

Neural Network Approximation: Three Hidden Layers Are Enough

Zuowei Shen^a, Haizhao Yang^b, Shijun Zhang^a

^a*Department of Mathematics, National University of Singapore*

^b*Department of Mathematics, Purdue University*

Abstract

A three-hidden-layer neural network with super approximation power is introduced. This network is built with the floor function ($\lfloor x \rfloor$), the exponential function (2^x), the step function ($\mathbb{1}_{x \geq 0}$), or their compositions as the activation function in each neuron and hence we call such networks as Floor-Exponential-Step (FLES) networks. For any width hyper-parameter $N \in \mathbb{N}^+$, it is shown that FLES networks with width $\max\{d, N\}$ and three hidden layers can uniformly approximate a Hölder continuous function f on $[0, 1]^d$ with an exponential approximation rate $3\lambda(2\sqrt{d})^\alpha 2^{-\alpha N}$, where $\alpha \in (0, 1]$ and $\lambda > 0$ are the Hölder order and constant, respectively. More generally for an arbitrary continuous function f on $[0, 1]^d$ with a modulus of continuity $\omega_f(\cdot)$, the constructive approximation rate is $2\omega_f(2\sqrt{d})2^{-N} + \omega_f(2\sqrt{d}2^{-N})$. Moreover, we extend such a result to general bounded continuous functions on a bounded set $E \subseteq \mathbb{R}^d$. As a consequence, this new class of networks overcomes the curse of dimensionality in approximation power when the variation of $\omega_f(r)$ as $r \rightarrow 0$ is moderate (e.g., $\omega_f(r) \lesssim r^\alpha$ for Hölder continuous functions), since the major term to be concerned in our approximation rate is essentially \sqrt{d} times a function of N independent of d within the modulus of continuity. Finally, we extend our analysis to derive similar approximation results in the L^p -norm for $p \in [1, \infty)$ via replacing Floor-Exponential-Step activation functions by continuous activation functions.

Keywords:

Exponential Convergence, Curse of Dimensionality, Deep Neural Network,

Email addresses: matzuows@nus.edu.sg (Zuowei Shen), haizhao@purdue.edu (Haizhao Yang), zhangshijun@u.nus.edu (Shijun Zhang)

1. Introduction

This paper studies the approximation power of neural networks and shows that three hidden layers are enough for neural networks to achieve super approximation capacity. In particular, leveraging the power of advanced yet simple activation functions, we will introduce new theories and network architectures with only three hidden layers achieving exponential convergence and avoiding the curse of dimensionality simultaneously for (Hölder) continuous functions with an explicit approximation bound. The theories established in this paper would provide new insights to explain why deeper neural networks are better than one-hidden-layer neural networks for large-scale and high-dimensional problems. The approximation theories here are constructive (i.e., with explicit formulas to specify network parameters) and quantitative (i.e., results valid for essentially arbitrary width and/or depth without lower bound constraints) with explicit error bounds working for three-hidden-layer networks with arbitrary width.

Constructive approximation with quantitative results and explicit error bounds would provide important guides for deciding the network sizes in deep learning. For example, the (nearly) optimal approximation rates of deep ReLU networks with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ for a Lipschitz continuous function and a C^s function f on $[0, 1]^d$ are $\mathcal{O}(\sqrt{d}N^{-2/d}L^{-2/d})$ and $\mathcal{O}(\|f\|_{C^s}(\frac{N}{\ln N})^{-2s/d}(\frac{L}{\ln L})^{-2s/d})$ Shen et al. (2020); Lu et al. (2020), respectively. For results in terms of the number of nonzero parameters, the reader is referred to Yarotsky (2017); Schmidt-Hieber (2020); Petersen and Voigtlaender (2018); Yarotsky (2018); Gühring et al. (2019); Yarotsky and Zhevnerchuk (2020) and the reference therein. Obviously, the curse of dimensionality exists in ReLU networks for these generic functions and, therefore, ReLU networks would need to be exponentially large in d to maintain a reasonably good approximation accuracy. The curse could be lessened when target function spaces are smaller. To name a few, Poggio et al. (2017); Barron and Klusowski (2018); E et al. (2019); Montanelli et al. (2020); Chen et al. (2019a); Hutzenthaler et al. (2020) and the reference therein for ReLU networks. The limitation of ReLU networks motivated the work in Shen et al. (2021) to introduce Floor-ReLU networks built with either a Floor ($\lfloor x \rfloor$) or ReLU ($\max\{0, x\}$) activation function in each neuron. It was shown

by construction in Shen et al. (2021) that Floor-ReLU networks with width $\max\{d, 5N + 13\}$ and depth $64dL + 3$ can uniformly approximate a Hölder continuous function f on $[0, 1]^d$ with a root-exponential approximation rate $3\lambda d^{\alpha/2} N^{-\alpha\sqrt{L}}$ without the curse of dimensionality, where $\alpha \in (0, 1]$ and $\lambda > 0$ are the Hölder order and constant, respectively.

The most important message of Shen et al. (2021) (and probably also of Yarotsky and Zhevnerchuk (2020)) is that the combination of simple activation functions can create super approximation power. In the Floor-ReLU networks mentioned above, the power of depth is fully reflected in the approximation rate $3\lambda d^{\alpha/2} N^{-\alpha\sqrt{L}}$ that is root-exponential in depth. However, the power of width is much weaker and the approximation rate is polynomial in width if depth is fixed. This seems to be inconsistent with recent development of network optimization theory Jacot et al. (2018); Du et al. (2019); Mei et al. (2018); Wu et al. (2018); Chen et al. (2019b); Lu et al. (2020); Luo and Yang (2020), where larger width instead of depth can ease the challenge of highly nonconvex optimization. The mystery of the power of width and depth remains and it motivates us to demonstrate that width can also enable super approximation power when armed with appropriate activation functions.

In particular, we explore the floor function, the exponential function (2^x), the step function ($\mathbf{1}_{x \geq 0}$), or their compositions as activation functions to build fully-connected feed-forward neural networks. These networks are called Floor-Exponential-Step (FLES) networks. As we shall prove by construction, Theorem 1.1 below shows that FLES networks with width $\max\{d, N\}$ and three hidden layers can uniformly approximate a continuous function f on $[0, 1]^d$ with an exponential approximation rate $2\omega_f(2\sqrt{d})2^{-N} + \omega_f(2\sqrt{d}2^{-N})$, where $\omega_f(\cdot)$ is the modulus of continuity defined as

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

In particular, there are three kinds of activation functions denoted as σ_1 , σ_2 , and σ_3 in FLES networks (see Figure 1 for an illustration):

$$\sigma_1(x) := \lfloor x \rfloor, \quad \sigma_2(x) := 2^x, \quad \text{and} \quad \sigma_3(x) := \mathcal{T}(x - \lfloor x \rfloor - \frac{1}{2}), \quad \text{for any } x \in \mathbb{R},$$

where

$$\mathcal{T}(x) := \mathbf{1}_{x \geq 0} = \begin{cases} 1, & x \geq 0, \\ 0, & x < 0, \end{cases} \quad \text{for any } x \in \mathbb{R}.$$

Theorem 1.1. *Given an arbitrary continuous function f defined on $[0, 1]^d$, for any $N \in \mathbb{N}^+$, there exist $a_1, a_2, \dots, a_N \in [0, \frac{1}{2})$ such that*

$$|\phi(\mathbf{x}) - f(\mathbf{x})| \leq 2\omega_f(2\sqrt{d})2^{-N} + \omega_f(2\sqrt{d})2^{-N},$$

for any $\mathbf{x} = (x_1, x_2, \dots, x_d) \in [0, 1]^d$, where ϕ is defined by a formula in a_1, a_2, \dots, a_N as follows.

$$\phi(\mathbf{x}) = 2\omega_f(2\sqrt{d}) \sum_{j=1}^N 2^{-j} \sigma_3 \left(a_j \cdot \sigma_2 \left(1 + \sum_{i=1}^d 2^{(i-1)N} \sigma_1(2^{N-1} x_i) \right) \right) + f(\mathbf{0}) - \omega_f(2\sqrt{d}).$$

We remark that ϕ in Theorem 1.1 is essentially determined by N parameters a_1, a_2, \dots, a_N , which can be trained by a $(\sigma_1, \sigma_2, \sigma_3)$ -activated network with width $\max\{d, N\}$, three hidden layers, and $2(d + N + 1)$ nonzero parameters. See Figure 1 for an illustration.

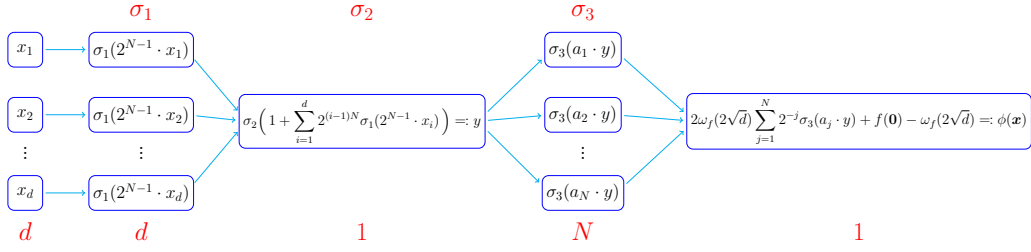


Figure 1: An illustration of the desired three-hidden-layer network in Theorem 1.1 for any $\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}$. Each of the red functions “ σ_1 ”, “ σ_2 ”, and “ σ_3 ” above the network is the activation function of the corresponding hidden layer. The number of neurons in each hidden layer is indicated by the red number below it.

The rate in $\omega_f(2\sqrt{d})2^{-N}$ implicitly depends on N through the modulus of continuity of f , while the rate in $2\omega_f(2\sqrt{d})2^{-N}$ is explicit in N . Simplifying the implicit approximation rate to make it explicitly depend on N is challenging in general. However, if f is a Hölder continuous function on $[0, 1]^d$ of order $\alpha \in (0, 1]$ with a Hölder constant $\lambda > 0$, i.e., $f(\mathbf{x})$ satisfying

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq \lambda \|\mathbf{x} - \mathbf{y}\|_2^\alpha, \quad \text{for any } \mathbf{x}, \mathbf{y} \in [0, 1]^d, \quad (1)$$

then $\omega_f(r) \leq \lambda r^\alpha$ for any $r \geq 0$. Therefore, in the case of Hölder continuous functions, the approximation rate is simplified to $3\lambda(2\sqrt{d})^\alpha 2^{-\alpha N}$ as shown in the following corollary. In the special case of Lipschitz continuous functions with a Lipschitz constant $\lambda > 0$, the approximation rate is simplified to $6\lambda\sqrt{d}2^{-N}$.

Corollary 1.2. *Given any Hölder continuous function f on $[0, 1]^d$ of order $\alpha \in (0, 1]$ with a Hölder constant $\lambda > 0$, for any $N \in \mathbb{N}^+$, there exists a_1, a_2, \dots, a_N such that*

$$|\phi(\mathbf{x}) - f(\mathbf{x})| \leq 3\lambda(2\sqrt{d})^\alpha 2^{-\alpha N}, \quad \text{for any } \mathbf{x} = (x_1, x_2, \dots, x_d) \in [0, 1]^d,$$

where ϕ is defined by a formula in a_1, a_2, \dots, a_N as follows.

$$\phi(\mathbf{x}) = 2\omega_f(2\sqrt{d}) \sum_{j=1}^N 2^{-j} \sigma_3 \left(a_j \cdot \sigma_2 \left(1 + \sum_{i=1}^d 2^{(i-1)N} \sigma_1(2^{N-1} x_i) \right) \right) + f(\mathbf{0}) - \omega_f(2\sqrt{d}).$$

First, Theorem 1.1 and Corollary 1.2 show that the approximation capacity of three-hidden-layer neural networks with simple activation functions for continuous functions can be exponentially improved by increasing the network width, and the approximation error can be explicitly characterized in terms of the width $\mathcal{O}(N)$. Second, this new class of networks overcomes the curse of dimensionality in the approximation power when the modulus of continuity is moderate, since the approximation order is essentially \sqrt{d} times a function of N independent of d within the modulus of continuity. Therefore, three hidden layers are enough for neural networks to achieve exponential convergence and avoid the curse of dimensionality for generic functions. The width is also powerful in network approximation.

The rest of this paper is organized as follows. In Section 2, we discuss the application scope of our theory, study the connection between the approximation error and the Vapnik-Chervonenkis (VC) dimension, establish Corollary 2.3 to extend our analysis to general bounded continuous functions on a bounded set, and compare related works in the literature. We will prove Theorem 1.1 and Corollary 2.3 in Section 3. In Section 4, we explore alternative continuous activation functions other than σ_1 , σ_2 , and σ_3 for super approximation power. Finally, we conclude this paper in Section 5.

2. Discussion

In this section, we will further interpret our results and discuss related research in the field of neural network approximation.

