detailed proof of Theorem 2.4.

Proof of Theorem 2.4. Recall that the VC-dimension of a function set is defined as the size of the largest set of points that this class of functions can shatter. So our goal is to find a subset of \mathcal{F} to shatter $\mathcal{O}(\varepsilon^{-d/\alpha})$ points in $[0, 1]^d$, which can be divided into two steps.

- Construct $\{f_\chi : \chi \in \mathcal{B}\} \subseteq \text{Hölder}([0, 1]^d, \alpha, 1)$ that scatters $\mathcal{O}(\varepsilon^{-d/\alpha})$ points, where \mathcal{B} is a set defined later.
- Design $\phi_\chi \in \mathcal{F}$, for each $\chi \in \mathcal{B}$, based on f_χ and Equation (2.2) such that $\{\phi_\chi : \chi \in \mathcal{B}\} \subseteq \mathcal{F}$ also shatters $\mathcal{O}(\varepsilon^{-d/\alpha})$ points.

^①To make this equation valid for any $N, L \in \mathbb{N}^+$ with $N \geq 2$, one needs to choose $C_1(\alpha, d)$ and $C_2(\alpha, d)$ carefully based on Theorem 2.4 and Corollary 1.3.

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

296 **Step 1:** Construct $\{f_\chi : \chi \in \mathcal{B}\} \subseteq \text{Hölder}([0, 1]^d, \alpha, 1)$ that scatters $\mathcal{O}(\varepsilon^{-d/\alpha})$ points.

297 Let $K = \lfloor (9\varepsilon/2)^{-1/\alpha} \rfloor \in \mathbb{N}^+$ and divide $[0, 1]^d$ into K^d non-overlapping sub-cubes
298 $\{Q_\beta\}_\beta$ as follows:

$$299 \quad Q_\beta := \{\mathbf{x} = [x_1, x_2, \dots, x_d]^T \in [0, 1]^d : x_i \in [\frac{\beta_i}{K}, \frac{\beta_i+1}{K}], \ i = 1, 2, \dots, d\},$$

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

301 Define a function ζ_Q on $[0, 1]^d$ corresponding to $Q = Q(\mathbf{x}_0, \eta) \subseteq [0, 1]^d$ such that:

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

305 Define

$$306 \quad \mathcal{B} := \{\chi : \chi \text{ is a map from } \{0, 1, \dots, K-1\}^d \text{ to } \{-1, 1\}\}.$$

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

$$308 \quad f_\chi(\mathbf{x}) := \sum_{\beta \in \{0, 1, \dots, K-1\}^d} \chi(\beta) \zeta_{Q_\beta}(\mathbf{x}),$$

309 where $\zeta_{Q_\beta}(\mathbf{x})$ is the associated function introduced just above. It is easy to check that
310 $\{f_\chi : \chi \in \mathcal{B}\} \subseteq \text{Hölder}([0, 1]^d, \alpha, 1)$ can shatter $K^d = \mathcal{O}(\varepsilon^{-d/\alpha})$ points in $[0, 1]^d$.

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

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

$$313 \quad \|\phi_\chi - f_\chi\|_{L^\infty([0, 1]^d)} \leq \varepsilon + \varepsilon/81.$$

314 Let $\mu(\cdot)$ denote the Lebesgue measure of a set. Then, for each $\chi \in \mathcal{B}$, there exists
315 $\mathcal{H}_\chi \subseteq [0, 1]^d$ with $\mu(\mathcal{H}_\chi) = 0$ such that

$$316 \quad |\phi_\chi(\mathbf{x}) - f_\chi(\mathbf{x})| \leq \frac{82}{81}\varepsilon, \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \mathcal{H}_\chi.$$

317 Set $\mathcal{H} = \cup_{\chi \in \mathcal{B}} \mathcal{H}_\chi$, then we have $\mu(\mathcal{H}) = 0$ and

$$318 \quad |\phi_\chi(\mathbf{x}) - f_\chi(\mathbf{x})| \leq \frac{82}{81}\varepsilon, \quad \text{for any } \chi \in \mathcal{B} \text{ and } \mathbf{x} \in [0, 1]^d \setminus \mathcal{H}. \quad (2.4)$$

319 Since Q_β has a sidelength $\frac{1}{K} = \frac{1}{\lfloor (9\varepsilon/2)^{-1/\alpha} \rfloor}$, we have, for each $\beta \in \{0, 1, \dots, K-1\}^d$ and
320 any $\mathbf{x} \in \frac{1}{10}Q_\beta$ ^②,

$$321 \quad |f_\chi(\mathbf{x})| = |\zeta_{Q_\beta}(\mathbf{x})| \geq \frac{9}{10}|\zeta_{Q_\beta}(\mathbf{x}_{Q_\beta})| = \frac{9}{10}\left(\frac{1}{2\lfloor (9\varepsilon/2)^{-1/\alpha} \rfloor}\right)^\alpha/2 \geq \frac{81}{80}\varepsilon, \quad (2.5)$$

^② $\frac{1}{10}Q_\beta$ denotes the closed cube whose sidelength is 1/10 of that of Q_β and which shares the same center of Q_β .

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

Note that $(\frac{1}{10}Q_\beta) \setminus \mathcal{H}$ is not empty, since $\mu((\frac{1}{10}Q_\beta) \setminus \mathcal{H}) > 0$ for each $\beta \in \{0, 1, \dots, K-1\}^d$. Together with Equation (2.4) and (2.5), there exists $\mathbf{x}_\beta \in (\frac{1}{10}Q_\beta) \setminus \mathcal{H}$ such that, for each $\beta \in \{0, 1, \dots, K-1\}^d$ and each $\chi \in \mathcal{B}$,

$$|f_\chi(\mathbf{x}_\beta)| \geq \frac{81}{80}\varepsilon > \frac{82}{81}\varepsilon \geq |f_\chi(\mathbf{x}_\beta) - \phi_\chi(\mathbf{x}_\beta)|,$$

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

$$\text{VCDim}(\mathcal{F}) \geq \text{VCDim}(\{\phi_\chi : \chi \in \mathcal{B}\}) \geq K^d = \lfloor (9\varepsilon/2)^{-1/\alpha} \rfloor^d \geq (9\varepsilon)^{-d/\alpha}, \quad (2.6)$$

where the last inequality comes from the fact $\lfloor x \rfloor \geq x/2 \geq x/(2^{1/\alpha})$ for any $x \in [1, \infty)$ and $\alpha \in (0, 1]$. So we finish the proof. \square

2.4 Approximation in irregular domain

We extend our analysis to general continuous functions defined on any irregular bounded set in \mathbb{R}^d . The key idea is to extend the target function to a hypercube while preserving the modulus of continuity. For a general set $E \subseteq \mathbb{R}^d$, the modulus of continuity of $f \in C(E)$ is defined via

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

In particular, $\omega_f(\cdot)$ is short of $\omega_f^E(\cdot)$ in the case of $E = [0, 1]^d$. Then, Theorem 1.1 can be generalized to $f \in C(E)$ for any bounded set $E \subseteq [-R, R]^d$ with $R > 0$, as shown in the following theorem.

Theorem 2.5. *Given a continuous function $f \in C(E)$ with $E \subseteq [-R, R]^d$ and $R > 0$, for any $N \in \mathbb{N}^+$, $L \in \mathbb{N}^+$, and $p \in [1, \infty]$, there exists a function ϕ implemented by a ReLU network with width $C_1 \max \{d \lfloor N^{1/d} \rfloor, N+2\}$ and depth $11L + C_2$ such that*

$$\|f - \phi\|_{L^p(E)} \leq 131(2R)^{d/p} \sqrt{d} \omega_f^E \left(2R(N^2 L^2 \log_3(N+2))^{-1/d} \right),$$

where $C_1 = 16$ and $C_2 = 18$ if $p \in [1, \infty)$; $C_1 = 3^{d+3}$ and $C_2 = 18 + 2d$ if $p = \infty$.

Proof. Given any $f \in C(E)$, by Lemma 4.2 of [22] via setting $S = \mathbb{R}^d$, there exists $g \in C(\mathbb{R}^d)$ such that

- $g(\mathbf{x}) = f(\mathbf{x})$ for any $\mathbf{x} \in E \subseteq [-R, R]^d$;
- $\omega_g^S(r) = \omega_f^E(r)$ for any $r \geq 0$.

Define

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

By applying Theorem 1.1 to $\tilde{g} \in C([0, 1]^d)$, there exists a function $\tilde{\phi}$ implemented by a ReLU network with width $C_1 \max\{d\lfloor N^{1/d} \rfloor, N + 2\}$ and depth $11L + C_2$ such that

$$\|\tilde{\phi} - \tilde{g}\|_{L^p([0, 1]^d)} \leq 131\sqrt{d}\omega_{\tilde{g}}\left((N^2L^2\log_3(N + 2))^{-1/d}\right),$$

where $C_1 = 16$ and $C_2 = 18$ if $p \in [1, \infty)$; $C_1 = 3^{d+3}$ and $C_2 = 18 + 2d$ if $p = \infty$.

Recall that $f(\mathbf{x}) = g(\mathbf{x}) = \tilde{g}(\frac{\mathbf{x}+R}{2R})$ for any $\mathbf{x} \in E \subseteq [-R, R]^d$ and

$$\omega_{\tilde{g}}(r) \leq \omega_{\tilde{g}}^S(r) = \omega_{\tilde{g}}^S(2Rr) = \omega_f^E(2Rr), \quad \text{for any } r \geq 0.$$

Define $\phi(\mathbf{x}) := \tilde{\phi}(\frac{\mathbf{x}+R}{2R}) = \tilde{\phi} \circ \mathcal{L}(\mathbf{x})$ for any $\mathbf{x} \in \mathbb{R}^d$, where \mathcal{L} is an affine linear map given by $\mathcal{L}(\mathbf{x}) = \frac{\mathbf{x}+R}{2R}$. Clearly, ϕ can be implemented by a ReLU network with width $C_1 \max\{d\lfloor N^{1/d} \rfloor, N + 2\}$ and depth $11L + C_2$, where $C_1 = 16$ and $C_2 = 18$ if $p \in [1, \infty)$; $C_1 = 3^{d+3}$ and $C_2 = 18 + 2d$ if $p = \infty$. Moreover, for any $\mathbf{x} \in E \subseteq [-R, R]^d$, we have $\frac{\mathbf{x}+R}{2R} \in [0, 1]^d$, implying

$$\begin{aligned} \|\phi - f\|_{L^p(E)} &= \|\phi - g\|_{L^p(E)} = \|\tilde{\phi} \circ \mathcal{L} - \tilde{g} \circ \mathcal{L}\|_{L^p(E)} \\ &\leq \|\tilde{\phi} \circ \mathcal{L} - \tilde{g} \circ \mathcal{L}\|_{L^p([-R, R]^d)} = (2R)^{d/p} \|\tilde{\phi} - \tilde{g}\|_{L^p([0, 1]^d)} \\ &\leq 131(2R)^{d/p} \sqrt{d}\omega_{\tilde{g}}\left((N^2L^2\log_3(N + 2))^{-1/d}\right) \\ &\leq 131(2R)^{d/p} \sqrt{d}\omega_f^E\left(2R(N^2L^2\log_3(N + 2))^{-1/d}\right). \end{aligned}$$

With the discussion above, we have proved Theorem 2.5. \square

3 Proof of Theorem 2.1

We will prove Theorem 2.1 in this section. We first present the key ideas in Section 3.1. The detailed proof is presented in Section 3, based on two propositions in Section 3.1, the proofs of which can be founded in Section 4.

3.1 Key ideas of proving Theorem 2.1

Given an arbitrary $f \in C([0, 1]^d)$, our goal is to construct an almost piecewise constant function ϕ implemented by a ReLU network to approximate f well. To this end, we introduce a piecewise constant function $f_p \approx f$ serving as an intermediate approximant in our construction in the sense that

$$f \approx f_p \text{ on } [0, 1]^d \quad \text{and} \quad f_p \approx \phi \text{ on } [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta).$$

The approximation in $f \approx f_p$ is a simple and standard technique in constructive approximation. The most technical part is to design a deep ReLU network with the desired width and depth to implement a function ϕ with $\phi \approx f_p$ outside $\Omega([0, 1]^d, K, \delta)$. See Figure 3 for an illustration. The introduction of the trifling region is to ease the construction of ϕ , which is a continuous piecewise linear 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 $[0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)$.

Now let us discuss the detailed steps of construction.

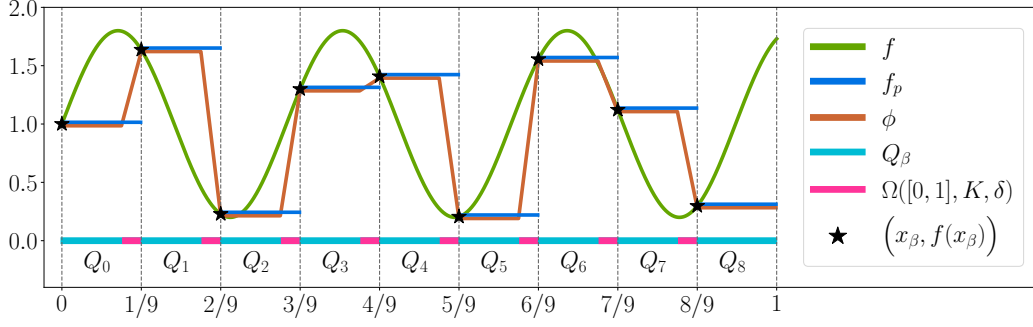


Figure 3: An illustration of f , f_p , ϕ , x_β , Q_β , and the trifling region $\Omega([0, 1]^d, K, \delta)$ in the one-dimensional case for $\beta \in \{0, 1, \dots, K-1\}^d$, where $K = N^2 L^2 \log_3(N+2)$ and $d = 1$ with $N = 1$ and $L = 3$. f is the target function; f_p is the piecewise constant function approximating f ; ϕ is a function, implemented by a ReLU network, approximating f ; and x_β is a representative of Q_β . The measure of $\Omega([0, 1]^d, K, \delta)$ can be arbitrarily small as we shall see in the proof of Theorem 1.1.

(1) First, divide $[0, 1]^d$ into a union of important regions $\{Q_\beta\}_\beta$ and the trifling region $\Omega([0, 1]^d, K, \delta)$, where each Q_β is associated with a representative $\mathbf{x}_\beta \in Q_\beta$ such that $f(\mathbf{x}_\beta) = f_p(\mathbf{x}_\beta)$ for each index vector $\beta \in \{0, 1, \dots, K-1\}^d$, where $K = \mathcal{O}((N^2 L^2 \ln N)^{1/d})$ is the partition number per dimension (see Figure 5 for examples for $d = 1$ and $d = 2$).

(2) Next, we design a vector function $\Phi_1(\mathbf{x})$ constructed via

$$\Phi_1(\mathbf{x}) = [\phi_1(x_1), \phi_1(x_2), \dots, \phi_1(x_d)]^T$$

to project the whole cube Q_β to a d -dimensional index β for each β , where each one-dimensional function ϕ_1 is a step function implemented by a ReLU network.

(3) The third step is to solve a point fitting problem. To be precise, we construct a function ϕ_2 implemented by a ReLU network to map β approximately to $f_p(\mathbf{x}_\beta) = f(\mathbf{x}_\beta)$. Then $\phi_2 \circ \Phi_1(\mathbf{x}) = \phi_2(\beta) \approx f_p(\mathbf{x}_\beta) = f(\mathbf{x}_\beta)$ for any $\mathbf{x} \in Q_\beta$ and each β , implying $\phi := \phi_2 \circ \Phi_1 \approx f_p \approx f$ on $[0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)$. We would like to point out that we only need to care about the values of ϕ_2 at a set of points $\{0, 1, \dots, K-1\}^d$ in the construction of ϕ_2 according to our design $\phi = \phi_2 \circ \Phi_1$ as illustrated in Figure 4. Therefore, it is not necessary to care about the values of ϕ_2 sampled outside the set $\{0, 1, \dots, K-1\}^d$, which is a key point to ease the design of a ReLU network to implement ϕ_2 as we shall see later.

Finally, we discuss how to implement Φ_1 and ϕ_2 by deep ReLU networks with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ using two propositions as we shall prove in Section 4.2 and 4.3 later. We first construct a ReLU network with desired width and depth by Proposition 3.1 to implement a one-dimensional step function ϕ_1 . Then Φ_1 can be attained via defining

$$\Phi_1(\mathbf{x}) = [\phi_1(x_1), \phi_1(x_2), \dots, \phi_1(x_d)]^T, \quad \text{for any } \mathbf{x} = [x_1, x_2, \dots, x_d]^T \in \mathbb{R}^d.$$

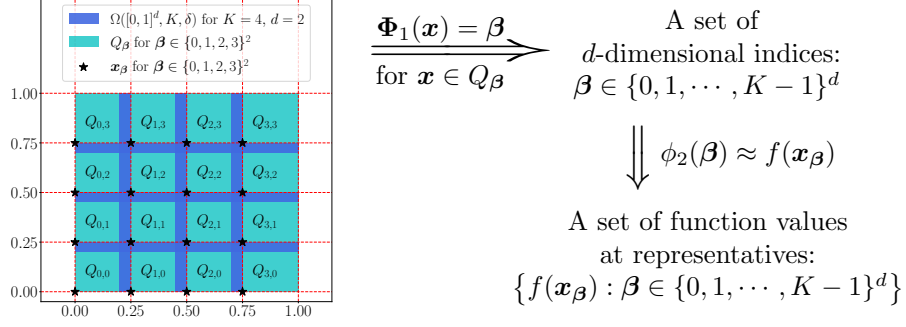


Figure 4: An illustration of the desired function $\phi = \phi_2 \circ \Phi_1$. Note that $\phi \approx f$ on $[0,1]^d \setminus \Omega([0,1]^d, K, \delta)$, since $\phi(\mathbf{x}) = \phi_2 \circ \Phi_1(\mathbf{x}) = \phi_2(\beta) \approx f(\mathbf{x}_\beta)$ for any $\mathbf{x} \in Q_\beta$ and each $\beta \in \{0, 1, \dots, K-1\}^d$.

408 **Proposition 3.1.** *For any $N, L, d \in \mathbb{N}^+$ and $\delta \in (0, \frac{1}{3K}]$ with*

409
$$K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor \lfloor n^{1/d} \rfloor, \quad \text{where } n = \lfloor \log_3(N+2) \rfloor,$$

410 *there exists a one-dimensional function ϕ implemented by a ReLU network with width*
 411 *$8\lfloor N^{1/d} \rfloor + 3$ and depth $2\lfloor L^{1/d} \rfloor + 5$ such that*

412
$$\phi(x) = k, \quad \text{if } x \in \left[\frac{k}{K}, \frac{k+1}{K} - \delta \cdot 1_{\{k \leq K-2\}} \right] \text{ for } k = 0, 1, \dots, K-1.$$

413 The construction of ϕ_2 is a direct result of Proposition 3.2 below, the proof of which
 414 relies on the bit extraction technique in [3].

415 **Proposition 3.2.** *Given any $\varepsilon > 0$ and arbitrary $N, L, J \in \mathbb{N}^+$ with $J \leq N^2 L^2 \lfloor \log_3(N+2) \rfloor$,*
 416 *assume $y_j \geq 0$ for $j = 0, 1, \dots, J-1$ are samples with*

417
$$|y_j - y_{j-1}| \leq \varepsilon, \quad \text{for } j = 1, 2, \dots, J-1.$$

418 *Then there exists $\phi \in \mathcal{NN}$ (#input = 1; width $\leq 16N + 30$; depth $\leq 6L + 10$; #output = 1)*
 419 *such that*

420 (i) $|\phi(j) - y_j| \leq \varepsilon$ for $j = 0, 1, \dots, J-1$.

421 (ii) $0 \leq \phi(x) \leq \max\{y_j : j = 0, 1, \dots, J-1\}$ for any $x \in \mathbb{R}$.

422 With the above propositions ready, let us prove Theorem 2.1 in Section 3.

423 3.2 Constructive proof

424 We essentially construct an almost piecewise constant function implemented by a
 425 ReLU network with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ to approximate f . We may assume f
 426 is not a constant function since it is a trivial case. Then $\omega_f(r) > 0$ for any $r > 0$. It is
 427 clear that $|f(\mathbf{x}) - f(\mathbf{0})| \leq \omega_f(\sqrt{d})$ for any $\mathbf{x} \in [0,1]^d$. Define $\tilde{f} = f - f(\mathbf{0}) + \omega_f(\sqrt{d})$, then
 428 $0 \leq \tilde{f}(\mathbf{x}) \leq 2\omega_f(\sqrt{d})$ for any $\mathbf{x} \in [0,1]^d$.

429 Let $M = N^2 L$, $n = \lfloor \log_3(N+2) \rfloor$, $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor^2 \lfloor n^{1/d} \rfloor$, and δ be an arbitrary
 430 number in $(0, \frac{1}{3K}]$. The proof can be divided into four steps as follows:

1. Normalize f as \tilde{f} , divide $[0, 1]^d$ into a union of sub-cubes $\{Q_\beta\}_{\beta \in \{0,1,\dots,K-1\}^d}$ and the trifling region $\Omega([0, 1]^d, K, \delta)$, and denote \mathbf{x}_β as the vertex of Q_β with minimum $\|\cdot\|_1$ norm;
2. Construct a sub-network to implement a vector function Φ_1 projecting the whole cube Q_β to the d -dimensional index β for each β , i.e., $\Phi_1(\mathbf{x}) = \beta$ for all $\mathbf{x} \in Q_\beta$;
3. Construct a sub-network to implement a function ϕ_2 mapping the index β approximately to $\tilde{f}(\mathbf{x}_\beta)$. This core step can be further divided into three sub-steps:
 - 3.1. Construct a sub-network to implement ψ_1 bijectively mapping the index set $\{0, 1, \dots, K-1\}^d$ to an auxiliary set $\mathcal{A}_1 \subseteq \{\frac{j}{2K^d} : j = 0, 1, \dots, 2K^d\}$ defined later (see Figure 6 for an illustration);
 - 3.2. Determine a continuous piecewise linear function g with a set of breakpoints $\mathcal{A}_1 \cup \mathcal{A}_2 \cup \{1\}$ satisfying: 1) assign the values of g at breakpoints in \mathcal{A}_1 based on $\{\tilde{f}(\mathbf{x}_\beta)\}_\beta$, i.e., $g \circ \psi_1(\beta) = \tilde{f}(\mathbf{x}_\beta)$; 2) assign the values of g at breakpoints in $\mathcal{A}_2 \cup \{1\}$ to reduce the variation of g for applying Proposition 3.2;
 - 3.3. Apply Proposition 3.2 to construct a sub-network to implement a function ψ_2 approximating g well on $\mathcal{A}_1 \cup \mathcal{A}_2 \cup \{1\}$. Then the desired function ϕ_2 is given by $\phi_2 = \psi_2 \circ \psi_1$ satisfying $\phi_2(\beta) = \psi_2 \circ \psi_1(\beta) \approx g \circ \psi_1(\beta) = \tilde{f}(\mathbf{x}_\beta)$;
4. Construct the final target network to implement the desired function ϕ such that $\phi(\mathbf{x}) = \phi_2 \circ \Phi_1(\mathbf{x}) + f(\mathbf{0}) - \omega_f(\sqrt{d}) \approx \tilde{f}(\mathbf{x}_\beta) + f(\mathbf{0}) - \omega_f(\sqrt{d}) = f(\mathbf{x}_\beta)$ for $\mathbf{x} \in Q_\beta$.

The details of these steps can be found below.

Step 1: Divide $[0, 1]^d$ into $\{Q_\beta\}_{\beta \in \{0,1,\dots,K-1\}^d}$ and $\Omega([0, 1]^d, K, \delta)$.