2.1. Application scope of our theory in machine learning

Let $\phi(\mathbf{x}; \boldsymbol{\theta})$ denote a function computed by a (fully-connected) network with $\boldsymbol{\theta}$ as the set of parameters. Given a target function f , consider the

expected error/risk of $\phi(\mathbf{x}; \boldsymbol{\theta})$

$$R_{\mathcal{D}}(\boldsymbol{\theta}) := \mathbb{E}_{\mathbf{x} \sim U(\mathcal{X})} [\ell(\phi(\mathbf{x}; \boldsymbol{\theta}), f(\mathbf{x}))]$$

with a loss function typically taken as $\ell(y, y') = \frac{1}{2}|y - y'|^2$, where $U(\mathcal{X})$ is an unknown data distribution over \mathcal{X} . For example, when $\ell(y, y') = \frac{1}{2}|y - y'|^2$ and U is a uniform distribution over $\mathcal{X} = [0, 1]^d$,

$$R_{\mathcal{D}}(\boldsymbol{\theta}) = \int_{[0,1]^d} \frac{1}{2} |\phi(\mathbf{x}; \boldsymbol{\theta}) - f(\mathbf{x})|^2 d\mathbf{x}.$$

The expected risk minimizer $\boldsymbol{\theta}_{\mathcal{D}}$ is defined as

$$\boldsymbol{\theta}_{\mathcal{D}} := \arg \min_{\boldsymbol{\theta}} R_{\mathcal{D}}(\boldsymbol{\theta}).$$

It is unachievable in practice since f and $U(\mathcal{X})$ are not available. Instead, we only have samples of f .

Given samples $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n$, the empirical risk is defined as

$$R_{\mathcal{S}}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \ell(\phi(\mathbf{x}_i; \boldsymbol{\theta}), f(\mathbf{x}_i)).$$

And we usually use it to approximate/model the expected risk $R_{\mathcal{D}}(\boldsymbol{\theta})$. The goal of supervised learning is to identify the empirical risk minimizer

$$\boldsymbol{\theta}_{\mathcal{S}} = \arg \min_{\boldsymbol{\theta}} R_{\mathcal{S}}(\boldsymbol{\theta}), \tag{2}$$

to obtain $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{S}}) \approx f(\mathbf{x})$. When a numerical optimization method is applied to solve (2), it may result in a numerical solution (denoted as $\boldsymbol{\theta}_{\mathcal{N}}$) that is not a global minimizer. Hence, the actually learned function generated by a neural network is $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}})$. The discrepancy between the target function f and the actually learned function $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}})$ is measured by an inference error

$$R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) = \mathbb{E}_{\mathbf{x} \sim U(\mathcal{X})} [\ell(\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}}), f(\mathbf{x}))] \stackrel{e.g.}{=} \int_{[0,1]^d} \frac{1}{2} |\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}}) - f(\mathbf{x})|^2 d\mathbf{x},$$

where the second equality holds when $\ell(y, y') = \frac{1}{2}|y - y'|^2$ and U is a uniform distribution over $\mathcal{X} = [0, 1]^d$.

Since $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}})$ is the expected inference error over all possible data samples, it can quantify how good the learned function $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}})$ is. Note that

$$\begin{aligned}
 R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) &= \underbrace{[R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}})]}_{\text{GE}} + \underbrace{[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}})]}_{\text{OE}} + \underbrace{[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}})]}_{\leq 0 \text{ by Eq. (2)}} + \underbrace{[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}}) - R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})]}_{\text{GE}} + \underbrace{R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})}_{\text{AE}} \\
 &\leq \underbrace{R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})}_{\text{Approximation error (AE)}} + \underbrace{[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}})]}_{\text{Optimization error (OE)}} + \underbrace{[R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}})] + [R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}}) - R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})]}_{\text{Generalization error (GE)}}, \quad (3)
 \end{aligned}$$

where the inequality comes from the fact that $[R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{D}})] \leq 0$ since $\boldsymbol{\theta}_{\mathcal{S}}$ is a global minimizer of $R_{\mathcal{S}}(\boldsymbol{\theta})$. Constructive approximation provides an upper bound of $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$ in terms of the network size, e.g., in terms of the network width and depth, or in terms of the number of parameters. The second term of Equation (3) is bounded by the optimization error of the numerical algorithm applied to solve the empirical loss minimization problem in Equation (2). Note that one only needs to make $R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{N}}) - R_{\mathcal{S}}(\boldsymbol{\theta}_{\mathcal{S}})$ small, but not $\boldsymbol{\theta}_{\mathcal{N}} - \boldsymbol{\theta}_{\mathcal{S}}$. The study of the bounds for the third and fourth terms is referred to as the generalization error analysis of neural networks. See Figure 2 for the intuitions of these three errors.

One of the key targets in the area of deep learning is to develop algorithms to reduce $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}})$. The constructive approximation established in this paper and in the literature provides upper bounds of the approximation error $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$ for several function spaces, which is crucial to estimate an upper bound of $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}})$. Instead of deriving an approximator to attain the approximation error bound, deep learning algorithms aim to identify a solution $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}})$ reducing the generalization and optimization errors in Equation (3). Solutions minimizing both generalization and optimization errors will lead to a good solution only if we also have a good upper bound estimate of $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$ as shown in Equation (3). Independent of whether our analysis here leads to a good approximator, which is an interesting topic to pursue, the theory here does provide a key ingredient in the error analysis of deep learning algorithms.

Theorem 1.1 and Corollary 1.2 provide an upper bound of $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$. This bound only depends on the given budget of neurons and layers of FLES networks. Hence, this bound is independent of the empirical loss minimization in Equation (2) and the optimization algorithm used to compute the numerical solution of Equation (2). In other words, Theorem 1.1 and Corollary 1.2 quantify the approximation power of FLES networks with a given size. Designing efficient optimization algorithms and analyzing the generalization

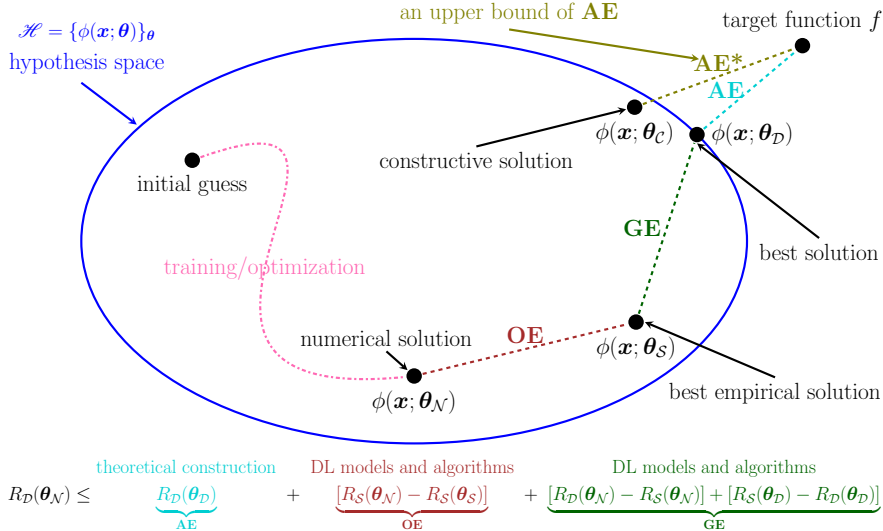


Figure 2: The intuitions of the approximation error (AE), the optimization error (OE), and the generalization error (GE). DL is short of deep learning. One needs to control AE, OE, and GE in order to bound the discrepancy between the target function f and the numerical solution $\phi(\mathbf{x}; \boldsymbol{\theta}_N)$ (what we can get in practice), measured by $R_{\mathcal{D}}(\boldsymbol{\theta}_N) = \mathbb{E}_{\mathbf{x} \sim U(\mathcal{X})} [\ell(\phi(\mathbf{x}; \boldsymbol{\theta}_N), f(\mathbf{x}))]$.

bounds for FLES networks are two other separate future directions.

2.2. Connection between approximation error and VC-dimension

The approximation error and the Vapnik-Chervonenkis (VC) dimension are two important measures of the capacity (complexity) of a set of functions. In this section, we discuss the connection between them.

Let us first present the definitions of VC-dimension and related concepts. Assume H is a class of functions mapping from a general domain \mathcal{X} to $\{0, 1\}$. We say H shatters a set of points $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\} \subseteq \mathcal{X}$ 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 $|\cdot|$ means the size of a set. The above 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 all i . For a general function set \mathcal{F} with its elements mapping from \mathcal{X} to \mathbb{R} , we say \mathcal{F} shatters $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\} \subseteq \mathcal{X}$ if $\mathcal{T} \circ \mathcal{F}$ does, where

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

For any $m \in \mathbb{N}^+$, the growth function of H is defined 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 2.1 (VC-dimension). Assume H is 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,

$$\text{VCDim}(H) := \sup \{m \in \mathbb{N}^+ : \Pi_H(m) = 2^m\}$$

in the case that $\{m \in \mathbb{N}^+ : \Pi_H(m) = 2^m\}$ is not empty. If $\{m \in \mathbb{N}^+ : \Pi_H(m) = 2^m\} = \emptyset$, we may define $\text{VCDim}(H) = 0$.

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 $\text{VCDim}(\mathcal{F}) := \text{VCDim}(\mathcal{T} \circ \mathcal{F})$,¹ 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 “the VC-dimension of a network (architecture)” means the VC-dimension of the function set that consists of all functions computed by this network (architecture).

As shown in [Yarotsky \(2018, 2017\)](#); [Shen et al. \(2019, 2020\)](#); [Lu et al. \(2020\)](#); [Shen et al. \(2021\)](#); [Zhang \(2020\)](#); [Shen et al. \(to appear\)](#), VC-dimension essentially determines the lower bound of the approximation errors of networks. For simplicity, we use $\text{H\"older}([0, 1]^d, \alpha, \lambda)$ as an example, where $\text{H\"older}([0, 1]^d, \alpha, \lambda)$ denotes the space of H\"older continuous functions of order $\alpha \in (0, 1]$ and a H\"older constant $\lambda > 0$. Without loss of generality, we assume $\lambda = 1$. [Theorem 2.2](#) below shows that the best possible approximation error of functions in $\text{H\"older}([0, 1]^d, \alpha, 1)$ approximated by functions in \mathcal{F} is bounded by a formula characterized by $\text{VCDim}(\mathcal{F})$.

Theorem 2.2 (Theorem 2.4 of [Shen et al. \(to appear\)](#) or Theorem 4.17 of [Zhang \(2020\)](#)). *Assume \mathcal{F} is a function set with all elements defined on $[0, 1]^d$. Given any $\varepsilon > 0$, suppose $\text{VCDim}(\mathcal{F}) \geq 1$ and*

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

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

¹One may also define $\text{VCDim}(\mathcal{F}) := \text{VCDim}(\widehat{\mathcal{T}} \circ \mathcal{F})$, where $\widehat{\mathcal{T}}(t) := \begin{cases} 1, & t > 0, \\ 0, & t \leq 0. \end{cases}$

This theorem investigates the connection between VC-dimension of \mathcal{F} and the approximation errors of functions in Hölder($[0, 1]^d, \alpha, 1$) approximated by elements of \mathcal{F} . Denote the best approximation error of functions in Hölder($[0, 1]^d, \alpha, 1$) approximated by the elements of \mathcal{F} as

$$\mathcal{E}_{\alpha,d}(\mathcal{F}) := \sup_{f \in \text{Hölder}([0,1]^d, \alpha, 1)} \left(\inf_{\phi \in \mathcal{F}} \|\phi - f\|_{L^\infty([0,1]^d)} \right).$$

Then, Theorem 2.2 implies that

$$\text{VCDim}(\mathcal{F})^{-\alpha/d}/9 \leq \mathcal{E}_{\alpha,d}(\mathcal{F}),$$

which means that the best possible approximation error is controlled by $\text{VCDim}(\mathcal{F})^{-\alpha/d}/9$. A typical application of this theorem is to prove the optimality of approximation errors when using ReLU networks to approximate functions in Hölder($[0, 1]^d, \alpha, 1$). It is shown in Harvey et al. (2017) that the VC-dimension of $\mathcal{F}_{N,L}$ is bounded by

$$\text{VCDim}(\mathcal{F}_{N,L}) \leq \mathcal{O}(N^2 L \cdot L \cdot \ln(N^2 L)) \leq \mathcal{O}(N^2 L^2 \ln(NL)),$$

where $\mathcal{F}_{N,L}$ is the space consisting of all functions implemented by ReLU networks with width N and depth L . It is shown in Section 4.4.1 of Zhang (2020) that

$$C_1(\alpha, d) \cdot \left(N^2 L^2 \ln(NL) \right)^{-\alpha/d} \leq \mathcal{E}_{\alpha,d}(\mathcal{F}_{N,L}) \leq C_2(\alpha, d) \cdot \left(N^2 L^2 \right)^{-\alpha/d},$$

where $C_1(\alpha, d)$ and $C_2(\alpha, d)$ are two positive constants determined by α and d .

Finally, we would like to point out that a large VC-dimension of the hypothesis space \mathcal{F} is a **necessary** condition of a good approximation error, but cannot guarantee a good approximation error, which also relies on other properties of the hypothesis space \mathcal{F} . For example, it is easy to check by Proposition 4.2 that

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

However, $\{\phi : \phi(x) = \cos(ax), a \in \mathbb{R}\}$ cannot achieve a good approximation error when approximating Hölder continuous functions. Designing a hypothesis space with a large VC-dimension is the first step for a good approximation

toll, but to realize the desired approximation power requires refined design of the hypothesis space, which is also the philosophy we followed in this paper. Our initial goal is to design a network architecture with a fixed depth (e.g., three hidden layers) to generate a hypothesis space with a sufficiently large VC-dimension (∞). As we shall see later, Proposition 3.2 implies that the VC-dimension of FLES networks is infinity, which is a necessary condition for our FLES networks to attain super approximation power.

2.3. Further interpretation of our theory

In the interpretation of our theory, three more aspects are important to discuss. The first one is whether it is possible to extend our theory to functions on a more general domain, e.g. $E \subseteq [-R, R]^d$ for any $R > 0$, because $R > 1$ may cause an implicit curse of dimensionality in some existing theory. The second one is how bad the modulus of continuity would be since it is related to a high-dimensional function f that may lead to an implicit curse of dimensionality in our approximation rate. The last one is the discussion of overcoming the zero derivative in training FLES networks.

First, we can generalize Theorem 1.1 to the function space $C(E)$ with $E \subseteq [-R, R]^d$ for any $R > 0$ in the following corollary with the modulus of continuity $\omega_f^E(\cdot)$ defined as follows. For an arbitrary set $E \subseteq \mathbb{R}^d$, $\omega_f^E(r)$ is defined via

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

As defined earlier, in the case $E = [0, 1]^d$, $\omega_f^E(r)$ is abbreviated to $\omega_f(r)$. The proof of this corollary will be presented in Section 3.2.

Corollary 2.3. *Given an arbitrary bounded continuous function f on $E \subseteq [-R, R]^d$ where R is an arbitrary positive real number, for any $N \in \mathbb{N}^+$, there exist $a_1, a_2, \dots, a_N \in [0, \frac{1}{2})$ such that*

$$|\phi(\mathbf{x}) - f(\mathbf{x})| \leq 2\omega_f^E(3R\sqrt{d})2^{-N} + \omega_f^E(3R\sqrt{d}2^{-N}),$$

for any $\mathbf{x} = (x_1, x_2, \dots, x_d) \in E$, where ϕ is defined by a formula in a_1, a_2, \dots, a_N as follows.

$$\phi(\mathbf{x}) = 2\omega_f^E(3R\sqrt{d}) \sum_{j=1}^N 2^{-j} \sigma_3 \left(a_j \cdot \sigma_2 \left(1 + \sum_{i=1}^d 2^{(i-1)N} \sigma_1 \left(2^N \frac{x_i + R}{3R} \right) \right) \right) + C_f,$$

where C_f is a constant determined by f .

Hence, the volume of the function domain $E \subseteq [-R, R]^d$ only has a mild influence on the approximation rate of our FLES networks. FLES networks can still avoid the curse of dimensionality and achieve exponential convergence for continuous functions on $E \subseteq [-R, R]^d$ when $R > 1$. For example, in the case of Hölder continuous functions of order $\alpha \in (0, 1]$ with a constant $\lambda > 0$ on $E \subseteq [-R, R]^d$, our approximation rate becomes $3\lambda(3R\sqrt{d}2^{-N})^\alpha$.

Second, most interesting continuous functions in practice have a good modulus of continuity such that there is no implicit curse of dimensionality hidden in $\omega_f(\cdot)$. For example, we have discussed the case of Hölder continuous functions previously. We would like to remark that the class of Hölder continuous functions implicitly depends on d through its definition in Equation (1), but this dependence is moderate since the ℓ^2 -norm in Equation (1) is the square root of a sum with d terms. Let us now discuss several cases of $\omega_f(\cdot)$ when we cannot achieve exponential convergence or cannot avoid the curse of dimensionality. The first example is $\omega_f(r) = \frac{1}{\ln(1/r)}$ for small $r > 0$, which leads to an approximation rate

$$3(N \ln 2 - \frac{1}{2} \ln d - \ln 2)^{-1}, \quad \text{for large } N \in \mathbb{N}^+.$$

Apparently, the above approximation rate still avoids the curse of dimensionality but there is no exponential convergence, which has been canceled out by “ln” in $\omega_f(\cdot)$. The second example is $\omega_f(r) = \frac{1}{\ln^{1/d}(1/r)}$ for small $r > 0$, which leads to an approximation rate

$$3(N \ln 2 - \frac{1}{2} \ln d - \ln 2)^{-1/d}, \quad \text{for large } N \in \mathbb{N}^+.$$

The power $1/d$ further weakens the approximation rate and hence the curse of dimensionality exists. The last example we would like to discuss is $\omega_f(r) = r^{\alpha/d}$ for small $r > 0$, which results in the approximation rate

$$3\lambda(2\sqrt{d})^{\alpha/d}2^{-\alpha N/d}, \quad \text{for large } N \in \mathbb{N}^+,$$

which achieves the exponential convergence and avoids the curse of dimensionality when we use very wide networks. Though we have provided several examples of immoderate $\omega_f(\cdot)$, to the best of our knowledge, we are not aware of useful continuous functions with $\omega_f(\cdot)$ that is immoderate.

Finally, we would like to point out that the training of FLES networks in practice may encounter two issues. First, network weights in our main theorems require high-precision computation that might not be available in

existing computer systems when the dimension d and the network size parameter N are large. But there is no theoretical evidence to exclude the possibility that similar approximation results can be achieved with reasonable weights in practical computation. Second, the vanishing gradient of piecewise constant activation functions makes standard SGD infeasible. There are two possible directions to solve the optimization problem for FLES networks: 1) gradient-free optimization methods, e.g., Nelder-Mead method [Nelder and Mead \(1965\)](#), genetic algorithm [Holland \(1992\)](#), simulated annealing [Kirkpatrick et al. \(1983\)](#), particle swarm optimization [Kennedy and Eberhart \(1995\)](#), and consensus-based optimization [Pinnau et al. \(2017\)](#); [Carrillo et al. \(2019\)](#); 2) applying optimization algorithms for quantized networks that also have piecewise constant activation functions [Lin et al. \(2019\)](#); [Boo et al. \(2020\)](#); [Bengio et al. \(2013\)](#); [Wang et al. \(2018\)](#); [Hubara et al. \(2017\)](#); [Yin et al. \(2019\)](#). For example, an empirical way is to use a straight-through estimator (STE) by setting the incoming gradients to the activation function (e.g., Floor) equal to its outgoing gradients, disregarding the derivative of the activation function itself. It would be interesting future work to explore efficient learning algorithms based on the FLES network.

2.4. Kolmogorov-Arnold Superposition Theorem

A closely related research topic is the Kolmogorov-Arnold representation theorem (KST) [Kolmogorov \(1956\)](#); [Arnold \(1957\)](#); [Kolmogorov \(1957\)](#) and its approximation in a form of modern neural networks. Our FLES networks admit super approximation power with a fixed number of layers for continuous functions and the KST exactly represent continuous functions using two hidden layers and $\mathcal{O}(d)$ neurons. More specifically, given any $f \in C([0, 1]^d)$, the KST shows that there exist continuous functions $\phi_q : \mathbb{R} \rightarrow \mathbb{R}$ and $\psi_{q,p} : [0, 1] \rightarrow \mathbb{R}$ such that

$$f(\mathbf{x}) = \sum_{q=0}^{2d} \phi_q \left(\sum_{p=1}^d \psi_{q,p}(x_p) \right), \quad \text{for any } \mathbf{x} = (x_1, \dots, x_d) \in [0, 1]^d. \quad (4)$$

Note that the activation functions $\{\phi_q\}$ (also called outer functions) of the neural network in Equation (4) have to depend on the target function f , though $\{\psi_{q,p}\}$ (also called inner functions) can be independent of f . The modulus of continuity of $\{\psi_{q,p}\}$ can be constructed such that they moderately depend on d , but the modulus of continuity of $\{\phi_q\}$ would be exponentially bad in d . In sum, the outer functions are too pathological such that there is

no existing numerical algorithms to evaluate these activation functions, even though they are shown to exist by iterative construction [Braun and Griebel \(2009\)](#).

There has been an active research line to develop more practical network approximation based on KST [Kůrková \(1991, 1992\)](#); [Maiorov and Pinkus \(1999\)](#); [Guliyev and Ismailov \(2018\)](#); [Montanelli and Yang \(2020\)](#); [Igel'nik and Parikh \(2003\)](#); [Schmidt-Hieber \(2021\)](#) by relaxing the exact representation to network approximation with an ε -error. The key issue these KST-related networks attempting to address is the f -dependency of the activation functions and the main goal is to construct neural networks conquering the curse of dimensionality in a more practical way computationally. The main idea of these variants is to apply computable activation functions independent of f to construct neural networks to approximate the outer and inner functions of the KST, resulting in a larger network that can approximate a continuous function with the desired accuracy. Using this idea, the seminal work in [Kůrková \(1992\)](#) applied sigmoid activation functions and constructed two-hidden-layer networks to approximate $f \in C([0, 1]^d)$. Though the activation functions are independent of f , the number of neurons scales exponentially in d and the curse of dimensionality exists. Cubic-splines and piecewise linear functions have also been used to approximate the outer and inner functions of KST in [Igel'nik and Parikh \(2003\)](#); [Montanelli and Yang \(2020\)](#); [Schmidt-Hieber \(2021\)](#), resulting in cubic-spline networks or deep ReLU networks to approximate $f \in C([0, 1]^d)$. But the approximation bounds in these works still suffer from the curse of dimensionality unless f has simple outer functions in the KST. It is still an open problem to characterize the class of functions with a moderate outer function in KST.

To the best of our knowledge, the most successful construction of neural networks with f -independent activation functions conquering the curse of dimensionality is in [Maiorov and Pinkus \(1999\)](#); [Guliyev and Ismailov \(2018\)](#), where a two-hidden-layer network with $\mathcal{O}(d)$ neurons can approximate $f \in C([0, 1]^d)$ within an arbitrary error ε . Let us briefly summarize their main ideas to obtain such an exciting result here. 1) Identify a dense and countable subset $\{u_k\}_{k=1}^\infty$ of $C([-1, 1])$, e.g., polynomials with rational coefficients. 2) Construct an activation function ϱ to “store” all $u_k(x)$ for $x \in [-1, 1]$. For example, divide the domain of $\varrho(x)$ into countable pieces and each piece is a connected interval of length 2 associated with a u_k . In particular, let $\varrho(x + 4k + 1) = a_k + b_k x + c_k u_k(x)$ for any $x \in [-1, 1]$ with carefully chosen constants a_k , b_k , and c_k such that $\varrho(x)$ can be a sigmoid function. 3) By

construction, there exists a one-hidden-layer network with width 3 and $\varrho(x)$ as the activation function to approximate any outer or inner function in KST with an arbitrary accuracy parameter δ . Only the parameters of the one-hidden-layer network depend on the target function and accuracy. 4) Replace the inner and outer function in KST with these one-hidden-layer networks to achieve a two-hidden-layer network with $\varrho(x)$ as the activation function and width $\mathcal{O}(d)$ to approximate an arbitrary $f \in C([0, 1]^d)$ within an arbitrary error ε . Unfortunately, the construction of the parameters of this magic network relies on the evaluation of the outer and inner functions of KST, which is not computationally feasible even if computation with arbitrary precision is allowed.

We would like to remark that, though the approximation rate of FLES networks in this paper is relatively worse than the approximation rate in [Maiorov and Pinkus \(1999\)](#); [Guliyev and Ismailov \(2018\)](#), our activation functions are much simpler and there are explicit formulas to specify the parameters of FLES networks. If computation with an arbitrary precision is allowed and the target function f can be arbitrarily sampled, we can specify all the weights in FLES networks. Besides, our approximation rate is sufficiently attractive since it is exponential and avoids the curse of dimensionality. For a large dimension d , the width parameter of our FLES network can be chosen as $N = d$, which leads to a FLES network of size $\mathcal{O}(d)$ with an approximation accuracy $\mathcal{O}(2^{-d})$ for Lipschitz continuous functions. $\mathcal{O}(2^{-d})$ is sufficiently attractive. In practice, when d is very large, N could be much smaller than d and our approximation rate is still attractive.

Finally, we list several KST-related results in [Table 1](#) for a quick comparison.² As shown in [Table 1](#), there exists a trade-off between the complexity of activation functions and the network size when the approximation error is fixed. A key advantage of our FLES networks is to use simple and explicit activation functions to attain an exponential convergence rate.

2.5. Discussion on the literature

In this section, we will discuss other recent development of neural network approximation. Our discussion will be divided into mainly three parts according to the analysis methodology in the references: 1) functions admitting integral representations; 2) linear approximation; 3) bit extraction.

²The result in [Shen et al. \(2019\)](#) is for Hölder functions, but can be easily generalized to general continuous functions.

Table 1: A comparison of several KST-related results for approximating $f \in C([0, 1]^d)$.

paper	number of hidden layers	width	activation function(s)	error	remark
Kolmogorov (1956); Arnold (1957); Kolmogorov (1957)	2	$2d + 1$	f -dependent	0	original KST
Maiorov and Pinkus (1999); Guliyev and Ismailov (2018)	2	$\mathcal{O}(d)$	f -independent	arbitrary error ε	based on KST
Shen et al. (2019)	3	$\mathcal{O}(dN)$	ReLU	$\mathcal{O}(\omega_f(N^{-2/d}))$	not based on KST
this paper	3	$\max\{d, N\}$	$(\sigma_1, \sigma_2, \sigma_3)$	$2\omega_f(\sqrt{d})2^{-N} + \omega_f(\sqrt{d}2^{-N})$	not based on KST

In the seminal work of Barron (1993), its variants or generalization Barron and Klusowski (2018); E et al. (2019); Chen and Wu (2019); Montanelli et al. (2020), and related references therein, d -dimensional functions of the following form were considered:

$$f(\mathbf{x}) = \int_{\tilde{\Omega}} a(\mathbf{w})K(\mathbf{w} \cdot \mathbf{x})d\mu(\mathbf{w}),$$

where $\tilde{\Omega} \subseteq \mathbb{R}^d$, $\mu(\mathbf{w})$ is a Lebesgue measure in \mathbf{w} , and $\mathbf{x} \in \Omega \subseteq \mathbb{R}^d$. The above integral representation is equivalent to the expectation of a high-dimensional random function when \mathbf{w} is treated as a random variable. By the law of large number theory, the average of N samples of the integrand leads to an approximation of $f(\mathbf{x})$ with an approximation error bounded by $\frac{C_f \sqrt{\mu(\Omega)}}{\sqrt{N}}$ measured in $L^2(\Omega, \mu)$ (Equation (6) of Barron (1993)), where $\mathcal{O}(N)$ is the total number of parameters in the network, C_f is a d -dimensional integral with an integrand related to f , and $\mu(\Omega)$ is the Lebesgue measure of Ω . As discussed in Barron (1993), $\mu(\Omega)$ and C_f would be exponential in d and standard smoothness properties of f alone are not enough to remove the exponential dependence of C_f on d . Therefore, the curse of dimensionality exists in the whole approximation error while the curse does not exist in the approximation rate in N .