Define $\mathbf{x}_\beta := \beta/K$ and

$$Q_\beta := \left\{ \mathbf{x} = [x_1, \dots, x_d]^T \in [0, 1]^d : x_i \in \left[\frac{\beta_i}{K}, \frac{\beta_i+1}{K} - \delta \cdot 1_{\{\beta_i \leq K-2\}} \right], i = 1, \dots, d \right\}$$

for each d -dimensional index $\beta = [\beta_1, \dots, \beta_d]^T \in \{0, 1, \dots, K-1\}^d$. Recall that $\Omega([0, 1]^d, K, \delta)$ is the trifling region defined in Equation (2.1). Apparently, \mathbf{x}_β is the vertex of Q_β with minimum $\|\cdot\|_1$ norm and

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

see Figure 5 for illustrations.

Step 2: Construct Φ_1 mapping $\mathbf{x} \in Q_\beta$ to β .

By Proposition 3.1, there exists $\phi_1 \in \mathcal{NN}(\text{width} \leq 4\lfloor N^{1/d} \rfloor + 3; \text{depth} \leq 4\lfloor L^{1/d} \rfloor + 5)$ such that

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

It follows that $\phi_1(x_i) = \beta_i$ if $\mathbf{x} = [x_1, x_2, \dots, x_d]^T \in Q_\beta$ for each $\beta = [\beta_1, \beta_2, \dots, \beta_d]^T$.

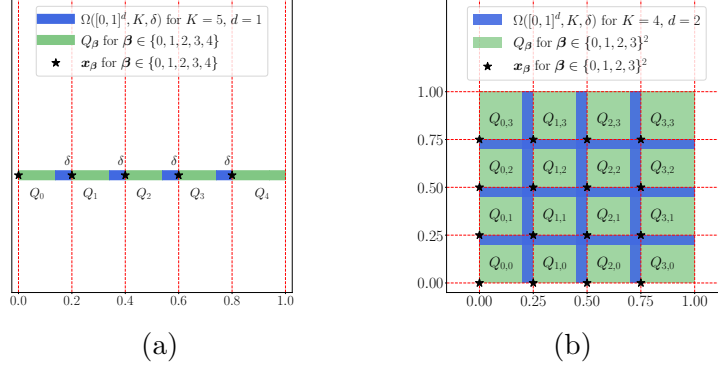


Figure 5: Illustrations of $\Omega([0, 1]^d, K, \delta)$, Q_β , and \mathbf{x}_β for $\beta \in \{0, 1, \dots, K-1\}^d$. (a) $K = 5$ and $d = 1$. (b) $K = 4$ and $d = 2$.

By defining

$$\Phi_1(\mathbf{x}) := [\phi_1(x_1), \phi_1(x_2), \dots, \phi_1(x_d)]^T, \quad \text{for any } \mathbf{x} = [x_1, x_2, \dots, x_d]^T \in \mathbb{R}^d,$$

we have $\Phi_1(\mathbf{x}) = \beta$ if $\mathbf{x} \in Q_\beta$ for $\beta \in \{0, 1, \dots, K-1\}^d$.

Step 3: Construct ϕ_2 mapping β approximately to $\tilde{f}(\mathbf{x}_\beta)$.

The construction of the sub-network implementing ϕ_2 is essentially based on Proposition 3.2. To meet the requirements of applying Proposition 3.2, we first define two auxiliary set \mathcal{A}_1 and \mathcal{A}_2 as

$$\mathcal{A}_1 := \left\{ \frac{i}{K^{d-1}} + \frac{k}{2K^d} : i = 0, 1, \dots, K^{d-1}-1 \quad \text{and} \quad k = 0, 1, \dots, K-1 \right\}$$

and

$$\mathcal{A}_2 := \left\{ \frac{i}{K^{d-1}} + \frac{K+k}{2K^d} : i = 0, 1, \dots, K^{d-1}-1 \quad \text{and} \quad k = 0, 1, \dots, K-1 \right\}.$$

Clearly, $\mathcal{A}_1 \cup \mathcal{A}_2 \cup \{1\} = \left\{ \frac{j}{2K^d} : j = 0, 1, \dots, 2K^d \right\}$ and $\mathcal{A}_1 \cap \mathcal{A}_2 = \emptyset$. See Figure 5 for an illustration of \mathcal{A}_1 and \mathcal{A}_2 . Next, we further divide this step into three sub-steps.

Step 3.1: Construct ψ_1 bijectively mapping $\{0, 1, \dots, K-1\}^d$ to \mathcal{A}_1 .

Inspired by the binary representation, we define

$$\psi_1(\mathbf{x}) := \frac{x_d}{2K^d} + \sum_{i=1}^{d-1} \frac{x_i}{K^i}, \quad \text{for any } \mathbf{x} = [x_1, x_2, \dots, x_d]^T \in \mathbb{R}^d. \quad (3.1)$$

Then ψ_1 is a linear function bijectively mapping the index set $\{0, 1, \dots, K-1\}^d$ to

$$\begin{aligned} & \left\{ \frac{\beta_d}{2K^d} + \sum_{i=1}^{d-1} \frac{\beta_i}{K^i} : \beta \in \{0, 1, \dots, K-1\}^d \right\} \\ &= \left\{ \frac{i}{K^{d-1}} + \frac{k}{2K^d} : i = 0, 1, \dots, K^{d-1}-1 \quad \text{and} \quad k = 0, 1, \dots, K-1 \right\} = \mathcal{A}_1. \end{aligned}$$

Step 3.2: Construct g to satisfy $g \circ \psi_1(\beta) = \tilde{f}(\mathbf{x}_\beta)$ and to meet the requirements of applying Proposition 3.2.

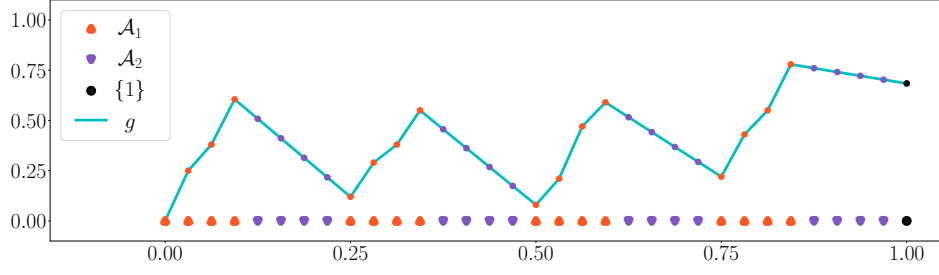


Figure 6: An illustration of \mathcal{A}_1 , \mathcal{A}_2 , $\{1\}$, and g for $d = 2$ and $K = 4$.

Let $g : [0, 1] \rightarrow \mathbb{R}$ be a continuous piecewise linear function with a set of breakpoints $\{\frac{j}{2K^d} : j = 0, 1, \dots, 2K^d\} = \mathcal{A}_1 \cup \mathcal{A}_2 \cup \{1\}$ and the values of g at these breakpoints satisfy the following properties:

- The values of g at the breakpoints in \mathcal{A}_1 are set as

$$g(\psi_1(\beta)) = \tilde{f}(\mathbf{x}_\beta), \quad \text{for any } \beta \in \{0, 1, \dots, K-1\}^d; \quad (3.2)$$

- At the breakpoint 1, let $g(1) = \tilde{f}(\mathbf{1})$, where $\mathbf{1} = [1, 1, \dots, 1]^T \in \mathbb{R}^d$;
- The values of g at the breakpoints in \mathcal{A}_2 are assigned to reduce the variation of g , which is a requirement of applying Proposition 3.2. Note that

$$\left\{ \frac{i}{K^{d-1}} - \frac{K+1}{2K^d}, \frac{i}{K^{d-1}} \right\} \subseteq \mathcal{A}_1 \cup \{1\}, \quad \text{for } i = 1, 2, \dots, K^{d-1},$$

implying the values of g at $\frac{i}{K^{d-1}} - \frac{K+1}{2K^d}$ and $\frac{i}{K^{d-1}}$ have been assigned for $i = 1, 2, \dots, K^{d-1}$. Thus, the values of g at the breakpoints in \mathcal{A}_2 can be successfully assigned by letting g linear on each interval $[\frac{i}{K^{d-1}} - \frac{K+1}{2K^d}, \frac{i}{K^{d-1}}]$ for $i = 1, 2, \dots, K^{d-1}$, since $\mathcal{A}_2 \subseteq \cup_{i=1}^{K^{d-1}} [\frac{i}{K^{d-1}} - \frac{K+1}{2K^d}, \frac{i}{K^{d-1}}]$.

Apparently, such a function g exists (see Figure 6 for an example) and satisfies

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

and

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

Step 3.3: Construct ψ_2 approximating g well on $\mathcal{A}_1 \cup \mathcal{A}_2 \cup \{1\}$.

Note that

$$2K^d = 2\left(\lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor^2 \lfloor n^{1/d} \rfloor\right)^d \leq 2(N^2 L^2 n) \leq N^2 \lceil \sqrt{2}L \rceil^2 \lceil \log_3(N+2) \rceil.$$

By Proposition 3.2 (set $y_j = g(\frac{j}{2K^2})$ and $\varepsilon = \omega_f(\frac{\sqrt{d}}{K}) > 0$ therein), there exists

$$\tilde{\psi}_2 \in \mathcal{NN}(\#input = 1; \text{ width } \leq 16N + 30; \text{ depth } \leq 6\lceil \sqrt{2}L \rceil + 10; \#output = 1)$$

505 such that

$$506 \quad |\tilde{\psi}_2(j) - g(\frac{j}{2K^d})| \leq \omega_f(\frac{\sqrt{d}}{K}), \quad \text{for } j = 0, 1, \dots, 2K^d - 1,$$

507 and

$$508 \quad 0 \leq \tilde{\psi}_2(x) \leq \max\{g(\frac{j}{2K^d}) : j = 0, 1, \dots, 2K^d - 1\} \leq 2\omega_f(\sqrt{d}), \quad \text{for any } x \in \mathbb{R}.$$

509 By defining $\psi_2(x) := \tilde{\psi}_2(2K^d x)$ for any $x \in \mathbb{R}$, we have $\psi_2 \in \mathcal{NN}(\#input = 1; \text{width} \leq$
 510 $16N + 30; \text{depth} \leq 6\lceil\sqrt{2}L\rceil + 10; \#output = 1)$,

$$511 \quad 0 \leq \psi_2(x) = \tilde{\psi}_2(2K^d x) \leq 2\omega_f(\sqrt{d}), \quad \text{for any } x \in \mathbb{R}, \quad (3.3)$$

512 and

$$513 \quad |\psi_2(\frac{j}{2K^d}) - g(\frac{j}{2K^d})| = |\tilde{\psi}_2(j) - g(\frac{j}{2K^d})| \leq \omega_f(\frac{\sqrt{d}}{K}), \quad \text{for } j = 0, 1, \dots, 2K^d - 1. \quad (3.4)$$

514 Let us end Step 3 by defining the desired function ϕ_2 as $\phi_2 := \psi_2 \circ \psi_1$. Note that
 515 $\psi_1 : \mathbb{R}^d \rightarrow \mathbb{R}$ is a linear function and $\psi_2 \in \mathcal{NN}(\#input = 1; \text{width} \leq 16N + 30; \text{depth} \leq$
 516 $6\lceil\sqrt{2}L\rceil + 10; \#output = 1)$. Thus, $\phi_2 \in \mathcal{NN}(\#input = 1; \text{width} \leq 16N + 30; \text{depth} \leq$
 517 $6\lceil\sqrt{2}L\rceil + 10; \#output = 1)$. By Equation (3.2) and (3.4), we have

$$518 \quad |\phi_2(\beta) - \tilde{f}(\mathbf{x}_\beta)| = |\psi_2(\psi_1(\beta)) - g(\psi_1(\beta))| \leq \omega_f(\frac{\sqrt{d}}{K}), \quad (3.5)$$

519 for any $\beta \in \{0, 1, \dots, K - 1\}^d$. Equation (3.3) and $\phi_2 = \psi_2 \circ \psi_1$ implies

$$520 \quad 0 \leq \phi_2(\mathbf{x}) \leq 2\omega_f(\sqrt{d}), \quad \text{for any } \mathbf{x} \in \mathbb{R}^d. \quad (3.6)$$

521 **Step 4:** Construct the final network to implement the desired function ϕ .