Linear approximation is an efficient approximation tool for smooth functions that computes the approximant of a target function via a linear projection to a Hilbert space or a Banach space as the approximant space. Typical examples include approximation via orthogonal polynomials, Fourier series expansion, etc. Inspired by the seminal work in Yarotsky (2017), where deep ReLU networks were constructed to approximate polynomials with exponential convergence, subsequent works in E and Wang (2018); Opschoor et al. (2019); Montanelli and Du (2019); Chen and Wu (2019); Montanelli et al. (2020); Yarotsky and Zhevnerchuk (2020); Lu et al. (2020); Montanelli and Yang (2020); Yang and Wang (2020) have constructed deep ReLU networks to approximate various smooth function spaces. The main idea of these works is to approximate smooth functions via (piecewise) polynomial

approximation first and then construct deep ReLU networks to approximate the ensemble of polynomials. But the curse of dimensionality exists since polynomial approximation cannot avoid the curse. Finally, a different approach is used in [Li et al. \(to appear\)](#). The authors of [Li et al. \(to appear\)](#) use a dynamic system based approach to obtain a universal approximation property of residual neural networks.

The bit extraction proposed in [Bartlett et al. \(1998\)](#) has been a very important technique to develop nearly optimal approximation rates of deep ReLU neural networks [Yarotsky \(2018\)](#); [Shen et al. \(2020\)](#); [Lu et al. \(2020\)](#); [Yang and Wang \(2020\)](#); [Zhang \(2020\)](#); [Shen et al. \(to appear\)](#) and the optimality is based on the nearly optimal VC-dimension bound of ReLU networks in [Harvey et al. \(2017\)](#). The bit extraction was also applied in [Shen et al. \(2021\)](#); [Schmidt-Hieber \(2021\)](#) and this paper to develop network approximation theories. In the first step, an efficient projection map in a form of a ReLU, or a Floor-ReLU, or a FLES network is constructed to project high-dimensional points to one-dimensional points such that the high-dimensional approximation problem is reduced to a one-dimensional approximation problem. In the second step, the one-dimensional approximation problem is solved by constructing a ReLU, or a Floor-ReLU, or a FLES network, which can be efficiently compressed via the bit extraction. Although shallower neural networks can also carry out the above two steps, bit extraction can take full advantage of the power of depth and construct deep neural networks with a nearly optimal number of parameters or neurons to fulfill the above two steps.

3. Theoretical Analysis

In this section, we first introduce basic notations in this paper in [Section 3.1](#). Then we prove [Theorem 1.1](#) and [Corollary 2.3](#) in [Section 3.2](#) based on [Theorem 3.1](#), which is proved in [Section 3.3](#).

3.1. Notations

The main notations of this paper are listed as follows.

- Vectors and matrices are denoted in a bold font. Standard vectorization is adopted in the matrix and vector computation. For example, adding a scalar and a vector means adding the scalar to each entry of the vector.

- Let \mathbb{R} denote the set of real numbers.
- Let \mathbb{Z} , \mathbb{N} , and \mathbb{N}^+ denote the set of integers, natural numbers, all positive integers, respectively, i.e., $\mathbb{Z} = \{0, 1, 2, \dots\} \cup \{-1, -2, -3, \dots\}$, $\mathbb{N} = \{0, 1, 2, \dots\}$, and $\mathbb{N}^+ = \{1, 2, 3, \dots\}$.
- For any $p \in [1, \infty)$, the p -norm of a vector $\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$ is defined by

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

- The floor function (Floor) is defined as $\lfloor x \rfloor := \max\{n : n \leq x, n \in \mathbb{Z}\}$ for any $x \in \mathbb{R}$.
- For $\theta \in [0, 1)$, 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., $\text{bin}0.\theta_1\theta_2\cdots\theta_L := \sum_{\ell=1}^L \theta_{\ell} 2^{-\ell}$.
- 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 .

3.2. Proof of Theorem 1.1 and Corollary 2.3

In this section, we will prove Theorem 1.1 and Corollary 2.3. To this end, we first introduce Theorem 3.1 that works only for $[0, 1)^d$, regarded as a weaker variant of Theorem 1.1.

Theorem 3.1. *Given an arbitrary continuous function f on $[0, 1]^d$, for any $N \in \mathbb{N}^+$, there exist $a_1, a_2, \dots, a_N \in [0, \frac{1}{2})$ such that*

$$|\phi(\mathbf{x}) - f(\mathbf{x})| \leq 2\omega_f(\sqrt{d})2^{-N} + \omega_f(\sqrt{d}2^{-N}),$$

for any $\mathbf{x} = (x_1, x_2, \dots, x_d) \in [0, 1)^d$, where ϕ is defined by a formula in a_1, a_2, \dots, a_N as follows.

$$\phi(\mathbf{x}) = 2\omega_f(\sqrt{d}) \sum_{j=1}^N 2^{-j} \sigma_3 \left(a_j \cdot \sigma_2 \left(1 + \sum_{i=1}^d 2^{(i-1)N} \sigma_1(2^N x_i) \right) \right) + f(\mathbf{0}) - \omega_f(\sqrt{d}).$$

We will prove Theorem 1.1 and Corollary 2.3 based on Theorem 3.1, the proof of which can be found in Section 3.3.

First, let us prove Theorem 1.1 by assuming Theorem 3.1 is true.

Proof of Theorem 1.1. Given any $f \in C([0, 1]^d)$, by Lemma 4.2 of Shen et al. (2020) via setting $E = [0, 1]^d$ and $S = \mathbb{R}^d$, there exists $g \in C(\mathbb{R}^d)$ such that

- $g(\mathbf{x}) = f(\mathbf{x})$ for any $\mathbf{x} \in E = [0, 1]^d$;
- $\omega_g^S(r) = \omega_f^E(r) = \omega_f(r)$ for any $r \geq 0$.

Define $\tilde{g}(\mathbf{x}) := g(2\mathbf{x})$ for any $\mathbf{x} \in \mathbb{R}^d$. By applying Theorem 3.1 to $\tilde{g} \in C(\mathbb{R}^d)$, there exist $a_1, a_2, \dots, a_N \in [0, \frac{1}{2})$ such that

$$|\tilde{\phi}(\mathbf{x}) - \tilde{g}(\mathbf{x})| \leq 2\omega_{\tilde{g}}^S(\sqrt{d})2^{-N} + \omega_{\tilde{g}}^S(\sqrt{d}2^{-N}), \quad \text{for any } \mathbf{x} \in [0, 1]^d, \quad (5)$$

where

$$\tilde{\phi}(\mathbf{x}) = 2\omega_{\tilde{g}}^S(\sqrt{d}) \sum_{j=1}^N 2^{-j} \sigma_3 \left(a_j \cdot \sigma_2 \left(1 + \sum_{i=1}^d 2^{(i-1)N} \sigma_1(2^N x_i) \right) \right) + \tilde{g}(\mathbf{0}) - \omega_{\tilde{g}}^S(\sqrt{d}).$$

Note that $f(\mathbf{x}) = g(\mathbf{x}) = \tilde{g}(\frac{\mathbf{x}}{2})$ for any $\mathbf{x} \in E = [0, 1]^d$ and

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

Define $\phi(\mathbf{x}) := \tilde{\phi}(2\mathbf{x})$ for any $\mathbf{x} \in \mathbb{R}^d$. Then by Equation (5), for any $\mathbf{x} \in [0, 1]^d = E$, we have $\frac{\mathbf{x}}{2} \in [0, \frac{1}{2}]^d \subseteq [0, 1]^d$, implying

$$\begin{aligned} |\phi(\mathbf{x}) - f(\mathbf{x})| &= |\phi(\mathbf{x}) - g(\mathbf{x})| = |\tilde{\phi}(\frac{\mathbf{x}}{2}) - \tilde{g}(\frac{\mathbf{x}}{2})| \\ &\leq 2\omega_{\tilde{g}}^S(\sqrt{d})2^{-N} + \omega_{\tilde{g}}^S(\sqrt{d}2^{-N}) \\ &= 2\omega_f(2\sqrt{d})2^{-N} + \omega_f(2\sqrt{d}2^{-N}), \end{aligned}$$

where $\phi(\mathbf{x}) := \tilde{\phi}(\frac{\mathbf{x}}{2})$ can be represented by

$$2\omega_f(2\sqrt{d}) \sum_{j=1}^N 2^{-j} \sigma_3 \left(a_j \cdot \sigma_2 \left(1 + \sum_{i=1}^d 2^{(i-1)N} \sigma_1(2^{N-1} x_i) \right) \right) + f(\mathbf{0}) - \omega_f(2\sqrt{d}).$$

With the discussion above, we have proved Theorem 1.1. \square

Next, we present the proof of Corollary 2.3 below.

Proof of Corollary 2.3. Given any bounded continuous function $f \in C(E)$, by Lemma 4.2 of Shen et al. (2020) 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

$$\tilde{g}(\mathbf{x}) := g(3R\mathbf{x} - R), \quad \text{for any } \mathbf{x} \in \mathbb{R}^d.$$

By applying Theorem 3.1 to $\tilde{g} \in C(\mathbb{R}^d)$, there exist $a_1, a_2, \dots, a_N \in [0, \frac{1}{2})$ such that

$$|\tilde{\phi}(\mathbf{x}) - \tilde{g}(\mathbf{x})| \leq 2\omega_{\tilde{g}}^S(\sqrt{d})2^{-N} + \omega_{\tilde{g}}^S(\sqrt{d}2^{-N}), \quad \text{for any } \mathbf{x} \in [0, 1]^d, \quad (6)$$

where

$$\tilde{\phi}(\mathbf{x}) = 2\omega_{\tilde{g}}^S(\sqrt{d}) \sum_{j=1}^N 2^{-j} \sigma_3 \left(a_j \cdot \sigma_2 \left(1 + \sum_{i=1}^d 2^{(i-1)N} \sigma_1(2^N x_i) \right) \right) + \tilde{g}(\mathbf{0}) - \omega_{\tilde{g}}^S(\sqrt{d}).$$

Note that $f(\mathbf{x}) = g(\mathbf{x}) = \tilde{g}(\frac{\mathbf{x}+R}{3R})$ for any $\mathbf{x} \in E \subseteq [-R, R]^d$ and

$$\omega_g^S(r) = \omega_f^E(3Rr) = \omega_f^E(3Rr), \quad \text{for any } r \geq 0.$$

Define $\phi(\mathbf{x}) := \tilde{\phi}(\frac{\mathbf{x}+R}{3R})$ for any $\mathbf{x} \in \mathbb{R}^d$. Then by Equation (6), for any $\mathbf{x} \in E \subseteq [-R, R]^d$, we have $\frac{\mathbf{x}+R}{3R} \in [0, \frac{2}{3}]^d \subseteq [0, 1]^d$, implying

$$\begin{aligned} |\phi(\mathbf{x}) - f(\mathbf{x})| &= |\phi(\mathbf{x}) - g(\mathbf{x})| = \left| \tilde{\phi}\left(\frac{\mathbf{x}+R}{3R}\right) - \tilde{g}\left(\frac{\mathbf{x}+R}{3R}\right) \right| \\ &\leq 2\omega_{\tilde{g}}^S(\sqrt{d})2^{-N} + \omega_{\tilde{g}}^S(\sqrt{d}2^{-N}) \\ &= 2\omega_f(3R\sqrt{d})2^{-N} + \omega_f(3R\sqrt{d}2^{-N}), \end{aligned}$$

where $\phi(\mathbf{x}) = \tilde{\phi}(\frac{\mathbf{x}+R}{3R})$ can be represented by

$$2\omega_f(3R\sqrt{d}) \sum_{j=1}^N 2^{-j} \sigma_3 \left(a_j \cdot \sigma_2 \left(1 + \sum_{i=1}^d 2^{(i-1)N} \sigma_1\left(2^N \frac{x_i+R}{3R}\right) \right) \right) + C_f,$$

where $C_f = \tilde{g}(\mathbf{0}) - \omega_{\tilde{g}}^S(\sqrt{d})$ is a constant essentially determined by f . With the discussion above, we have proved Corollary 2.3. \square

3.3. Proof of Theorem 3.1

To prove Theorem 3.1, we first present the proof sketch. Shortly speaking, we construct piecewise constant functions to approximate continuous functions. There are five key steps in our construction.