522 Define $\phi := \phi_2 \circ \Phi_1 + f(\mathbf{0}) - \omega_f(\sqrt{d})$. Since $\phi_1 \in \mathcal{NN}(\text{width} \leq 8\lceil N^{1/d} \rceil + 3; \text{depth} \leq$
 523 $2\lceil L^{1/d} \rceil + 5)$, we have $\Phi_1 \in \mathcal{NN}(\#input = d; \text{width} \leq 8d\lceil N^{1/d} \rceil + 3d; \text{depth} \leq 2L +$
 524 $5; \#output = d)$. It follows from the fact $\lceil\sqrt{2}L\rceil \leq \lceil\frac{3}{2}L\rceil \leq \frac{3}{2}L + \frac{1}{2}$ that $6\lceil\sqrt{2}L\rceil + 10 \leq 9L + 13$,
 525 implying

$$526 \quad \begin{aligned} \phi_2 &\in \mathcal{NN}(\#input = 1; \text{width} \leq 16N + 30; \text{depth} \leq 6\lceil\sqrt{2}L\rceil + 10; \#output = 1) \\ &\subseteq \mathcal{NN}(\#input = 1; \text{width} \leq 16N + 30; \text{depth} \leq 9L + 13; \#output = 1). \end{aligned}$$

527 Thus, $\phi = \phi_2 \circ \Phi_1 + f(\mathbf{0}) - \omega_f(\sqrt{d})$ is in

$$528 \quad \mathcal{NN}(\text{width} \leq \max\{8d\lceil N^{1/d} \rceil + 3d, 16N + 30\}; \text{depth} \leq (2L + 5) + (9L + 13) = 11L + 18).$$

529 Now let us estimate the approximation error. Note that $f = \tilde{f} + f(\mathbf{0}) - \omega_f(\sqrt{d})$. By
 530 Equation (3.5), for any $\mathbf{x} \in Q_\beta$ and $\beta \in \{0, 1, \dots, K - 1\}^d$, we have

$$\begin{aligned} |f(\mathbf{x}) - \phi(\mathbf{x})| &= |\tilde{f}(\mathbf{x}) - \phi_2(\Phi_1(\mathbf{x}))| = |\tilde{f}(\mathbf{x}) - \phi_2(\beta)| \\ &\leq |\tilde{f}(\mathbf{x}) - \tilde{f}(\mathbf{x}_\beta)| + |\tilde{f}(\mathbf{x}_\beta) - \phi_2(\beta)| \\ 531 \quad &\leq \omega_f(\frac{\sqrt{d}}{K}) + \omega_f(\frac{\sqrt{d}}{K}) \leq 2\omega_f\left(64\sqrt{d}\left(N^2L^2\log_3(N+2)\right)^{-1/d}\right), \end{aligned}$$

where the last inequality comes from the fact

$$K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor^2 \lfloor n^{1/d} \rfloor \geq \frac{N^{2/d} L^{2/d} n^{1/d}}{32} = \frac{N^{2/d} L^{2/d} \lfloor \log_3(N+2) \rfloor^{1/d}}{32} \geq \frac{(N^2 L^2 \log_3(N+2))^{1/d}}{64},$$

for any $N, L \in \mathbb{N}^+$. Recall the fact $\omega_f(j \cdot r) \leq j \cdot \omega_f(r)$ for any $j \in \mathbb{N}^+$ and $r \in [0, \infty)$.

Therefore, for any $\mathbf{x} \in \bigcup_{\beta \in \{0,1,\dots,K-1\}^d} Q_\beta = [0,1]^d \setminus \Omega([0,1]^d, K, \delta)$, we have

$$\begin{aligned} |f(\mathbf{x}) - \phi(\mathbf{x})| &\leq 2\omega_f\left(64\sqrt{d}(N^2 L^2 \log_3(N+2))^{-1/d}\right) \\ &\leq 2\lfloor 64\sqrt{d} \rfloor \omega_f\left((N^2 L^2 \log_3(N+2))^{-1/d}\right) \\ &\leq 130\sqrt{d} \omega_f\left((N^2 L^2 \log_3(N+2))^{-1/d}\right). \end{aligned}$$

It remains to show the upper bound of ϕ . By Equation (3.6) and $\phi = \phi_2 \circ \Phi_1 + f(\mathbf{0}) - \omega_f(\sqrt{d})$, it holds that $\|\phi\|_{L^\infty(\mathbb{R}^d)} \leq |f(\mathbf{0})| + \omega_f(\sqrt{d})$. Thus, we finish the proof.

4 Proofs of propositions in Section 3

In this section, we will prove the propositions in Section 3. We first introduce several basic results of ReLU networks. Next, we prove Proposition 3.1 and 3.2 based on these basic results.

4.1 Basic results of ReLU networks

To simplify the proofs of two propositions in Section 3, we introduce three lemmas below, which are basic results of ReLU networks

Lemma 4.1. *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 $\phi \in \mathcal{NN}(\text{\#input} = 1; \text{widthvec} = [2N_1, 2N_2 + 1]; \text{\#output} = 1)$ satisfying the following conditions.*

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

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

Lemma 4.2. *Given any $N, L, d \in \mathbb{N}^+$, it holds that*

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

Lemma 4.3. *Given any $L \in \mathbb{N}^+$, for arbitrary $N_1, N_2, \dots, N_L \in \mathbb{N}^+$, it holds that*

$$\mathcal{NN}(\text{\#input} = 1; \text{widthvec} = [N_1, N_2, \dots, N_L]; \text{\#output} = 1) \subseteq \text{CPwL}\left(\mathbb{R}, \prod_{\ell=1}^L (2N_\ell)\right). \quad (4.1)$$

Lemma 4.1 is a part of Theorem 3.2 in [27] or Lemma 2.2 in [21]. Lemma 4.1 is Theorem 3.1 in [27] or Lemma 3.4 in [21]. Now let us present the proof of Lemma 4.3.

557 *Proof of Lemma 4.3.* We use the mathematics induction to prove this lemma. First,
 558 consider the case $L = 1$. Given any $f \in \mathcal{NN}(\#input = 1; \text{widthvec} = [N_1]; \#output = 1)$,
 559 there exist $w_{0,j}, b_{0,j}, w_{1,j}, b_1$ for $j = 1, 2, \dots, N_1$ such that

$$560 \quad f(x) = \sum_{j=1}^{N_1} w_{1,j} \sigma(w_{0,j}x + b_{0,j}) + b_1, \quad \text{for any } x \in \mathbb{R}.$$

561 In above equation, each σ (ReLU) produces at most one breakpoint for f . Thus, f has
 562 at most N_1 breakpoints, implying $f \in \text{CPwL}(\mathbb{R}, N_1)$. Therefore, Equation (4.1) holds for
 563 $L = 1$.

564 Now assume Equation (4.1) holds for $n = k \in \mathbb{N}^+$, we would like to show it is also
 565 true for $n = k + 1$. Given any

$$566 \quad f \in \mathcal{NN}(\#input = 1; \text{widthvec} = [N_1, N_2, \dots, N_{k+1}]; \#output = 1),$$

567 there exist $w_{k,j}, b_{k,j}, w_{k+1,j}, b_{k+1} \in \mathbb{R}$ and

$$568 \quad f_j \in \mathcal{NN}(\#input = 1; \text{widthvec} = [N_1, N_2, \dots, N_k]; \#output = 1),$$

569 for $j = 1, 2, \dots, N_{k+1}$ such that

$$570 \quad f(x) = \sum_{j=1}^{N_{k+1}} w_{k+1,j} \sigma(w_{k,j} f_j(x) + b_{k,j}) + b_{k+1}, \quad \text{for any } x \in \mathbb{R}.$$

571 It is easy to verify that

$$572 \quad g \in \text{CPwL}(\mathbb{R}, n) \implies w_2 \sigma(w_1 g + b_1) + b_2 \in \text{CPwL}(\mathbb{R}, 2n),$$

573 for any $w_1, b_1, w_2, b_2 \in \mathbb{R}$ and $n \in \mathbb{N}^+$ with $n \geq 2$. By the induction hypothesis, we have

$$574 \quad f_j \in \mathcal{NN}(\#input = 1; \text{widthvec} = [N_1, N_2, \dots, N_k]; \#output = 1) \subseteq \text{CPwL}\left(\mathbb{R}, \prod_{\ell=1}^k (2N_\ell)\right),$$

575 for $j = 1, 2, \dots, N_{k+1}$. Therefore,

$$576 \quad f(x) = \sum_{j=1}^{N_{k+1}} w_{k+1,j} \sigma(w_{k,j} f_j(x) + b_{k,j}) + b_{k+1}$$

577 has at most

$$578 \quad N_{k+1} \cdot \left(2 \prod_{\ell=1}^k (2N_\ell)\right) = \prod_{\ell=1}^{k+1} (2N_\ell)$$

579 breakpoints, implying $f \in \text{CPwL}\left(\mathbb{R}, \prod_{\ell=1}^{k+1} (2N_\ell)\right)$. Thus, Equation (4.1) holds for $L = k+1$,
 580 which means we finish the induction process. So we complete the proof. \square

581 **Lemma 4.4.** For any $n \in \mathbb{N}^+$, it holds that

$$582 \quad \text{CPwL}(\mathbb{R}, n) \subseteq \mathcal{NN}(\#input = 1; \text{widthvec} = [n + 1]; \#output = 1). \quad (4.2)$$

583 *Proof.* We use the mathematics induction to prove Equation (4.2). First, consider the
 584 case $n = 1$. Given any $f \in \text{CPwL}(\mathbb{R}, n)$, there exist $a_1, a_2, x_0 \in \mathbb{R}$ such that

$$585 \quad f(x) = \begin{cases} a_1(x - x_0) + f(x_0), & \text{if } x \geq x_0, \\ a_2(x_0 - x) + f(x_0), & \text{if } x < x_0. \end{cases}$$

586 Thus, $f(x) = a_1\sigma(x - x_0) + a_2\sigma(x_0 - x) + f(x_0)$ for any $x \in \mathbb{R}$, implying $f \in \mathcal{NN}(\#input =$
 587 $1; \text{widthvec} = [2]; \#output = 1)$. Thus, Equation (4.2) holds for $n = 1$.

588 Now assume Equation (4.2) holds for $n = k \in \mathbb{N}^+$, we would like to show it is also
 589 true for $n = k + 1$. Given any $f \in \text{CPwL}(\mathbb{R}, k + 1)$, we may assume the biggest breakpoint
 590 of f is x_0 since it is trivial for the case that f has no breakpoint. Denote the slopes of
 591 the linear pieces left and right next to x_0 by a_1 and a_2 , respectively. Define

$$592 \quad \tilde{f}(x) := f(x) - (a_2 - a_1)\sigma(x - x_0), \quad \text{for any } x \in \mathbb{R}.$$

593 Then \tilde{f} has at most k breakpoints. By the induction hypothesis, we have

$$594 \quad \tilde{f} \in \text{CPwL}(\mathbb{R}, k) \subseteq \mathcal{NN}(\#input = 1; \text{widthvec} = [k + 1]; \#output = 1).$$

595 Thus, there exist $w_{0,j}, b_{0,j}, w_{1,j}, b_1$ for $j = 1, 2, \dots, k + 1$ such that

$$596 \quad \tilde{f}(x) = \sum_{j=1}^{k+1} w_{1,j}\sigma(w_{0,j}x + b_{0,j}) + b_1, \quad \text{for any } x \in \mathbb{R}.$$

597 Therefore, for any $x \in \mathbb{R}$, we have

$$598 \quad f(x) = (a_2 - a_1)\sigma(x - x_0) + \tilde{f}(x) = (a_2 - a_1)\sigma(x - x_0) + \sum_{j=1}^k w_{1,j}\sigma(w_{0,j}x + b_{0,j}) + b_1,$$

599 implying $f \in \mathcal{NN}(\#input = 1; \text{widthvec} = [k + 1]; \#output = 1)$. Thus, Equation (4.2)
 600 holds for $k + 1$, which means we finish the induction process. So we complete the proof. \square

601 4.2 Proof of Proposition 3.1

602 The setting $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor^2 \lfloor n^{1/d} \rfloor = \mathcal{O}(N^{2/d} L^{2/d} n^{1/d})$ is not neat here, but it is
 603 very convenient for later use. Now, let us present the detailed proof of Proposition 3.1.

604 Denote $K = \widetilde{M} \cdot \widetilde{L}$, where $\widetilde{M} = \lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor$, $n = \lfloor \log_3(N + 2) \rfloor$, and $\widetilde{L} = \lfloor L^{1/d} \rfloor \lfloor n^{1/d} \rfloor$.

605 Consider the sample set

$$606 \quad \{(1, \widetilde{M} - 1), (2, 0)\} \cup \left\{ \left(\frac{m}{\widetilde{M}}, m \right) : m = 0, 1, \dots, \widetilde{M} - 1 \right\} \\ \cup \left\{ \left(\frac{m+1}{\widetilde{M}} - \delta, m \right) : m = 0, 1, \dots, \widetilde{M} - 2 \right\}.$$

607 Its size is

$$608 \quad 2\widetilde{M} + 1 = 2\lfloor N^{1/d} \rfloor^2 \lfloor L^{1/d} \rfloor + 1 = \lfloor N^{1/d} \rfloor \cdot \left((2\lfloor N^{1/d} \rfloor \lfloor L^{1/d} \rfloor - 1) + 1 \right) + 1.$$

609 By Lemma 4.1 (set $N_1 = \lfloor N^{1/d} \rfloor$ and $N_2 = 2\lfloor N^{1/d} \rfloor \lfloor L^{1/d} \rfloor - 1$ therein), there exists

$$610 \quad \phi_1 \in \mathcal{NN}(\text{widthvec} = [2\lfloor N^{1/d} \rfloor, 2(2\lfloor N^{1/d} \rfloor \lfloor L^{1/d} \rfloor - 1) + 1]) \\ = \mathcal{NN}(\text{widthvec} = [2\lfloor N^{1/d} \rfloor, 4\lfloor N^{1/d} \rfloor \lfloor L^{1/d} \rfloor - 1])$$

611 such that

612 • $\phi_1(\frac{\widetilde{M}-1}{\widetilde{M}}) = \phi_1(1) = \widetilde{M} - 1$ and $\phi_1(\frac{m}{\widetilde{M}}) = \phi_1(\frac{m+1}{\widetilde{M}} - \delta) = m$ for $m = 0, 1, \dots, \widetilde{M} - 2$.

613 • ϕ_1 is linear on $[\frac{\widetilde{M}-1}{\widetilde{M}}, 1]$ and each interval $[\frac{m}{\widetilde{M}}, \frac{m+1}{\widetilde{M}} - \delta]$ for $m = 0, 1, \dots, \widetilde{M} - 2$.

614 Then, for $m = 0, 1, \dots, \widetilde{M} - 1$, we have

$$615 \quad \phi_1(x) = m, \quad \text{for any } x \in [\frac{m}{\widetilde{M}}, \frac{m+1}{\widetilde{M}} - \delta \cdot 1_{\{m \leq \widetilde{M}-2\}}]. \quad (4.3)$$

616 Now consider the another sample set

$$617 \quad \begin{aligned} & \{(\frac{1}{\widetilde{M}}, \widetilde{L} - 1), (2, 0)\} \cup \{(\frac{\ell}{\widetilde{M}\widetilde{L}}, \ell) : \ell = 0, 1, \dots, \widetilde{L} - 1\} \\ & \cup \{(\frac{\ell+1}{\widetilde{M}\widetilde{L}} - \delta, \ell) : \ell = 0, 1, \dots, \widetilde{L} - 2\}. \end{aligned}$$

618 Its size is

$$619 \quad 2\widetilde{L} + 1 = 2\lfloor L^{1/d} \rfloor \lfloor n^{1/d} \rfloor + 1 = \lfloor n^{1/d} \rfloor \cdot ((2\lfloor L^{1/d} \rfloor - 1) + 1) + 1.$$

620 By Lemma 4.1 (set $N_1 = \lfloor n^{1/d} \rfloor$ and $N_2 = 2\lfloor L^{1/d} \rfloor - 1$ therein), there exists

$$621 \quad \begin{aligned} \phi_2 & \in \mathcal{NN}(\text{widthvec} = [2\lfloor n^{1/d} \rfloor, 2(2\lfloor L^{1/d} \rfloor - 1) + 1]) \\ & = \mathcal{NN}(\text{widthvec} = [2\lfloor n^{1/d} \rfloor, 4\lfloor L^{1/d} \rfloor - 1]) \end{aligned}$$

622 such that

623 • $\phi_2(\frac{\widetilde{L}-1}{\widetilde{M}\widetilde{L}}) = \phi_2(\frac{1}{\widetilde{M}}) = \widetilde{L} - 1$ and $\phi_2(\frac{\ell}{\widetilde{M}\widetilde{L}}) = \phi_2(\frac{\ell+1}{\widetilde{M}\widetilde{L}} - \delta) = \ell$ for $\ell = 0, 1, \dots, \widetilde{L} - 2$.

624 • ϕ_2 is linear on $[\frac{\widetilde{L}-1}{\widetilde{M}\widetilde{L}}, \frac{1}{\widetilde{M}}]$ and each interval $[\frac{\ell}{\widetilde{M}\widetilde{L}}, \frac{\ell+1}{\widetilde{M}\widetilde{L}} - \delta]$ for $\ell = 0, 1, \dots, \widetilde{L} - 2$.

625 It follows that, for $m = 0, 1, \dots, \widetilde{M} - 1$ and $\ell = 0, 1, \dots, \widetilde{L} - 1$,

$$626 \quad \phi_2(x - \frac{m}{\widetilde{M}}) = \ell, \quad \text{for any } x \in [\frac{m\widetilde{L}+\ell}{\widetilde{M}\widetilde{L}}, \frac{m\widetilde{L}+\ell+1}{\widetilde{M}\widetilde{L}} - \delta \cdot 1_{\{\ell \leq \widetilde{L}-2\}}]. \quad (4.4)$$

627 $K = \widetilde{M} \cdot \widetilde{L}$ implies any $k \in \{0, 1, \dots, K-1\}$ can be unique represented by $k = m\widetilde{L} + \ell$ for
628 $m = 0, 1, \dots, \widetilde{M} - 1$ and $\ell = 0, 1, \dots, \widetilde{L} - 1$. Then the desired function ϕ can be implemented
629 by a ReLU network shown in Figure 7.

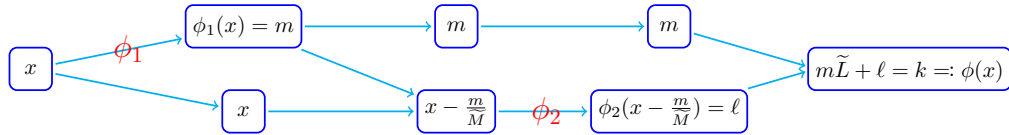


Figure 7: An illustration of the network architecture implementing ϕ based on Equation (4.3) and (4.4) for $x \in [\frac{k}{K}, \frac{k}{K} - \delta \cdot 1_{\{k \leq K-2\}}] = [\frac{m\widetilde{L}+\ell}{\widetilde{M}\widetilde{L}}, \frac{m\widetilde{L}+\ell+1}{\widetilde{M}\widetilde{L}} - \delta \cdot 1_{\{m \leq \widetilde{M}-2 \text{ or } \ell \leq \widetilde{L}-2\}}]$, where $k = m\widetilde{L} + \ell$ for $m = 0, 1, \dots, \widetilde{M} - 1$ and $\ell = 0, 1, \dots, \widetilde{L} - 1$.

630 Clearly,

$$631 \quad \phi(x) = k, \quad \text{if } x \in [\frac{k}{K}, \frac{k}{K} - \delta \cdot 1_{\{k \leq K-2\}}], \quad \text{for any } k \in \{0, 1, \dots, K-1\}.$$

632 By Lemma 4.2, we have

$$633 \quad \begin{aligned} \phi_1 &\in \mathcal{NN}(\#input = 1; \text{widthvec} = [2\lfloor N^{1/d} \rfloor, 4\lfloor N^{1/d} \rfloor \lfloor L^{1/d} \rfloor - 1]; \#output = 1) \\ &\subseteq \mathcal{NN}(\#input = 1; \text{width} \leq 8\lfloor N^{1/d} \rfloor + 2; \text{depth} \leq \lfloor L^{1/d} \rfloor + 1; \#output = 1) \end{aligned}$$

634 and

$$635 \quad \begin{aligned} \phi_2 &\in \mathcal{NN}(\#input = 1; \text{widthvec} = [2\lfloor n^{1/d} \rfloor, 4\lfloor L^{1/d} \rfloor - 1]; \#output = 1) \\ &\subseteq \mathcal{NN}(\#input = 1; \text{width} \leq 8\lfloor n^{1/d} \rfloor + 2; \text{depth} \leq \lfloor L^{1/d} \rfloor + 1; \#output = 1). \end{aligned}$$

636 Recall that $n = \lfloor \log_3(N + 2) \rfloor \leq N$. It follows from Figure 7 that ϕ can be implemented
637 by a ReLU network with width

$$638 \quad \max \{8\lfloor N^{1/d} \rfloor + 2 + 1, 8\lfloor n^{1/d} \rfloor + 2 + 1\} = 8\lfloor N^{1/d} \rfloor + 3$$

639 and depth

$$640 \quad (\lfloor L^{1/d} \rfloor + 1) + 2 + (\lfloor L^{1/d} \rfloor + 1) + 1 = 2\lfloor L^{1/d} \rfloor + 5.$$

641 So we finish the proof.

642 4.3 Proof of Proposition 3.2

643 The proof of Proposition 3.2 is based on the bit extraction technique in [3, 7]. In
644 fact, we modify this technique to extract the sum of many bits rather than one bit and
645 this modification can be summarized in Lemma 4.5 and 4.6 below.

646 **Lemma 4.5.** *For any $n \in \mathbb{N}^+$, there exists a function ϕ in*

$$647 \quad \mathcal{NN}(\#input = 2; \text{width} \leq (n + 1)2^{n+1}; \text{depth} \leq 3; \#output = 1)$$

648 *such that: Given any $\theta_j \in \{0, 1\}$ for $j = 1, 2, \dots, n$, we have*

$$649 \quad \phi(\text{bin } 0.\theta_1\theta_2\cdots\theta_n, i) = \sum_{j=1}^i \theta_j, \quad \text{for any } i \in \{0, 1, 2, \dots, n\}. \quad \textcircled{3}$$

650 *Proof.* Define $\theta = \text{bin } 0.\theta_1\theta_2\cdots\theta_n$. Clearly,

$$651 \quad \theta_j = \lfloor 2^j \theta \rfloor / 2 - \lfloor 2^{j-1} \theta \rfloor, \quad \text{for any } j \in \{1, 2, \dots, n\}.$$

652 We shall use a ReLU network to replace $\lfloor \cdot \rfloor$. Let $g \in \text{CPwL}(\mathbb{R}, 2^{n+1} - 2)$ be the function
653 matching the set of samples

$$654 \quad \bigcup_{k=0}^{2^n-1} \{(k, k), (k + 1 - \delta, k)\}, \quad \text{where } \delta = 2^{-(n+1)}.$$

655 Then $g(x) = \lfloor x \rfloor$ for any $x \in \bigcup_{k=0}^{2^n-1} [k, k + 1 - \delta]$. Note that

$$656 \quad 2^j \theta \in \bigcup_{k=0}^{2^n-1} [k, k + 1 - \delta], \quad \text{for any } j \in \{1, 2, \dots, n\}.$$

^③By convention, $\sum_{j=n}^m a_j = 0$ if $n > m$, no matter what a_j is for each j .

657 Thus,

$$658 \quad \theta_j = \lfloor 2^j \theta \rfloor / 2 - \lfloor 2^{j-1} \theta \rfloor = g(2^j \theta) / 2 - g(2^{j-1} \theta), \quad \text{for any } j \in \{1, 2, \dots, n\}. \quad (4.5)$$

659 It is easy to design a ReLU network to output $\theta_1, \theta_2, \dots, \theta_n$ by Equation (4.5) when using
 660 $\theta = \text{bin}0.\theta_1\theta_2\cdots\theta_n$ as the input. However, it is highly non-trivial to construct a ReLU
 661 network to output $\sum_{j=1}^i \theta_j$ with another input i , since many operations like multiplication
 662 and comparison are not allowed in designing ReLU networks. Now let us establish a
 663 formula to represent $\sum_{j=1}^i \theta_j$ in a form of a ReLU FNN as follows.