1. Normalize f as \tilde{f} satisfying $\tilde{f}(\mathbf{x}) \in [0, 1]$ for any $\mathbf{x} \in [0, 1]^d$, divide $[0, 1]^d$ into a set of non-overlapping cubes $\{Q_\beta\}_{\beta \in \{0, 1, \dots, J-1\}^d}$, and denote \mathbf{x}_β as the vertex of Q_β with minimum $\|\cdot\|_1$ norm, where J is an integer determined later. See Figure 3 for the illustrations of Q_β and \mathbf{x}_β .
2. Construct a vector-valued function $\Phi_1 : \mathbb{R}^d \rightarrow \mathbb{R}^d$ projecting the whole cube Q_β to the index β , i.e., $\Phi_1(\mathbf{x}) = \beta$ for all $\mathbf{x} \in Q_\beta$ and each $\beta \in \{0, 1, \dots, J-1\}^d$.
3. Construct a linear function $\phi_2 : \mathbb{R}^d \rightarrow \mathbb{R}$ bijectively mapping $\beta \in \{0, 1, \dots, J-1\}^d$ to $\phi_2(\beta) \in \{1, 2, \dots, J^d\}$.
4. Construct a function $\phi_3 : \mathbb{R} \rightarrow \mathbb{R}$ mapping $\phi_2(\beta) \in \{1, 2, \dots, J^d\}$ approximately to $\tilde{f}(\mathbf{x}_\beta)$, i.e., $\phi_3(\phi_2(\beta)) \approx \tilde{f}(\mathbf{x}_\beta)$ for each $\beta \in \{0, 1, \dots, J-1\}^d$.
5. Define $\tilde{\phi} := \phi_3 \circ \phi_2 \circ \Phi_1$. Then $\tilde{\phi}$ is a piecewise constant function mapping $\mathbf{x} \in Q_\beta$ to $\phi_3(\phi_2(\beta)) \approx \tilde{f}(\mathbf{x}_\beta)$ for each $\beta \in \{0, 1, \dots, J-1\}^d$, implying $\tilde{\phi} \approx \tilde{f}$. Finally, re-scale and shift $\tilde{\phi}$ to obtain the final function ϕ approximating f well.

Recall that

$$\sigma_1(x) := \lfloor x \rfloor, \quad \sigma_2(x) := 2^x, \quad \text{and} \quad \sigma_3(x) := \mathcal{T}(x - \lfloor x \rfloor - \frac{1}{2}), \quad \text{for any } x \in \mathbb{R},$$

where

$$\mathcal{T}(x) := \mathbf{1}_{x \geq 0} = \begin{cases} 1, & x \geq 0, \\ 0, & x < 0, \end{cases} \quad \text{for any } x \in \mathbb{R}.$$

Step 1 and 5 are straightforward. To implement Step 2, we introduce σ_1 since it can help to significantly simplify the construction of the vector-valued projecting function Φ_1 . The implementation of Step 3 is based on the J -ary representation, namely, define $\phi_2(\mathbf{x}) := 1 + \sum_{i=1}^d J^{i-1} x_i$. The most technical step above is Step 4, which is essentially a point fitting problem. Solving such a problem eventually relies on the bit extraction technique in Shen et al. (2020); Lu et al. (2020); Shen et al. (2021); Harvey et al. (2017); Bartlett et al. (1998); Zhang (2020); Yarotsky (2018). To extract sufficient many bits

with a limited neuron budget, we introduce two powerful activation functions σ_2 and σ_3 , as shown in the proposition below.

Proposition 3.2. *Given any $K \in \mathbb{N}^+$ and arbitrary $\theta_1, \theta_2, \dots, \theta_K \in \{0, 1\}$, it holds that*

$$\sigma_3(a \cdot \sigma_2(k)) = \sigma_3(2^k \cdot a) = \theta_k, \quad \text{for any } k \in \{1, 2, \dots, K\},$$

where

$$a = \sum_{j=1}^K 2^{-j-1} \cdot \theta_j \in [0, \frac{1}{2}).$$

Proof. Since $\theta_j \in \{0, 1\}$ for $j \in \{1, 2, \dots, K\}$, we have

$$0 \leq \sum_{j=1}^K 2^{-j-1} \cdot \theta_j \leq \sum_{j=1}^K 2^{-j-1} < \frac{1}{2},$$

implying $a \in [0, \frac{1}{2})$.

Next, fix $k \in \{1, 2, \dots, K\}$ for the proof below. It holds that

$$2^k \cdot a = 2^k \cdot \sum_{j=1}^K 2^{-j-1} \cdot \theta_j = \underbrace{\sum_{j=1}^{k-1} 2^{k-j-1} \cdot \theta_j}_{\text{an integer}} + \overbrace{\frac{1}{2}\theta_k}^{0 \text{ or } \frac{1}{2}} + \underbrace{\sum_{j=k+1}^K 2^{k-j-1} \cdot \theta_j}_{\text{in } [0, \frac{1}{2})}. \quad (7)$$

Clearly, the first term in Equation (7) $\sum_{j=1}^{k-1} 2^{k-j-1} \cdot \theta_j$ is a non-negative integer since $\theta_j \in \{0, 1\}$ for any $j \in \{1, 2, \dots, K\}$. As for the third term in Equation (7), we have

$$0 \leq \sum_{j=k+1}^K 2^{k-j-1} \cdot \theta_j \leq \sum_{j=k+1}^K 2^{k-j-1} < \frac{1}{2}$$

Therefore, by Equation (7), we have

$$2^k \cdot a \in \bigcup_{n \in \mathbb{N}} [n, n + \frac{1}{2}), \text{ if } \theta_k = 0 \quad \text{and} \quad 2^k \cdot a \in \bigcup_{n \in \mathbb{N}} [n + \frac{1}{2}, n + 1), \text{ if } \theta_k = 1. \quad (8)$$

Recall that $\sigma_3(x) = \mathcal{T}(x - \lfloor x \rfloor - \frac{1}{2})$, where $\mathcal{T}(x) = \begin{cases} 1, & x \geq 0, \\ 0, & x < 0. \end{cases}$ It is easy to verify that

$$\sigma_3(x) = 0 \text{ if } x \in \bigcup_{n \in \mathbb{N}} [n, n + \frac{1}{2}) \quad \text{and} \quad \sigma_3(x) = 1 \text{ if } x \in \bigcup_{n \in \mathbb{N}} [n + \frac{1}{2}, n + 1).$$

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

If $\theta_k = 0$, by Equation (8), we have

$$2^k \cdot a \in \bigcup_{n \in \mathbb{N}} [n, n + \frac{1}{2}) \implies \sigma_3(2^k \cdot a) = 0 = \theta_k.$$

Similarly, if $\theta_k = 1$, by Equation (8), we have

$$2^k \cdot a \in \bigcup_{n \in \mathbb{N}} [n + \frac{1}{2}, n + 1) \implies \sigma_3(2^k \cdot a) = 1 = \theta_k.$$

Since $k \in \{1, 2, \dots, K\}$ is arbitrary, we have $\sigma_3(a \cdot \sigma_2(k)) = \sigma_3(2^k \cdot a) = \theta_k$ for any $k \in \{1, 2, \dots, K\}$. So we finish the proof. \square

We would like to point out that Proposition 3.2 indicates that the VC-dimension of the function space

$$\{f : f(x) = \sigma_3(a \cdot x), \text{ for } a \in \mathbb{R}\}$$

is infinity, which implies that the VC-dimension of FLES networks is also infinity. As discussed previously in Section 2.2, having an infinite VC-dimension is a necessary condition for our FLES networks to attain super approximation power.

With Proposition 3.2 in hand, we are ready to prove Theorem 3.1.

Proof of Theorem 3.1. The proof consists of five steps.

Step 1: Set up.

Assume f is not a constant function since it is a trivial case. Then $\omega_f(r) > 0$ for any $r > 0$. Clearly, $|f(\mathbf{x}) - f(\mathbf{0})| \leq \omega_f(\sqrt{d})$ for any $\mathbf{x} \in [0, 1]^d$. Define

$$\tilde{f} := \frac{f - f(\mathbf{0}) + \omega_f(\sqrt{d})}{2\omega_f(\sqrt{d})}. \quad (9)$$

It follows that $\tilde{f}(\mathbf{x}) \in [0, 1]$ for any $\mathbf{x} \in [0, 1]^d$.

Set $J = 2^N$ and divide $[0, 1]^d$ into J^d cubes $\{Q_\beta\}_\beta$. To be exact, defined $\mathbf{x}_\beta := \beta/J$ and

$$Q_\beta := \left\{ \mathbf{x} = (x_1, x_2, \dots, x_d) : x_i \in \left[\frac{\beta_i}{J}, \frac{\beta_i+1}{J} \right) \text{ for } i = 1, 2, \dots, d \right\},$$

for each $\beta = (\beta_1, \beta_2, \dots, \beta_d) \in \{0, 1, \dots, J-1\}^d$. See Figure 3 for illustrations.

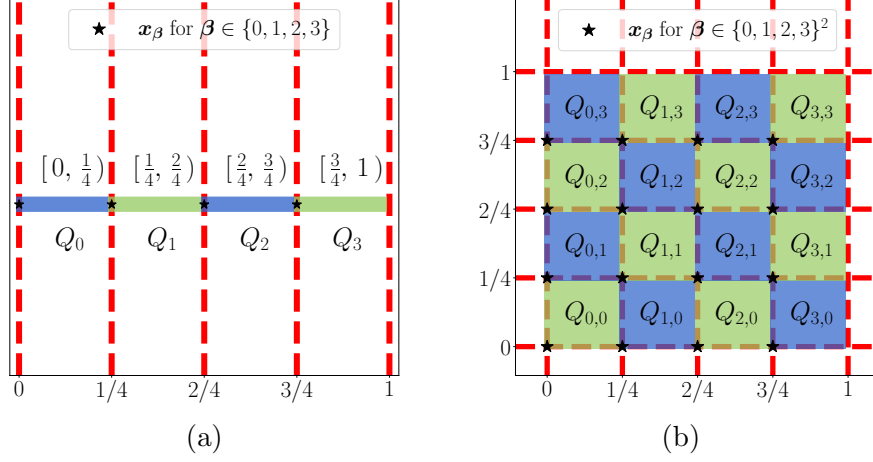


Figure 3: Illustrations of Q_β and \mathbf{x}_β for any $\beta \in \{0, 1, \dots, J-1\}^d$. (a) $J = 4$, $d = 1$. (b) $J = 4$, $d = 2$.

Step 2: Construct Φ_1 mapping $\mathbf{x} \in Q_\beta$ to β for each $\beta \in \{0, 1, \dots, J-1\}^d$.

Define

$$\Phi_1(\mathbf{x}) := \left(\sigma_1(Jx_1), \sigma_1(Jx_2), \dots, \sigma_1(Jx_d) \right) = \left(\lfloor Jx_1 \rfloor, \lfloor Jx_2 \rfloor, \dots, \lfloor Jx_d \rfloor \right),$$

for any $\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$. Then, for any $\mathbf{x} \in Q_\beta$ and each $\beta \in \{0, 1, \dots, J-1\}^d$, we have

$$\Phi_1(\mathbf{x}) = \left(\lfloor Jx_1 \rfloor, \lfloor Jx_2 \rfloor, \dots, \lfloor Jx_d \rfloor \right) = (\beta_1, \beta_2, \dots, \beta_d) = \beta. \quad (10)$$

Step 3: Construct ϕ_2 bijectively mapping $\beta \in \{0, 1, \dots, J-1\}^d$ to $\phi_2(\beta) \in \{1, 2, \dots, J^d\}$.

Inspired by the J -ary representation, we define a linear function

$$\phi_2(\mathbf{x}) := 1 + \sum_{i=1}^d J^{i-1} x_i, \quad \text{for each } \mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d.$$

Then ϕ_2 is a bijection from $\{0, 1, \dots, J-1\}^d$ to $\{1, 2, \dots, J^d\}$.

Step 4: Construct ϕ_3 mapping $\phi_2(\beta) \in \{1, 2, \dots, J^d\}$ approximately to $\tilde{f}(\mathbf{x}_\beta)$.

For each $k \in \{1, 2, \dots, J^d\}$, there exists a unique $\boldsymbol{\beta} \in \{0, 1, \dots, J-1\}^d$ such that $\phi_2(\boldsymbol{\beta}) = k$. Thus, define

$$\xi_k := \tilde{f}(\mathbf{x}_\beta) \in [0, 1], \quad \text{for any } k \in \{1, 2, \dots, J^d\} \text{ with } k = \phi_2(\boldsymbol{\beta}). \quad (11)$$

For each $k \in \{1, 2, \dots, J^d\}$, there exist $\theta_{k,1}, \theta_{k,2}, \dots, \theta_{k,N} \in \{0, 1\}$ such that

$$|\xi_k - \text{bin}0.\theta_{k,1}\theta_{k,2}\cdots\theta_{k,N}| \leq 2^{-N}. \quad (12)$$

For each $j \in \{1, 2, \dots, N\}$, by Proposition 3.2 (set $K = J^d$ therein), there exists $a_j \in [0, \frac{1}{2})$ such that

$$\sigma_3(2^k \cdot a_j) = \theta_{k,j}, \quad \text{for any } k \in \{1, 2, \dots, J^d\}.$$

Define

$$\phi_3(x) := \sum_{j=1}^N 2^{-j} \sigma_3(a_j \cdot \sigma_2(x)) = \sum_{j=1}^N 2^{-j} \sigma_3(2^x \cdot a_j), \quad \text{for any } x \in \mathbb{R}.$$

Then, for any $k \in \{1, 2, \dots, J^d\}$, we have

$$\phi_3(k) = \sum_{j=1}^N 2^{-j} \sigma_3(2^k \cdot a_j) = \sum_{j=1}^N 2^{-j} \cdot \theta_{k,j} = \text{bin}0.\theta_{k,1}\theta_{k,2}\cdots\theta_{k,N}. \quad (13)$$

Step 5: Define $\tilde{\phi} := \phi_3 \circ \phi_2 \circ \Phi_1$ approximating \tilde{f} well, and re-scale and shift $\tilde{\phi}$ to obtain ϕ approximating f well.