664 Define $\mathcal{T}(n) := \sigma(n+1) - \sigma(n) = \begin{cases} 1, & n \geq 0 \\ 0, & n < 0 \end{cases}$ for any integer n . Then, by Equation (4.5)
 665 and the fact $x_1 x_2 = \sigma(x_1 + x_2 - 1)$ for any $x_1, x_2 \in \{0, 1\}$, we have, for $i = 0, 1, 2, \dots, n$,

$$\begin{aligned} \sum_{j=1}^i \theta_j &= \sum_{j=1}^n \theta_j \cdot \mathcal{T}(i-j) = \sum_{j=1}^n \theta_j \cdot (\sigma(i-j+1) - \sigma(i-j)) \\ 666 \quad &= \sum_{j=1}^n \sigma(\theta_j + \sigma(i-j+1) - \sigma(i-j) - 1) \\ &= \sum_{j=1}^n \sigma(g(2^j \theta) / 2 - g(2^{j-1} \theta) + \sigma(i-j+1) - \sigma(i-j) - 1). \end{aligned}$$

667 Define

$$668 \quad z_{i,j} := \sigma(g(2^j \theta) / 2 - g(2^{j-1} \theta) + \sigma(i-j+1) - \sigma(i-j) - 1), \quad (4.6)$$

669 for any $i, j \in \{1, 2, \dots, n\}$. Then the goal is to design ϕ satisfying

$$670 \quad \phi(\theta, i) = \sum_{j=1}^i \theta_j = \sum_{j=1}^n z_{i,j}, \quad \text{for any } i \in \{0, 1, 2, \dots, n\}. \quad (4.7)$$

671 See Figure 8 for the network architecture implementing the desired function ϕ .

672 By Lemma 4.4, we have

$$673 \quad g \in \text{CPwL}(\mathbb{R}, 2^{n+1} - 2) \subseteq \mathcal{NN}(\# \text{input} = 1; \text{widthvec} = [2^{n+1} - 1]; \# \text{output} = 1),$$

674 implying

$$675 \quad g(2^j \cdot) \in \text{CPwL}(\mathbb{R}, 2^{n+1} - 2) \subseteq \mathcal{NN}(\# \text{input} = 1; \text{widthvec} = [2^{n+1} - 1]; \# \text{output} = 1),$$

676 for any $j = 0, 1, 2, \dots, n$. Clearly, the network in Figure 8 has width $(n+1)(2^{n+1} - 1) +$
 677 $(n+1) = (n+1)2^{n+1}$ and depth 3. So we finish the proof. \square

678 **Lemma 4.6.** *For any $n, L \in \mathbb{N}^+$, there exists a function ϕ in*

$$679 \quad \mathcal{NN}(\# \text{input} = 2; \text{width} \leq (n+3)2^{n+1} + 4; \text{depth} \leq 4L + 2; \# \text{output} = 1)$$

680 *such that: Given any $\theta_j \in \{0, 1\}$ for $j = 1, 2, \dots, Ln$, we have*

$$681 \quad \phi(\text{bin}0.\theta_1\theta_2\cdots\theta_{Ln}, k) = \sum_{j=1}^k \theta_j, \quad \text{for any } k \in \{1, 2, \dots, Ln\}.$$

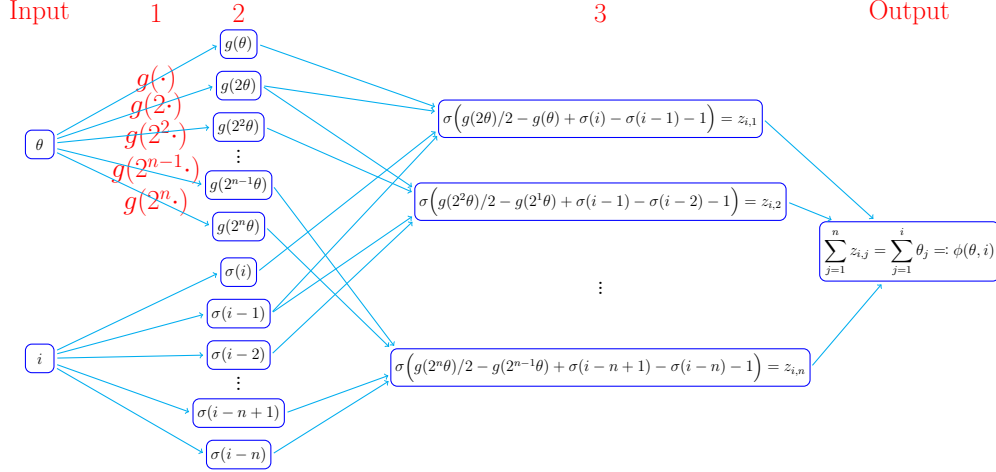


Figure 8: An illustration of the network implementing the desired function ϕ with the input $[\theta, i]^T = [\text{bin}0.\theta_1\theta_2\cdots\theta_n, i]^T$ for any $i \in \{0, 1, 2, \dots, n\}$ and $\theta_1, \theta_2, \dots, \theta_n \in \{0, 1\}$. $g(2^j \cdot)$ can be implemented by a one-hidden-layer network with width $2^{n+1} - 1$ for each $j \in \{0, 1, \dots, n\}$. The red numbers above the architecture indicate the order of hidden layers. The network architecture is essentially determined by Equation (4.6) and (4.7), which are valid no matter what $\theta_1, \theta_2, \dots, \theta_n \in \{0, 1\}$. Thus, the desired function ϕ is independent of $\theta_1, \theta_2, \dots, \theta_n \in \{0, 1\}$. We omit ReLU (σ) for a neuron if its output is non-negative without ReLU. Such a simplification are applied to similar figures in this paper.

682 *Proof.* Let $g_1 \in \text{CPwL}(\mathbb{R}, 2^{n+1} - 2)$ be the function matching the set of samples

$$683 \quad \bigcup_{i=0}^{2^n-1} \{(i, i), (i+1-\delta, i)\}, \quad \text{where } \delta = 2^{-(L_{n+1})}.$$

684 Then $g_1(x) = \lfloor x \rfloor$ for any $x \in \bigcup_{i=0}^{2^n-1} [i, i+1-\delta]$. Note that

$$685 \quad 2^n \cdot \text{bin}0.\theta_{\ell_{n+1}}\cdots\theta_{L_n} \in \bigcup_{i=0}^{2^n-1} [i, i+1-\delta], \quad \text{for any } \ell \in \{0, 1, \dots, L-1\}.$$

686 Thus, for any $\ell \in \{0, 1, \dots, L-1\}$, we have

$$687 \quad \text{bin}0.\theta_{\ell_{n+1}}\cdots\theta_{\ell_{n+n}} = \frac{\lfloor 2^n \cdot \text{bin}0.\theta_{\ell_{n+1}}\cdots\theta_{L_n} \rfloor}{2^n} = \frac{g_1(2^n \cdot \text{bin}0.\theta_{\ell_{n+1}}\cdots\theta_{L_n})}{2^n}. \quad (4.8)$$

688 Define $g_2(x) := 2^n x - g_1(2^n x)$ for any $x \in \mathbb{R}$. Then $g_2 \in \text{CPwL}(\mathbb{R}, 2^{n+1} - 2)$ and

$$689 \quad \begin{aligned} & \text{bin}0.\theta_{(\ell+1)_{n+1}}\cdots\theta_{L_n} = 2^n \left(\text{bin}0.\theta_{\ell_{n+1}}\cdots\theta_{L_n} - \text{bin}0.\theta_{\ell_{n+1}}\cdots\theta_{\ell_{n+n}} \right) \\ & = 2^n \left(\text{bin}0.\theta_{\ell_{n+1}}\cdots\theta_{L_n} - \frac{g_1(2^n \cdot \text{bin}0.\theta_{\ell_{n+1}}\cdots\theta_{L_n})}{2^n} \right) = g_2(\text{bin}0.\theta_{\ell_{n+1}}\cdots\theta_{L_n}). \end{aligned} \quad (4.9)$$

690 By Lemma 4.5, there exists

$$691 \quad \phi_1 \in \mathcal{NN}(\# \text{input} = 2; \text{width} \leq (n+1)2^{n+1}; \text{depth} \leq 3; \# \text{output} = 1)$$

such that: For any $\xi_1, \xi_2, \dots, \xi_n \in \{0, 1\}$, we have

$$\phi_1(\text{bin}0.\xi_1\xi_2\cdots\xi_n, i) = \sum_{j=1}^i \xi_j, \quad \text{for } i = 0, 1, 2, \dots, n.$$

It follows that

$$\phi_1(\text{bin}0.\theta_{\ell n+1}\theta_{\ell n+2}\cdots\theta_{\ell n+n}, i) = \sum_{j=1}^i \theta_{\ell n+j}, \quad \text{for } \ell = 0, 1, \dots, L-1 \text{ and } i = 0, 1, \dots, n. \quad (4.10)$$

Define $\phi_{2,\ell}(x) := \min\{\sigma(x - \ell n), n\}$ for any $x \in \mathbb{R}$ and $\ell \in \{0, 1, \dots, L-1\}$. For any $k \in \{1, 2, \dots, Ln\}$, there exists $k_1 \in \{0, 1, \dots, L-1\}$ and $k_2 \in \{1, 2, \dots, n\}$ such that $k = k_1 n + k_2$, implying

$$\begin{aligned} \sum_{i=1}^k \theta_i &= \sum_{i=1}^{k_1 n + k_2} \theta_i = \sum_{\ell=0}^{k_1-1} \left(\sum_{j=1}^n \theta_{\ell n+j} \right) + \sum_{\ell=k_1}^{k_1-1} \left(\sum_{j=1}^{k_2} \theta_{\ell n+j} \right) + \sum_{\ell=k_1+1}^{L-1} \left(\sum_{j=1}^0 \theta_{\ell n+j} \right) \\ &= \sum_{\ell=0}^{L-1} \left(\sum_{j=1}^{\min\{\sigma(k-\ell n), n\}} \theta_{\ell n+j} \right) = \sum_{\ell=0}^{L-1} \left(\sum_{j=1}^{\phi_{2,\ell}(k)} \theta_{\ell n+j} \right). \end{aligned} \quad (4.11)$$

Then, the desired function ϕ can be implemented by the network architecture in Figure 9.

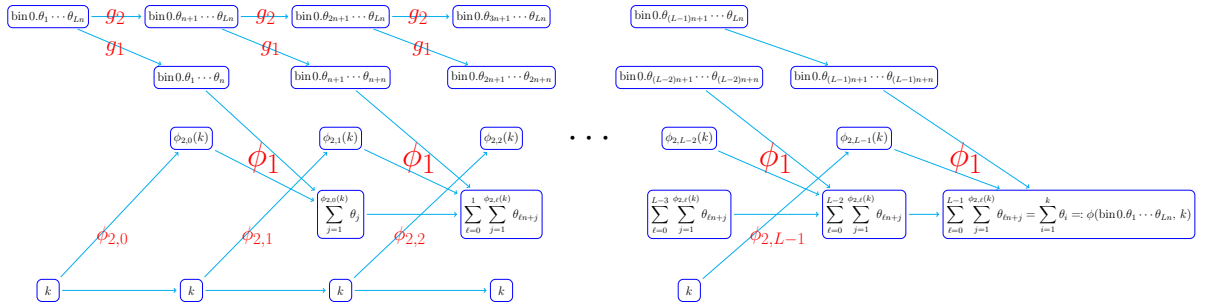


Figure 9: An illustration of the network implementing the desired function ϕ with the input $[\text{bin}0.\theta_1\theta_2\cdots\theta_{Ln}, k]^T$ for any $k \in \{1, 2, \dots, Ln\}$ and $\theta_1, \theta_2, \dots, \theta_{Ln} \in \{0, 1\}$. The network architecture is essentially determined by Equation (4.8), (4.9), (4.10), and (4.11), which are valid no matter what $\theta_1, \theta_2, \dots, \theta_{Ln} \in \{0, 1\}$. Thus, the desired function ϕ is independent of $\theta_1, \theta_2, \dots, \theta_{Ln} \in \{0, 1\}$. We omit ReLU (σ) for a neuron if its output is non-negative without ReLU.

By Lemma 4.4, we have

$$g_1, g_2 \in \text{CPwL}(\mathbb{R}, 2^{n+1} - 2) \subseteq \mathcal{NN}(\# \text{input} = 1; \text{widthvec} = [2^{n+1} - 1]; \# \text{output} = 1).$$

Recall that $\phi_1 \in \mathcal{NN}(\text{width} \leq (n+1)2^{n+1}; \text{depth} \leq 3)$. As shown in Figure 10, $\phi_{2,\ell}(x) \in \mathcal{NN}(\text{width} \leq 4; \text{depth} \leq 2)$ for $\ell = 0, 1, \dots, L-1$. Therefore, the network in Figure 9 has width

$$(2^{n+1} - 1) + (2^{n+1} - 1) + (n+1)2^{n+1} + 1 + 4 + 1 = (n+3)2^{n+1} + 4$$

and depth

$$2 + L(1 + 3) = 4L + 2.$$

So we finish the proof. \square

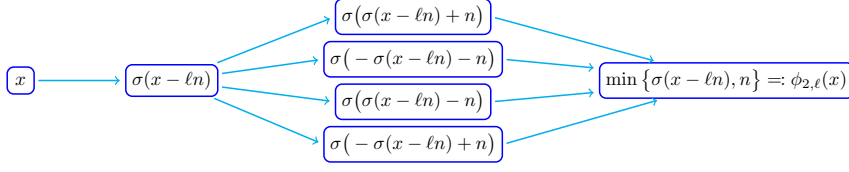


Figure 10: An illustration of the network implementing the desired function $\phi_{2,\ell}$ for each $\ell \in \{0, 1, \dots, L-1\}$, based on $\min\{x, n\} = \frac{1}{2}(\sigma(x-n) - \sigma(-x-n) - \sigma(x+n) - \sigma(-x+n))$.

Next, we introduce Lemma 4.7 to map indices to the partial sum of given bits.

Lemma 4.7. *Given any $N, L \in \mathbb{N}^+$ and arbitrary $\theta_{m,k} \in \{0, 1\}$ for $m = 0, 1, \dots, M-1$ and $k = 0, 1, \dots, Ln-1$, where $M = N^2L$ and $n = \lfloor \log_3(N+2) \rfloor$, there exists*

$$\phi \in \mathcal{NN}(\#input = 2; \text{ width } \leq 6N + 14; \text{ depth } \leq 5L + 4; \#output = 1)$$

such that

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

Proof. Define

$$y_m := \text{bin}0.\theta_{m,0}\theta_{m,1}\dots\theta_{m,Ln-1}, \quad \text{for } m = 0, 1, \dots, M-1.$$

Consider the sample set $\{(m, y_m) : m = 0, 1, \dots, M\}$, whose cardinality is

$$M + 1 = N((NL - 1) + 1) + 1.$$

By Lemma 4.1 (set $N_1 = N$ and $N_2 = NL - 1$ therein), there exists

$$\begin{aligned} \phi_1 &\in \mathcal{NN}(\#input = 1; \text{ widthvec } = [2N, 2(NL - 1) + 1]; \#output = 1) \\ &= \mathcal{NN}(\#input = 1; \text{ widthvec } = [2N, 2NL - 1]; \#output = 1) \end{aligned}$$

such that

$$\phi_1(m) = y_m, \quad \text{for } m = 0, 1, \dots, M-1.$$

By Lemma 4.5, there exists

$$\phi_2 \in \mathcal{NN}(\#input = 2; \text{ width } \leq (n+3)2^{n+1} + 4; \text{ depth } \leq 4L + 2; \#output = 1)$$

such that, for any $\xi_1, \xi_2, \dots, \xi_{Ln} \in \{0, 1\}$, we have

$$\phi_2(\text{bin}0.\xi_1\xi_2\dots\xi_{Ln}, k) = \sum_{j=1}^k \xi_j, \quad \text{for } k = 1, 2, \dots, Ln.$$

It follows that, for any $\xi_0, \xi_1, \dots, \xi_{Ln-1} \in \{0, 1\}$, we have

$$\phi_2(\text{bin}0.\xi_0\xi_1\dots\xi_{Ln-1}, k+1) = \sum_{j=0}^k \xi_j, \quad \text{for } k = 0, 1, \dots, Ln-1.$$

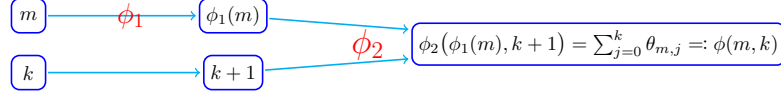


Figure 11: An illustration of the network implementing the desired function ϕ for $m = 0, 1, \dots, M-1$ and $k = 0, 1, \dots, Ln-1$.

Thus, for $m = 0, 1, \dots, M-1$ and $k = 0, 1, \dots, Ln-1$, we have

$$\phi_2(\phi_1(m), k+1) = \phi_2(y_m, k+1) = \phi_2(0.\theta_{m,0}\theta_{m,1}\dots\theta_{m,L-1}, k+1) = \sum_{j=0}^k \theta_{m,j}.$$

Hence, the desired function ϕ can be implemented by the network shown in Figure 11. By Lemma 4.2, $\phi_1 \in \mathcal{NN}(\text{widthvec} = [2N, 2NL-1]) \subseteq \mathcal{NN}(\text{width} \leq 4N+2; \text{depth} \leq L+1)$. It holds that

$$(n+3)2^{n+1} + 4 \leq 6 \cdot (3^n) + 2 = 6 \cdot (3^{\lfloor \log_3(N+2) \rfloor}) + 2 \leq 6(N+2) + 2 = 6N+14,$$

implying

$$\begin{aligned} \phi_2 &\in \mathcal{NN}(\#input = 2; \text{width} \leq (n+3)2^{n+1} + 4; \text{depth} \leq 4L+2; \#output = 1) \\ &\subseteq \mathcal{NN}(\#input = 2; \text{width} \leq 6N+14; \text{depth} \leq 4L+2; \#output = 1). \end{aligned}$$

Therefore, the network in Figure 11 is with width $\max\{(4N+2)+1, 6N+14\} = 6N+14$ and depth $(4L+2)+1+(L+1) = 5L+4$. So we finish the proof. \square

Next, we apply Lemma 4.7 to prove Lemma 4.8 below, which is a key intermediate conclusion to prove Proposition 3.2.

Lemma 4.8. For any $\varepsilon > 0$ and $N, L \in \mathbb{N}^+$, denote $M = N^2L$ and $n = \lfloor \log_3(N+2) \rfloor$. Assume $y_{m,k} \geq 0$ for $m = 0, 1, \dots, M-1$ and $k = 0, 1, \dots, Ln-1$ are samples with

$$|y_{m,k} - y_{m,k-1}| \leq \varepsilon, \quad \text{for } m = 0, 1, \dots, M-1 \quad \text{and} \quad k = 1, 2, \dots, Ln-1.$$

Then there exists $\phi \in \mathcal{NN}(\#input = 2; \text{width} \leq 16N+30; \text{depth} \leq 5L+7; \#output = 1)$ such that

(i) $|\phi(m, k) - y_{m,k}| \leq \varepsilon$ for $m = 0, 1, \dots, M-1$ and $k = 0, 1, \dots, Ln-1$;

(ii) $0 \leq \phi(x_1, x_2) \leq \max\{y_{m,k} : m = 0, 1, \dots, M-1 \text{ and } k = 0, 1, \dots, Ln-1\}$ for any $x_1, x_2 \in \mathbb{R}$.

Proof. Define

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

We will construct a function implemented by a ReLU network to map the index (m, k) to $a_{m,k}\varepsilon$ for $m = 0, 1, \dots, M-1$ and $k = 0, 1, \dots, Ln-1$.

754 Define $b_{m,0} := 0$ and $b_{m,k} := a_{m,k} - a_{m,k-1}$ for $m = 0, 1, \dots, M-1$ and $k = 1, 2, \dots, Ln-1$.
 755 Since $|y_{m,k} - y_{m,k-1}| \leq \varepsilon$ for all m and k , we have $b_{m,k} \in \{-1, 0, 1\}$. Hence, there exist $c_{m,k}$
 756 and $d_{m,k} \in \{0, 1\}$ such that $b_{m,k} = c_{m,k} - d_{m,k}$, which implies

$$\begin{aligned} a_{m,k} &= a_{m,0} + \sum_{j=1}^k (a_{m,j} - a_{m,j-1}) = a_{m,0} + \sum_{j=1}^k b_{m,j} = a_{m,0} + \sum_{j=0}^k b_{m,j} \\ &= a_{m,0} + \sum_{j=0}^k c_{m,j} - \sum_{j=0}^k d_{m,j}, \end{aligned}$$

758 for $m = 0, 1, \dots, M-1$ and $k = 0, 1, \dots, Ln-1$.

759 Consider the sample set

$$760 \quad \{(m, a_{m,0}) : m = 0, 1, \dots, M-1\} \cup \{(M, 0)\}.$$

761 Its size is $M+1 = N \cdot ((NL-1)+1) + 1$, by Lemma 4.1 (set $N_1 = N$ and $N_2 = NL-1$
 762 therein), there exists

$$763 \quad \psi_1 \in \mathcal{NN}(\text{widthvec} = [2N, 2(NL-1)+1]) = \mathcal{NN}(\text{widthvec} = [2N, 2NL-1])$$

764 such that

$$765 \quad \psi_1(m) = a_{m,0}, \quad \text{for } m = 0, 1, \dots, M-1.$$

766 By Lemma 4.7, there exist $\psi_2, \psi_3 \in \mathcal{NN}(\text{width} \leq 6N+14; \text{depth} \leq 5L+4)$ such that

$$767 \quad \psi_2(m, k) = \sum_{j=0}^k c_{m,j} \quad \text{and} \quad \psi_3(m, k) = \sum_{j=0}^k d_{m,j},$$

768 for $m = 0, 1, \dots, M-1$ and $k = 0, 1, \dots, Ln-1$. Hence, it holds that

$$769 \quad a_{m,k} = a_{m,0} + \sum_{j=0}^k c_{m,j} - \sum_{j=0}^k d_{m,j} = \psi_1(m) + \psi_2(m, k) - \psi_3(m, k), \quad (4.12)$$

770 for $m = 0, 1, \dots, M-1$ and $k = 0, 1, \dots, Ln-1$.

771 Define

$$772 \quad y_{\max} := \max\{y_{m,k} : m = 0, 1, \dots, M-1 \text{ and } k = 0, 1, \dots, Ln-1\}.$$

773 Then the desired function can be implemented by two sub-networks shown in Figure 12.

774 By Lemma 4.2,

$$\begin{aligned} 775 \quad \psi_1 &\in \mathcal{NN}(\#input = 1; \text{widthvec} = [2N, 2NL-1]; \#output = 1) \\ &\subseteq \mathcal{NN}(\#input = 1; \text{width} \leq 4N+2; \text{depth} \leq L+1; \#output = 1). \end{aligned}$$

776 Recall that $\psi_2, \psi_3 \in \mathcal{NN}(\text{width} \leq 6N+14; \text{depth} \leq 5L+4)$. Thus, $\phi_1 \in \mathcal{NN}(\text{width} \leq$
 777 $(4N+2) + 2(6N+14) = 16N+30; \text{depth} \leq (5L+4) + 1 = 5L+5)$ as shown in Figure 12.
 778 And it is clear that $\phi_2 \in \mathcal{NN}(\text{width} \leq 4; \text{depth} \leq 2)$, implying $\phi = \phi_2 \circ \phi_1 \in \mathcal{NN}(\text{width} \leq$
 779 $16N+30; \text{depth} \leq (5L+5) + 2 = 5L+7)$.