Define $\tilde{\phi} := \phi_3 \circ \phi_2 \circ \Phi_1$, by Equation (10), (11), (12), and (13), we have, for any $\mathbf{x} \in Q_\beta$ and each $\boldsymbol{\beta} \in \{0, 1, \dots, J-1\}^d$ with $k = \phi_2(\boldsymbol{\beta})$,

$$\begin{aligned} |\tilde{\phi}(\mathbf{x}) - \tilde{f}(\mathbf{x})| &\leq |\phi_3 \circ \phi_2 \circ \Phi_1(\mathbf{x}) - \tilde{f}(\mathbf{x}_\beta)| + |\tilde{f}(\mathbf{x}_\beta) - \tilde{f}(\mathbf{x})| \\ &\leq |\phi_3 \circ \phi_2(\boldsymbol{\beta}) - \tilde{f}(\mathbf{x}_\beta)| + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right) \leq |\phi_3(k) - \xi_k| + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right) \\ &\leq |\text{bin}0.\theta_{k,1}\theta_{k,2}\cdots\theta_{k,N} - \xi_k| + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right) \leq 2^{-N} + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right). \end{aligned}$$

Finally, define $\phi := 2\omega_f(\sqrt{d})\tilde{\phi} + f(\mathbf{0}) - \omega_f(\sqrt{d})$. Equation (9) implies $\omega_f(r) = 2\omega_f(\sqrt{d})\omega_{\tilde{f}}(r)$ for any $r \geq 0$, deducing

$$\begin{aligned} |\phi(\mathbf{x}) - f(\mathbf{x})| &= 2\omega_f(\sqrt{d})|\tilde{\phi}(\mathbf{x}) - \tilde{f}(\mathbf{x})| \leq 2\omega_f(\sqrt{d})(2^{-N} + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right)) \\ &= 2\omega_f(\sqrt{d})2^{-N} + \omega_f\left(\frac{\sqrt{d}}{J}\right) \\ &= 2\omega_f(\sqrt{d})2^{-N} + \omega_f(\sqrt{d}2^{-N}), \end{aligned}$$

for any $\mathbf{x} \in \bigcup_{\beta \in \{0,1,\dots,J-1\}^d} Q_\beta = [0,1]^d$. It follows from $J = 2^N$ and the definitions of Φ_1 , ϕ_2 , and ϕ_3 that

$$\begin{aligned}\phi(\mathbf{x}) &= 2\omega_f(\sqrt{d})\phi_3 \circ \phi_2 \circ \Phi_1(\mathbf{x}) + f(\mathbf{0}) - \omega_f(\sqrt{d}) \\ &= 2\omega_f(\sqrt{d})\phi_3\left(1 + \sum_{i=1}^d J^{i-1}\sigma_1(Jx_i)\right) + f(\mathbf{0}) - \omega_f(\sqrt{d}) \\ &= 2\omega_f(\sqrt{d})\sum_{j=1}^N 2^{-j}\sigma_3\left(a_j \cdot \sigma_2\left(1 + \sum_{i=1}^d 2^{(i-1)N}\sigma_1(2^N x_i)\right)\right) + f(\mathbf{0}) - \omega_f(\sqrt{d}).\end{aligned}$$

So we finish the proof. \square

4. Approximation with continuous activation functions

As discussed previously, our FLES networks can attain super approximation power. However, two activation functions in FLES networks are piecewise constant functions that would lead to challenges in numerical algorithm design. It is interesting to explore continuous activation functions achieving similar results. To this end, we introduce three new activation functions as follows. First, for any $\delta \in (0,1)$, we define

$$\varrho_{1,\delta}(x) := \begin{cases} n-1, & x \in [n-1, n-\delta], \\ (x-n+\delta)/\delta, & x \in (n-\delta, n], \end{cases} \quad \text{for any } n \in \mathbb{Z}.$$

In fact, $\varrho_{1,\delta}$ can be regarded as a ‘‘continuous version’’ of the floor function. Next, we define

$$\varrho_2(x) := 3^x, \quad \text{and} \quad \varrho_3(x) := \tilde{\mathcal{T}}(\cos(2\pi x)), \quad \text{for any } x \in \mathbb{R},$$

where

$$\tilde{\mathcal{T}}(x) := \begin{cases} 0, & x \in (\cos(\frac{4\pi}{9}), \infty), \\ 1 - x/\cos(\frac{4\pi}{9}), & x \in [0, \cos(\frac{4\pi}{9})], \\ 1, & x \in (-\infty, 0) \end{cases}$$

is a continuous piecewise linear function. ϱ_2 plays the same role of $\sigma_2(x) = 2^x$ and ϱ_3 is essentially a ‘‘continuous version’’ of σ_3 in FLES networks.

With these three activation functions in hand, we have the following theorem.

Theorem 4.1. *Let f be an arbitrary continuous function defined on $[0, 1]^d$. For any $\delta \in (0, 1)$, $N \in \mathbb{N}^+$, and $p \in [1, \infty)$, there exist $a_1, a_2, \dots, a_N \in [0, \frac{2}{9})$ such that*

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

where ϕ is defined by a formula in a_1, a_2, \dots, a_N as follows

$$\phi(\mathbf{x}) = 2\omega_f(\sqrt{d}) \sum_{j=1}^N 2^{-j} \varrho_3 \left(a_j \cdot \varrho_2 \left(1 + \sum_{i=1}^d 2^{(i-1)N} \varrho_{1,\delta}(2^N x_i) \right) \right) + f(\mathbf{0}) - \omega_f(\sqrt{d}).$$

The approximation error in Theorem 4.1 is characterized by L^p -norm for $p \in [1, \infty)$ instead of a pointwise error estimate in Theorem 1.1. By using ideas in Lu et al. (2020); Zhang (2020), we can extend this result to L^∞ -norm. However, this extension requires $2d + 3$ hidden layers instead of 3 hidden layers. Since our focus here is the approximation using three hidden layers, we will leave this extension as future work.

To prove Theorem 4.1, we need the following proposition.

Proposition 4.2. *Given any $K \in \mathbb{N}^+$ and arbitrary $\theta_1, \theta_2, \dots, \theta_K \in \{0, 1\}$, it holds that*

$$\varrho_3(a \cdot \varrho_2(k)) = \varrho_3(3^k \cdot a) = \theta_k, \quad \text{for any } k \in \{1, 2, \dots, K\},$$

where

$$a = \sum_{j=1}^K 3^{-j-1} \cdot \theta_j \in [0, \frac{2}{9}).$$

Proof. Since $\theta_j \in \{0, 1\}$ for $j \in \{1, 2, \dots, K\}$, we have

$$0 \leq \sum_{j=1}^K 3^{-j-1} \cdot \theta_j \leq \sum_{j=1}^K 3^{-j-1} < \frac{2}{9},$$

implying $a \in [0, \frac{2}{9})$.

Next, fix $k \in \{1, 2, \dots, K\}$ for the proof below. It holds that

$$3^k \cdot a = 3^k \cdot \sum_{j=1}^K 3^{-j-1} \cdot \theta_j = \underbrace{\sum_{j=1}^{k-1} 3^{k-j-1} \cdot \theta_j}_{\text{an integer}} + \overbrace{\frac{1}{3}\theta_k}^{0 \text{ or } \frac{1}{3}} + \underbrace{\sum_{j=k+1}^K 3^{k-j-1} \cdot \theta_j}_{\text{in } [0, \frac{2}{9})}. \quad (14)$$

Clearly, the first term $\sum_{j=1}^{k-1} 3^{k-j-1} \cdot \theta_j$ in Equation (14) is a non-negative integer since $\theta_j \in \{0, 1\}$ for any $j \in \{1, 2, \dots, K\}$. As for the third term in Equation (14), we have

$$0 \leq \sum_{j=k+1}^K 3^{k-j-1} \cdot \theta_j \leq \sum_{j=k+1}^K 3^{k-j-1} < \frac{2}{9}.$$

Recall that

$$\cos(2\pi x) \in (\cos(\frac{4\pi}{9}), 1], \quad \text{for any } x \in \bigcup_{n \in \mathbb{N}} [n, n + \frac{2}{9}),$$

and

$$\cos(2\pi x) \in [-1, \cos(\frac{2\pi}{3})] \subseteq [-1, 0], \quad \text{for any } x \in \bigcup_{n \in \mathbb{N}} [n + \frac{1}{3}, n + \frac{5}{9}).$$

Note that

$$\tilde{\mathcal{T}}(x) := \begin{cases} 0, & x \in (\cos(\frac{4\pi}{9}), \infty), \\ 1 - x / \cos(\frac{4\pi}{9}), & x \in [0, \cos(\frac{4\pi}{9})], \\ 1, & x \in (-\infty, 0). \end{cases}$$

Therefore, if $\theta_k = 0$, by Equation (14), we have

$$3^k \cdot a \in \bigcup_{n \in \mathbb{N}} [n, n + \frac{2}{9}) \implies \varrho_3(3^k \cdot a) = \tilde{\mathcal{T}}(\cos(2\pi \cdot 3^k \cdot a)) = 0 = \theta_k.$$

Similarly, if $\theta_k = 1$, by Equation (14), we have

$$3^k \cdot a \in \bigcup_{n \in \mathbb{N}} [n + \frac{1}{3}, n + \frac{5}{9}) \implies \varrho_3(3^k \cdot a) = \tilde{\mathcal{T}}(\cos(2\pi \cdot 3^k \cdot a)) = 1 = \theta_k.$$

Since $k \in \{1, 2, \dots, K\}$ is arbitrary, we have $\varrho_3(a \cdot \varrho_2(k)) = \varrho_3(3^k \cdot a) = \theta_k$ for any $k \in \{1, 2, \dots, K\}$. So we finish the proof. \square

Before proving Theorem 4.1, let us define a small region as follows to simplify the notation. Given any $J \in \mathbb{N}^+$ and $\delta \in (0, 1)$, define a small region $\Lambda([0, 1]^d, J, \delta)$ as

$$\Lambda([0, 1]^d, J, \delta) := \bigcup_{i=1}^d \left\{ \mathbf{x} = (x_1, \dots, x_d) \in [0, 1]^d : x_i \in \bigcup_{j=1}^{J-1} [\frac{j-\delta}{J}, \frac{j}{J}] \right\}.$$

In particular, $\Lambda([0, 1]^d, J, \delta) = \emptyset$ if $J = 1$. See Figure 4 for two examples.

With Proposition 4.2 in hand, we are ready to prove Theorem 4.1.

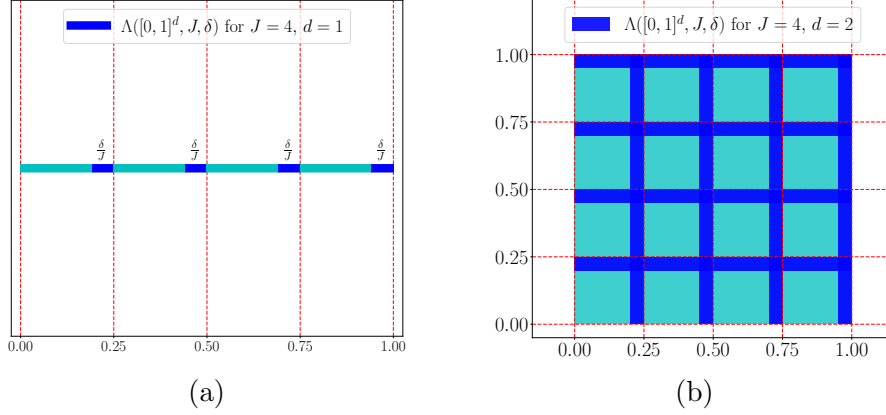


Figure 4: Illustrations of $\Lambda([0, 1]^d, J, \delta)$. (a) $J = 4, d = 1$. (b) $J = 4, d = 2$.

Proof of Theorem 4.1. The proof consists of five steps.

Step 1: Set up.

Assume f is not a constant function since it is a trivial case. Then $\omega_f(r) > 0$ for any $r > 0$. Clearly, $|f(\mathbf{x}) - f(\mathbf{0})| \leq \omega_f(\sqrt{d})$ for any $\mathbf{x} \in [0, 1]^d$. Define

$$\tilde{f} := \frac{f - f(\mathbf{0}) + \omega_f(\sqrt{d})}{2\omega_f(\sqrt{d})}. \quad (15)$$

It follows that $\tilde{f}(\mathbf{x}) \in [0, 1]$ for any $\mathbf{x} \in [0, 1]^d$.

Set $J = 2^N$ and divide $[0, 1]^d$ into J^d cubes $\{Q_\beta\}_\beta$ and a small region $\Lambda([0, 1]^d, J, \delta)$. To be exact, define $\mathbf{x}_\beta := \beta/J$ and

$$Q_\beta := \left\{ \mathbf{x} = (x_1, x_2, \dots, x_d) : x_i \in \left[\frac{\beta_i}{J}, \frac{\beta_i + 1 - \delta}{J} \right] \text{ for } i = 1, 2, \dots, d \right\},$$

for each $\beta = (\beta_1, \beta_2, \dots, \beta_d) \in \{0, 1, \dots, J-1\}^d$. See Figure 5 for illustrations.

Step 2: Construct Φ_1 mapping $\mathbf{x} \in Q_\beta$ to β for each $\beta \in \{0, 1, \dots, J-1\}^d$.

Define

$$\Phi_1(\mathbf{x}) := \left(\varrho_{1,\delta}(Jx_1), \varrho_{1,\delta}(Jx_2), \dots, \varrho_{1,\delta}(Jx_d) \right),$$

for any $\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$. Then, for any $\mathbf{x} \in Q_\beta$ and each $\beta \in \{0, 1, \dots, J-1\}^d$, we have

$$\Phi_1(\mathbf{x}) = \left(\varrho_{1,\delta}(Jx_1), \varrho_{1,\delta}(Jx_2), \dots, \varrho_{1,\delta}(Jx_d) \right) = (\beta_1, \beta_2, \dots, \beta_d) = \beta. \quad (16)$$

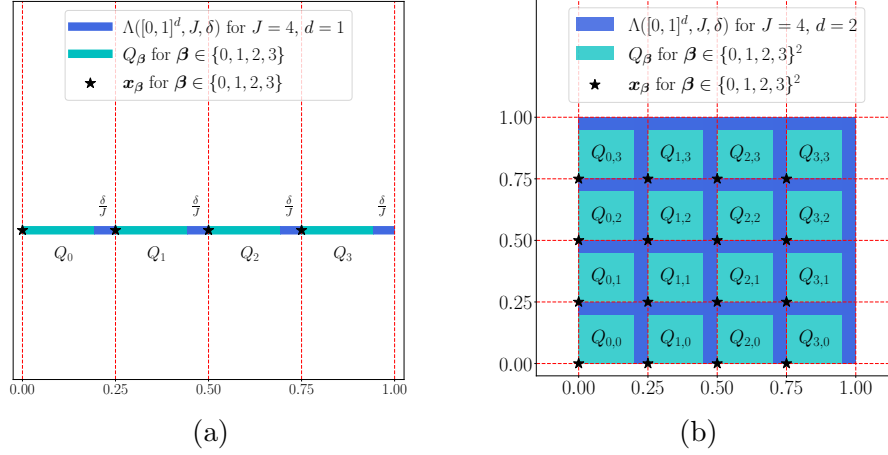


Figure 5: Illustrations of $\Lambda([0, 1]^d, J, \delta)$, Q_β , and \mathbf{x}_β for any $\beta \in \{0, 1, \dots, J-1\}^d$. (a) $J=4, d=1$. (b) $J=4, d=2$.

Step 3: Construct ϕ_2 bijectively mapping $\beta \in \{0, 1, \dots, J-1\}^d$ to $\phi_2(\beta) \in \{1, 2, \dots, J^d\}$.

Inspired by the J -ary representation, we define an affine linear map

$$\phi_2(\mathbf{x}) := 1 + \sum_{i=1}^d J^{i-1} x_i, \quad \text{for each } \mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d.$$

Then ϕ_2 is a bijection from $\{0, 1, \dots, J-1\}^d$ to $\{1, 2, \dots, J^d\}$.

Step 4: Construct ϕ_3 mapping $\phi_2(\beta) \in \{1, 2, \dots, J^d\}$ approximately to $\tilde{f}(\mathbf{x}_\beta)$.

For each $k \in \{1, 2, \dots, J^d\}$, there exists a unique $\beta \in \{0, 1, \dots, J-1\}^d$ such that $\phi_2(\beta) = k$. Thus, define

$$\xi_k := \tilde{f}(\mathbf{x}_\beta) \in [0, 1], \quad \text{for any } k \in \{1, 2, \dots, J^d\} \text{ with } k = \phi_2(\beta). \quad (17)$$

For each $k \in \{1, 2, \dots, J^d\}$, there exist $\theta_{k,1}, \theta_{k,2}, \dots, \theta_{k,N} \in \{0, 1\}$ such that

$$|\xi_k - \text{bin}0.\theta_{k,1}\theta_{k,2}\dots\theta_{k,N}| \leq 2^{-N}. \quad (18)$$

For each $j \in \{1, 2, \dots, N\}$, by Proposition 4.2 (set $K = J^d$ therein), there exists $a_j \in [0, \frac{2}{9})$ such that

$$\varrho_3(3^k \cdot a_j) = \theta_{k,j}, \quad \text{for any } k \in \{1, 2, \dots, J^d\}.$$

Define

$$\phi_3(x) := \sum_{j=1}^N 2^{-j} \varrho_3(a_j \cdot \varrho_2(x)) = \sum_{j=1}^N 2^{-j} \varrho_3(3^x \cdot a_j), \quad \text{for any } x \in \mathbb{R}.$$

Then we have

$$\varrho_3(x) \in [0, 1], \quad \text{for any } x \in \mathbb{R} \quad \implies \quad \phi_3(x) \in [0, 1], \quad \text{for any } x \in \mathbb{R}, \quad (19)$$

and

$$\phi_3(k) = \sum_{j=1}^N 2^{-j} \varrho_3(3^k \cdot a_j) = \sum_{j=1}^N 2^{-j} \cdot \theta_{k,j} = \text{bin}0.\theta_{k,1}\theta_{k,2}\cdots\theta_{k,N}, \quad (20)$$

for any $k \in \{1, 2, \dots, J^d\}$.

Step 5: Define $\tilde{\phi} := \phi_3 \circ \phi_2 \circ \Phi_1$ approximating \tilde{f} well, and re-scale and shift $\tilde{\phi}$ to obtain ϕ approximating f well.

Define $\tilde{\phi} := \phi_3 \circ \phi_2 \circ \Phi_1$, by Equation (16), (17), (18), and (20), we have, for any $\mathbf{x} \in Q_\beta$ and each $\beta \in \{0, 1, \dots, J-1\}^d$ with $k = \phi_2(\beta)$,

$$\begin{aligned} |\tilde{\phi}(\mathbf{x}) - \tilde{f}(\mathbf{x})| &= |\phi_3 \circ \phi_2 \circ \Phi_1(\mathbf{x}) - \tilde{f}(\mathbf{x}_\beta)| + |\tilde{f}(\mathbf{x}_\beta) - \tilde{f}(\mathbf{x})| \\ &\leq |\phi_3 \circ \phi_2(\beta) - \tilde{f}(\mathbf{x}_\beta)| + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right) \leq |\phi_3(k) - \xi_k| + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right) \\ &\leq |\text{bin}0.\theta_{k,1}\theta_{k,2}\cdots\theta_{k,N} - \xi_k| + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right) \leq 2^{-N} + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right). \end{aligned}$$

Finally, define $\phi := 2\omega_f(\sqrt{d})\tilde{\phi} + f(\mathbf{0}) - \omega_f(\sqrt{d})$. Equation (15) implies $\omega_f(r) = 2\omega_f(\sqrt{d})\omega_{\tilde{f}}(r)$ for any $r \geq 0$, deducing

$$\begin{aligned} |\phi(\mathbf{x}) - f(\mathbf{x})| &= 2\omega_f(\sqrt{d})|\tilde{\phi}(\mathbf{x}) - \tilde{f}(\mathbf{x})| \leq 2\omega_f(\sqrt{d})(2^{-N} + \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{J}\right)) \\ &= 2\omega_f(\sqrt{d})2^{-N} + \omega_f\left(\frac{\sqrt{d}}{J}\right), \end{aligned}$$

for any $\mathbf{x} \in \bigcup_{\beta \in \{0,1,\dots,J-1\}^d} Q_\beta$. By Equation (19) and the definition of

$$\phi = 2\omega_f(\sqrt{d})\phi_3 \circ \phi_2 \circ \Phi_1 + f(\mathbf{0}) - \omega_f(\sqrt{d}),$$

we have $\|\phi\|_{L^\infty(\mathbb{R}^d)} \leq |f(\mathbf{0})| + \omega_f(\sqrt{d})$. Let $\mu(\cdot)$ denote the Lebesgue measure. Note that $\|f\|_{L^\infty([0,1]^d)} \leq |f(\mathbf{0})| + \omega_f(\sqrt{d})$. It follows from $\mu(\Lambda([0,1]^d, J, \delta)) \leq$

$Jd\frac{\delta}{J} = d\delta$ that

$$\begin{aligned}
& \|\phi - f\|_{L^p([0,1]^d)}^p = \int_{[0,1]^d} |\phi(\mathbf{x}) - f(\mathbf{x})|^p d\mathbf{x} \\
&= \sum_{\beta \in \{0,1,\dots,J-1\}^d} \int_{Q_\beta} |\phi(\mathbf{x}) - f(\mathbf{x})|^p d\mathbf{x} + \int_{\Lambda([0,1]^d, J, \delta)} |\phi(\mathbf{x}) - f(\mathbf{x})|^p d\mathbf{x} \\
&\leq \sum_{\beta \in \{0,1,\dots,J-1\}^d} \mu(Q_\beta) \left(2\omega_f(\sqrt{d})2^{-N} + \omega_f\left(\frac{\sqrt{d}}{J}\right) \right)^p + (2|f(\mathbf{0})| + 2\omega_f(\sqrt{d}))^p d\delta \\
&\leq \left(2\omega_f(\sqrt{d})2^{-N} + \omega_f(\sqrt{d}2^{-N}) \right)^p + 2^p d\delta (|f(\mathbf{0})| + \omega_f(\sqrt{d}))^p.
\end{aligned}$$

By the definitions of Φ_1 , ϕ_2 , and ϕ_3 , we have

$$\begin{aligned}
\phi(\mathbf{x}) &= 2\omega_f(\sqrt{d})\phi_3 \circ \phi_2 \circ \Phi_1(\mathbf{x}) + f(\mathbf{0}) - \omega_f(\sqrt{d}) \\
&= 2\omega_f(\sqrt{d})\phi_3 \left(1 + \sum_{i=1}^d J^{i-1} \varrho_{1,\delta}(Jx_i) \right) + f(\mathbf{0}) - \omega_f(\sqrt{d}) \\
&= 2\omega_f(\sqrt{d}) \sum_{j=1}^N 2^{-j} \varrho_3 \left(a_j \cdot \varrho_2 \left(1 + \sum_{i=1}^d 2^{(i-1)N} \varrho_{1,\delta}(2^N x_i) \right) \right) + f(\mathbf{0}) - \omega_f(\sqrt{d}).
\end{aligned}$$

So we finish the proof. \square

5. Conclusion

This paper has introduced a theoretical framework to show that three hidden layers are enough for neural network approximation to achieve exponential convergence and avoid the curse of dimensionality for approximating functions as general as (Hölder) continuous functions. The key idea is to leverage the power of multiple simple activation functions: the floor function ($\lfloor x \rfloor$), the exponential function (2^x), the step function ($\mathbf{1}_{x \geq 0}$), or their compositions. This new class of networks is called the FLES network. Given a Lipschitz continuous function f on $[0,1]^d$, it was shown by construction that FLES networks with width $\max\{d, N\}$ and three hidden layers admit a uniform approximation rate $6\lambda\sqrt{d}2^{-N}$, where λ is the Lipschitz constant of f . More generally for an arbitrary continuous function f on $[0,1]^d$ with a modulus of continuity $\omega_f(\cdot)$, the constructive approximation rate is $2\omega_f(2\sqrt{d})2^{-N} + \omega_f(2\sqrt{d}2^{-N})$. We also extend such a result to general bounded continuous functions on a bounded set $E \subseteq \mathbb{R}^d$. The results

in this paper provide a theoretical lower bound of the power of FLES networks. Whether or not this bound is achievable in actual computation relies on advanced algorithm design as a separate line of research. Finally, we have also derived similar approximation results in the L^p -norm for $p \in [1, \infty)$ using continuous activation 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.

References

- Arnold, V.I., 1957. On functions of three variables. Dokl. Akad. Nauk SSSR 114, 679–681. doi:[978-3-642-01742-1_2](https://doi.org/10.1007/s00365-009-9054-2).
- Barron, A.R., 1993. Universal approximation bounds for superpositions of a sigmoidal function. IEEE Transactions on Information Theory 39, 930–945. doi:[10.1109/18.256500](https://doi.org/10.1109/18.256500).
- Barron, A.R., Klusowski, J.M., 2018. Approximation and estimation for high-dimensional deep learning networks. arXiv e-prints [arXiv:1809.03090](https://arxiv.org/abs/1809.03090).
- Bartlett, P., Maiorov, V., Meir, R., 1998. Almost linear VC-dimension bounds for piecewise polynomial networks. Neural Computation 10, 21592173. doi:[10.1162/089976698300017016](https://doi.org/10.1162/089976698300017016).
- Bengio, Y., Léonard, N., Courville, A., 2013. Estimating or propagating gradients through stochastic neurons for conditional computation. arXiv e-print [arXiv:1308.3432](https://arxiv.org/abs/1308.3432).
- Boo, Y., Shin, S., Sung, W., 2020. Quantized neural networks: Characterization and holistic optimization. arXiv e-print [arXiv:2006.00530](https://arxiv.org/abs/2006.00530).
- Braun, J., Griebel, M., 2009. On a constructive proof of Kolmogorov’s superposition theorem. Constructive Approximation 30, 653–675. doi:[10.1007/s00365-009-9054-2](https://doi.org/10.1007/s00365-009-9054-2).
- Carrillo, J.A.T., Jin, S., Li, L., Zhu, Y., 2019. A consensus-based global optimization method for high dimensional machine learning problems. arXiv e-print [arXiv:1909.09249](https://arxiv.org/abs/1909.09249).