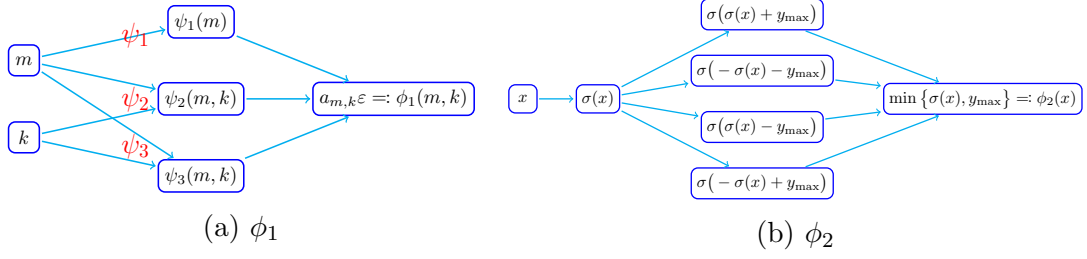


Figure 12: Illustrations of two sub-networks implementing the desired function $\phi = \phi_2 \circ \phi_1$ for $m = 0, 1, \dots, M-1$ and $k = 0, 1, \dots, Ln-1$, based on Equation (4.12) and the fact $\min\{x_1, x_2\} = \frac{x_1+x_2-|x_1-x_2|}{2} = \frac{\sigma(x_1+x_2)-\sigma(-x_1-x_2)-\sigma(x_1-x_2)-\sigma(-x_1+x_2)}{2}$.

Clearly, $0 \leq \phi(x_1, x_2) \leq y_{\max}$ for any $x_1, x_2 \in \mathbb{R}$, since $\phi(x_1, x_2) = \phi_2 \circ \phi_1(x_1, x_2) = \max\{\sigma(\phi_1(x_1, x_2)), y_{\max}\}$.

Note that $0 \leq a_{m,k}\varepsilon = \lfloor y_{m,k}/\varepsilon \rfloor \varepsilon \leq y_{\max}$. Then we have $\phi(m, k) = \phi_2 \circ \phi_1(m, k) = \phi_2(a_{m,k}\varepsilon) = \max\{\sigma(a_{m,k}\varepsilon), y_{\max}\} = a_{m,k}\varepsilon$. Therefore,

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

for $m = 0, 1, \dots, M-1$ and $k = 0, 1, \dots, Ln-1$. Hence, we finish the proof. \square

Finally, we apply Lemma 4.8 to prove Proposition 3.2.

Proof of Proposition 3.2. Denote $M = N^2L$, $n = \lfloor \log_3(N+2) \rfloor$, and $\widehat{L} = Ln$. We may assume $J = MLn = M\widehat{L}$ since we can set $y_{J-1} = y_J = y_{J+1} = \dots = y_{M\widehat{L}-1}$ if $J < M\widehat{L}$.

Consider the sample set

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

Its size is $2M + 1 = N \cdot ((2NL - 1) + 1) + 1$. By Lemma 4.1 (set $N_1 = N$ and $N_2 = NL - 1$ therein), there exist

$$\phi_1 \in \mathcal{NN}(\text{widthvec} = [2N, 2(2NL - 1) + 1]) = \mathcal{NN}(\text{widthvec} = [2N, 4NL - 1])$$

such that

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

It follows that

$$\phi_1(j) = m, \quad \text{and} \quad j - \widehat{L}\phi_1(j) = k, \quad \text{where } j = m\widehat{L} + k, \quad (4.13)$$

for $m = 0, 1, \dots, M-1$ and $k = 0, 1, \dots, \widehat{L} - 1$.

Note that any number j in $\{0, 1, \dots, J-1\}$ can be uniquely indexed as $j = m\widehat{L} + k$ for $m = 0, 1, \dots, M-1$ and $k = 0, 1, \dots, \widehat{L} - 1$. So we can denote $y_j = y_{m\widehat{L}+k}$ as $y_{m,k}$. Then by Lemma 4.8, there exists $\phi_2 \in \mathcal{NN}(\text{width} \leq 12N + 8; \text{depth} \leq 3L + 6)$ such that

$$|\phi_2(m, k) - y_{m,k}| \leq \varepsilon, \quad \text{for } m = 0, 1, \dots, M-1 \quad \text{and} \quad k = 0, 1, \dots, \widehat{L} - 1, \quad (4.14)$$

804 and

$$805 \quad 0 \leq \phi_2(x_1, x_2) \leq y_{\max}, \quad \text{for any } x_1, x_2 \in \mathbb{R}, \quad (4.15)$$

806 where $y_{\max} := \max\{y_{m,k} : m = 0, 1, \dots, M-1 \text{ and } k = 0, 1, \dots, \widehat{L}-1\} = \max\{y_j : j =$
 807 $0, 1, \dots, M\widehat{L}-1\}$.

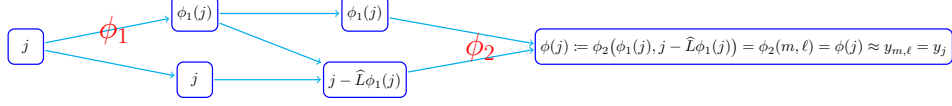


Figure 13: An illustration of the ReLU network implementing the desired function ϕ based Equation (4.13). The index $j \in \{0, 1, \dots, M\widehat{L}-1\}$ is uniquely represented by $j = mL+k$ for $m = 0, 1, \dots, M-1$ and $k = 0, 1, \dots, \widehat{L}-1$.

808 By Lemma 4.2,

$$809 \quad \begin{aligned} &\phi_1 \in \mathcal{NN}(\#input = 1; \text{widthvec} = [2N, 4NL - 1]; \#output = 1) \\ &\subseteq \mathcal{NN}(\#input = 1; \text{width} \leq 8N + 2; \text{depth} \leq L + 1; \#output = 1). \end{aligned}$$

810 Recall that $\phi_2 \in \mathcal{NN}(\text{width} \leq 16N + 30; \text{depth} \leq 5L + 7)$. So $\phi \in \mathcal{NN}(\text{width} \leq 16N +$
 811 $30; \text{depth} \leq (L + 1) + 2 + (5L + 7) = 6L + 10)$ as shown in Figure 13.

812 Equation (4.15) implies

$$813 \quad 0 \leq \phi(x) \leq y_{\max}, \quad \text{for any } x \in \mathbb{R},$$

814 since ϕ is given by $\phi(x) = \phi_2(\phi_1(x), x - L\phi_1(x))$.

815 Represent $j \in \{0, 1, \dots, M\widehat{L}-1\}$ via $j = mL + k$ for $m = 0, 1, \dots, M-1$ and $k =$
 816 $0, 1, \dots, \widehat{L}-1$. Then, by Equation (4.14), we have

$$817 \quad |\phi(j) - y_j| = |\phi_2(\phi_1(j), j - L\phi_1(j)) - y_j| = |\phi_2(m, k) - y_{m,k}| \leq \varepsilon,$$

818 for any $j \in \{0, 1, \dots, M\widehat{L}-1\} = \{0, 1, \dots, J-1\}$. So we finish the proof. \square

819 We would like to remark that the key idea in the proof of Proposition 3.2 is the bit
 820 extraction technique in Lemma 4.6, which allows us to store Ln bits in a binary number
 821 $\text{bin}0.\theta_1\theta_2\cdots\theta_{Ln}$ and extract each bit θ_i . The extraction operator can be efficiently carried
 822 out via a deep ReLU neural network demonstrating the power of depth.

823 5 Conclusion and future work

824 This paper aims at a quantitative and optimal approximation rate for ReLU net-
 825 works in terms of the width and depth to approximate continuous functions. It is
 826 shown by construction that ReLU networks with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ can
 827 approximate an arbitrary continuous function on $[0, 1]$ with an approximation rate
 828 $\mathcal{O}(\omega_f((N^2L^2 \ln N)^{-1/d}))$. By connecting the approximation property to VC-dimension,
 829 we prove that such a rate is optimal for Hölder continuous functions on $[0, 1]^d$ in terms

of the width and depth separately, and hence this rate is also optimal for the whole continuous function class. We also extend our analysis to general continuous functions on any bounded set in \mathbb{R}^d . 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 extend our conclusions to neural networks with other types of architectures (e.g., convolutional neural networks) and activation functions (e.g., tanh and sigmoid functions).

Acknowledgments

Z. Shen is supported by Tan Chin Tuan Centennial Professorship. H. Yang was partially supported by the US National Science Foundation under award DMS-1945029. S. Zhang is supported by a Postdoctoral Fellowship under NUS ENDOWMENT FUND (EXP WBS) (01 651).

References

- [1] M. ANTHONY AND P. L. BARTLETT, *Neural Network Learning: Theoretical Foundations*, Cambridge University Press, New York, NY, USA, 1st ed., 2009.
- [2] A. R. BARRON, *Universal approximation bounds for superpositions of a sigmoidal function*, IEEE Transactions on Information Theory, 39 (1993), pp. 930–945.
- [3] P. BARTLETT, V. MAIOROV, AND R. MEIR, *Almost linear VC dimension bounds for piecewise polynomial networks*, Neural Computation, 10 (1998), pp. 2159–2173.
- [4] 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.
- [5] G. CYBENKO, *Approximation by superpositions of a sigmoidal function*, MCSS, 2 (1989), pp. 303–314.
- [6] W. E AND Q. WANG, *Exponential convergence of the deep neural network approximation for analytic functions*, CoRR, abs/1807.00297 (2018).
- [7] 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.
- [8] K. HORNIK, M. STINCHCOMBE, AND H. WHITE, *Multilayer feedforward networks are universal approximators*, Neural Networks, 2 (1989), pp. 359 – 366.

- [9] 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.
- [10] K. KAWAGUCHI AND Y. BENGIO, *Depth with nonlinearity creates no bad local minima in resnets*, (2018).
- [11] M. J. KEARNS AND R. E. SCHAPIRE, *Efficient distribution-free learning of probabilistic concepts*, J. Comput. Syst. Sci., 48 (1994), pp. 464–497.
- [12] Q. LI, T. LIN, AND Z. SHEN, *Deep learning via dynamical systems: An approximation perspective*, arXiv e-prints, (2019).
- [13] S. LIANG AND R. SRIKANT, *Why deep neural networks?*, CoRR, abs/1610.04161 (2016).
- [14] J. LU, Z. SHEN, H. YANG, AND S. ZHANG, *Deep network approximation for smooth functions*, arXiv e-prints, (2020), p. arXiv:2001.03040.
- [15] 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.
- [16] H. MONTANELLI, H. YANG, AND Q. DU, *Deep ReLU networks overcome the curse of dimensionality for bandlimited functions*, Journal of Computational Mathematics, (to appear).
- [17] 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.
- [18] Q. N. NGUYEN AND M. HEIN, *The loss surface of deep and wide neural networks*, CoRR, abs/1704.08045 (2017).
- [19] P. PETERSEN AND F. VOIGTLAENDER, *Optimal approximation of piecewise smooth functions using deep ReLU neural networks*, Neural Networks, 108 (2018), pp. 296 – 330.
- [20] 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.
- [21] Z. SHEN, H. YANG, AND S. ZHANG, *Nonlinear approximation via compositions*, Neural Networks, 119 (2019), pp. 74 – 84.
- [22] —, *Deep network approximation characterized by number of neurons*, Communications in Computational Physics, 28 (2020), pp. 1768–1811.

- 898 [23] Z. SHEN, H. YANG, AND S. ZHANG, *Deep network with approximation error being*
899 *reciprocal of width to power of square root of depth*, arXiv e-prints, (2020).
- 900 [24] —, *Neural network approximation: Three hidden layers are enough*, arXiv e-
901 prints, (2020).
- 902 [25] D. YAROTSKY, *Error bounds for approximations with deep ReLU networks*, Neural
903 Networks, 94 (2017), pp. 103 – 114.
- 904 [26] D. YAROTSKY, *Optimal approximation of continuous functions by very deep ReLU*
905 *networks*, in Proceedings of the 31st Conference On Learning Theory, S. Bubeck,
906 V. Perchet, and P. Rigollet, eds., vol. 75 of Proceedings of Machine Learning Re-
907 search, PMLR, 06–09 Jul 2018, pp. 639–649.
- 908 [27] S. ZHANG, *Deep neural network approximation via function compositions*, PhD
909 Thesis, National University of Singapore, (2020). URL: [https://scholarbank.](https://scholarbank.nus.edu.sg/handle/10635/186064)
910 [nus.edu.sg/handle/10635/186064](https://scholarbank.nus.edu.sg/handle/10635/186064).