- Chen, L., Wu, C., 2019. A note on the expressive power of deep rectified linear unit networks in high-dimensional spaces. *Mathematical Methods in the Applied Sciences* 42, 3400–3404. doi:[10.1002/mma.5575](https://doi.org/10.1002/mma.5575).
- Chen, M., Jiang, H., Liao, W., Zhao, T., 2019a. Efficient approximation of deep ReLU networks for functions on low dimensional manifolds, in: Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., Garnett, R. (Eds.), *Advances in Neural Information Processing Systems* 32. Curran Associates, Inc., pp. 8174–8184.
- Chen, Z., Cao, Y., Zou, D., Gu, Q., 2019b. How much over-parameterization is sufficient to learn deep ReLU networks? CoRR arXiv:1911.12360. URL: <https://arxiv.org/abs/1911.12360>.
- Du, S.S., Zhai, X., Póczos, B., Singh, A., 2019. Gradient descent provably optimizes over-parameterized neural networks, in: *International Conference on Learning Representations*. URL: <https://openreview.net/forum?id=S1eK3i09YQ>.
- E, W., Ma, C., Wu, L., 2019. A priori estimates of the population risk for two-layer neural networks. *Communications in Mathematical Sciences* 17, 1407–1425. doi:[10.4310/CMS.2019.v17.n5.a11](https://doi.org/10.4310/CMS.2019.v17.n5.a11).
- E, W., Wang, Q., 2018. Exponential convergence of the deep neural network approximation for analytic functions. CoRR abs/1807.00297. URL: <http://arxiv.org/abs/1807.00297>, [arXiv:1807.00297](https://arxiv.org/abs/1807.00297).
- Gühring, I., Kutyniok, G., Petersen, P., 2019. Error bounds for approximations with deep ReLU neural networks in $W^{s,p}$ norms. arXiv e-prints [arXiv:1902.07896](https://arxiv.org/abs/1902.07896).
- Guliyev, N.J., Ismailov, V.E., 2018. Approximation capability of two hidden layer feedforward neural networks with fixed weights. *Neurocomputing* 316, 262–269. doi:[10.1016/j.neucom.2018.07.075](https://doi.org/10.1016/j.neucom.2018.07.075).
- Harvey, N., Liaw, C., Mehrabian, A., 2017. Nearly-tight VC-dimension bounds for piecewise linear neural networks, in: Kale, S., Shamir, O. (Eds.), *Proceedings of the 2017 Conference on Learning Theory*, PMLR, Amsterdam, Netherlands. pp. 1064–1068. URL: <http://proceedings.mlr.press/v65/harvey17a.html>.

- Holland, J.H., 1992. Genetic algorithms. *Scientific American* 267, 66–73. URL: <http://www.jstor.org/stable/24939139>.
- Hubara, I., Courbariaux, M., Soudry, D., El-Yaniv, R., Bengio, Y., 2017. Quantized neural networks: Training neural networks with low precision weights and activations. *J. Mach. Learn. Res.* 18, 68696898.
- Hutzenthaler, M., Jentzen, A., Wurstemberger, v.W., 2020. Overcoming the curse of dimensionality in the approximative pricing of financial derivatives with default risks. *Electron. J. Probab.* 25, 73 pp. doi:[10.1214/20-EJP423](https://doi.org/10.1214/20-EJP423).
- Igel'nik, B., Parikh, N., 2003. Kolmogorov's spline network. *IEEE Transactions on Neural Networks* 14, 725–733. doi:[10.1109/TNN.2003.813830](https://doi.org/10.1109/TNN.2003.813830).
- Jacot, A., Gabriel, F., Hongler, C., 2018. Neural tangent kernel: Convergence and generalization in neural networks, in: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (Eds.), *Advances in Neural Information Processing Systems*. Curran Associates, Inc.. volume 31, pp. 8571–8580. URL: <https://proceedings.neurips.cc/paper/2018/file/5a4be1fa34e62bb8a6ec6b91d2462f5a-Paper.pdf>.
- Kennedy, J., Eberhart, R., 1995. Particle swarm optimization, in: *Proceedings of ICNN'95 - International Conference on Neural Networks*, pp. 1942–1948 vol.4. doi:[10.1109/ICNN.1995.488968](https://doi.org/10.1109/ICNN.1995.488968).
- Kirkpatrick, S., Gelatt, C.D., Vecchi, M.P., 1983. Optimization by simulated annealing. *Science* 220, 671–680. doi:[10.1126/science.220.4598.671](https://doi.org/10.1126/science.220.4598.671).
- Kolmogorov, A.N., 1956. On the representation of continuous functions of several variables by superposition of continuous functions of a smaller number of variables. *Dokl. Akad. Nauk SSSR* 108, 179–182. doi:[10.1007/978-3-642-01742-1_5](https://doi.org/10.1007/978-3-642-01742-1_5).
- Kolmogorov, A.N., 1957. On the representation of continuous functions of several variables by superposition of continuous functions of one variable and addition. *Dokl. Akad. Nauk SSSR* 114, 953–956. URL: <http://mi.mathnet.ru/dan22050>.
- Kůrková, V., 1991. Kolmogorov's theorem is relevant. *Neural Computation* 3, 617–622. doi:[10.1162/neco.1991.3.4.617](https://doi.org/10.1162/neco.1991.3.4.617).

- Kůrková, V., 1992. Kolmogorov’s theorem and multilayer neural networks. *Neural Networks* 5, 501–506. doi:[10.1016/0893-6080\(92\)90012-8](https://doi.org/10.1016/0893-6080(92)90012-8).
- Li, Q., Lin, T., Shen, Z., to appear. Deep learning via dynamical systems: An approximation perspective. *Journal of European Mathematical Society* .
- Lin, Y., Lei, M., Niu, L., 2019. Optimization strategies in quantized neural networks: A review, in: 2019 International Conference on Data Mining Workshops (ICDMW), pp. 385–390. doi:[10.1109/ICDMW.2019.00063](https://doi.org/10.1109/ICDMW.2019.00063).
- Lu, J., Shen, Z., Yang, H., Zhang, S., 2020. Deep network approximation for smooth functions. *arXiv e-prints* [arXiv:2001.03040](https://arxiv.org/abs/2001.03040).
- Lu, Y., Ma, C., Lu, Y., Lu, J., Ying, L., 2020. A mean-field analysis of deep resnet and beyond: Towards provable optimization via overparameterization from depth. *CoRR* abs/2003.05508. [arXiv:2003.05508](https://arxiv.org/abs/2003.05508).
- Luo, T., Yang, H., 2020. Two-Layer Neural Networks for Partial Differential Equations: Optimization and Generalization Theory. *arXiv e-prints* [arXiv:2006.15733](https://arxiv.org/abs/2006.15733).
- Maierov, V., Pinkus, A., 1999. Lower bounds for approximation by MLP neural networks. *Neurocomputing* 25, 81–91. doi:[10.1016/S0925-2312\(98\)00111-8](https://doi.org/10.1016/S0925-2312(98)00111-8).
- Mei, S., Montanari, A., Nguyen, P.M., 2018. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences* 115, E7665–E7671. doi:[10.1073/pnas.1806579115](https://doi.org/10.1073/pnas.1806579115).
- Montanelli, H., Du, Q., 2019. New error bounds for deep ReLU networks using sparse grids. *SIAM Journal on Mathematics of Data Science* 1, 78–92. doi:[10.1137/18M1189336](https://doi.org/10.1137/18M1189336).
- Montanelli, H., Yang, H., 2020. Error bounds for deep ReLU networks using the Kolmogorov-Arnold superposition theorem. *Neural Networks* 129, 1–6. doi:[10.1016/j.neunet.2019.12.013](https://doi.org/10.1016/j.neunet.2019.12.013).
- Montanelli, H., Yang, H., Du, Q., 2020. Deep ReLU networks overcome the curse of dimensionality for bandlimited functions. *Journal of Computational Mathematics* .

- Nelder, J., Mead, R., 1965. A simplex method for function minimization. *Comput. J.* 7, 308–313. doi:[10.1093/comjnl/7.4.308](https://doi.org/10.1093/comjnl/7.4.308).
- Opschoor, J.A., Schwab, C., Zech, J., 2019. Exponential ReLU DNN expression of holomorphic maps in high dimension. Technical Report. Seminar for Applied Mathematics, ETH Zürich. Zurich. URL: <https://math.ethz.ch/sam/research/reports.html?id=839>.
- Petersen, P., Voigtlaender, F., 2018. Optimal approximation of piecewise smooth functions using deep ReLU neural networks. *Neural Networks* 108, 296–330. doi:[10.1016/j.neunet.2018.08.019](https://doi.org/10.1016/j.neunet.2018.08.019).
- Pinnau, R., Totzeck, C., Tse, O., Martin, S., 2017. A consensus-based model for global optimization and its mean-field limit. *Mathematical Models and Methods in Applied Sciences* 27, 183–204. doi:[10.1142/S0218202517400061](https://doi.org/10.1142/S0218202517400061).
- Poggio, T., Mhaskar, H.N., Rosasco, L., Miranda, B., Liao, Q., 2017. Why and when can deep—but not shallow—networks avoid the curse of dimensionality: A review. *International Journal of Automation and Computing* 14, 503–519. doi:[10.1007/s11633-017-1054-2](https://doi.org/10.1007/s11633-017-1054-2).
- Schmidt-Hieber, J., 2020. Nonparametric regression using deep neural networks with ReLU activation function. *Annals of Statistics* 48, 1875–1897. URL: <https://projecteuclid.org/euclid.aos/1597370649>.
- Schmidt-Hieber, J., 2021. The KolmogorovArnold representation theorem revisited. *Neural Networks* 137, 119–126. doi:[10.1016/j.neunet.2021.01.020](https://doi.org/10.1016/j.neunet.2021.01.020).
- Shen, Z., Yang, H., Zhang, S., 2019. Nonlinear approximation via compositions. *Neural Networks* 119, 74–84. doi:[10.1016/j.neunet.2019.07.011](https://doi.org/10.1016/j.neunet.2019.07.011).
- Shen, Z., Yang, H., Zhang, S., 2020. Deep network approximation characterized by number of neurons. *Communications in Computational Physics* 28, 1768–1811. doi:[10.4208/cicp.0A-2020-0149](https://doi.org/10.4208/cicp.0A-2020-0149).
- Shen, Z., Yang, H., Zhang, S., 2021. Deep network with approximation error being reciprocal of width to power of square root of depth. *Neural Computation* 33, 1005–1036. doi:[10.1162/neco_a_01364](https://doi.org/10.1162/neco_a_01364).

- Shen, Z., Yang, H., Zhang, S., to appear. Optimal approximation rate of relu networks in terms of width and depth. *Journal de Mathématiques Pures et Appliquées* .
- Wang, P., Hu, Q., Zhang, Y., Zhang, C., Liu, Y., Cheng, J., 2018. Two-step quantization for low-bit neural networks, in: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4376–4384. doi:[10.1109/CVPR.2018.00460](https://doi.org/10.1109/CVPR.2018.00460).
- Wu, L., Ma, C., E, W., 2018. How sgd selects the global minima in over-parameterized learning: A dynamical stability perspective, in: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (Eds.), *Advances in Neural Information Processing Systems 31*. Curran Associates, Inc., pp. 8279–8288. URL: <https://papers.nips.cc/paper/2018/hash/6651526b6fb8f29a00507de6a49ce30f-Abstract.html>.
- Yang, Y., Wang, Y., 2020. Approximation in shift-invariant spaces with deep ReLU neural networks. arXiv e-prints [arXiv:2005.11949](https://arxiv.org/abs/2005.11949).
- Yarotsky, D., 2017. Error bounds for approximations with deep ReLU networks. *Neural Networks* 94, 103–114. doi:[10.1016/j.neunet.2017.07.002](https://doi.org/10.1016/j.neunet.2017.07.002).
- Yarotsky, D., 2018. Optimal approximation of continuous functions by very deep ReLU networks, in: Bubeck, S., Perchet, V., Rigollet, P. (Eds.), *Proceedings of the 31st Conference On Learning Theory*, PMLR. pp. 639–649. URL: <http://proceedings.mlr.press/v75/yarotsky18a.html>.
- Yarotsky, D., Zhevnerchuk, A., 2020. The phase diagram of approximation rates for deep neural networks 33, 13005–13015. URL: https://proceedings.neurips.cc/paper_files/paper/2020/hash/979a3f14bae523dc5101c52120c535e9-Abstract.html.
- Yin, P., Lyu, J., Zhang, S., Osher, S.J., Qi, Y., Xin, J., 2019. Understanding straight-through estimator in training activation quantized neural nets URL: <https://openreview.net/forum?id=Skh4jRcKQ>.
- Zhang, S., 2020. Deep neural network approximation via function compositions. PhD Thesis, National University of Singapore URL: <https://scholarbank.nus.edu.sg/handle/10635/186064>.