

**DEEP NEURAL NETWORK  
APPROXIMATION VIA FUNCTION  
COMPOSITIONS**

**ZHANG SHIJUN**

*(B.Sc., Wuhan University, China)*

**A THESIS SUBMITTED  
FOR THE DEGREE OF DOCTOR OF PHILOSOPHY  
DEPARTMENT OF MATHEMATICS  
NATIONAL UNIVERSITY OF SINGAPORE  
2020**

Supervisors:

Professor Shen Zuowei, Main Supervisor

Assistant Professor Yang Haizhao, Co-Supervisor

Examiners:

Associate Professor Ji Hui

Assistant Professor Li Qianxiao

Professor Cai Jianfeng, Hong Kong University of Science and Technology



To my family

## DECLARATION

I hereby declare that the thesis is my original work and it has  
been written by me in its entirety. I have duly  
acknowledged all the sources of information which  
have been used in the thesis.

This thesis has also not been submitted for any degree  
in any university previously.

A handwritten signature in black ink, reading "Zhang Shijun" in a cursive style.

---

Zhang Shijun

December 30, 2020

---

# Acknowledgments

---

This dissertation would never be completed without the guidance of my supervisors, the help from my friends, and the support from my family.

First, I would like to express my sincere gratitude to my main supervisor professor Shen Zuowei, for his immense knowledge and research experience, and his guidance through each stage of my Ph.D. studies. His research philosophy plays a key role in defining the path of my research. During our discussions, I have gained a lot of knowledge and skills, especially the way of thinking and doing research. Without his constructive comments and suggestions, I would hardly complete my research work and this dissertation. It is my honor and pleasure to work with such an outstanding main supervisor.

Next, I would like to acknowledge my co-supervisor assistant professor Yang Haizhao, for his patience and enthusiasm, and his continuous support for my Ph.D. studies and related research. He convincingly guided and encouraged me to be professional. He has taught me a lot of things, including writing a research paper professionally, expressing my own opinions regarding a research problem precisely and clearly, *etc.* They benefited me a lot during my Ph.D. studies. Without his insightful feedback, I would hardly complete my research work and this dissertation. I am so lucky to do research with such a wonderful co-supervisor.

Besides my supervisors, I am also immensely grateful to associate professor Ji Hui, all members of the NUS Wavelet group, and all my classmates in the same student office as me. The numerous discussions with them in the group seminars or the student office helped me improve my knowledge in the research topics and develop a lot of skills for computer programming.

Finally, I would like to thank my family for their encouragement and understanding. My grandparents and parents raised me and supported me in achieving my pursuits. They kept me going on and this dissertation would not have been possibly finished without their input.

# Contents

|  |           |
|--|-----------|
| <b>Acknowledgments</b>                                   | <b>v</b>  |
| <b>Summary</b>   | <b>ix</b> |
| <b>1 Introduction</b>                                    | <b>1</b>  |
| 1.1 Contributions . . . . .                              | 2         |
| 1.2 Related work . . . . .                               | 5         |
| <b>2 Preliminaries</b>                                   | <b>9</b>  |
| 2.1 Notations . . . . .                                  | 9         |
| 2.1.1 Basic notations . . . . .                          | 9         |
| 2.1.2 Set notations . . . . .                            | 12        |
| 2.1.3 Neural network notations . . . . .                 | 13        |
| 2.2 Architecture of neural networks . . . . .            | 14        |
| 2.2.1 Descriptions . . . . .                             | 15        |
| 2.2.2 Compositions and combinations . . . . .            | 17        |
| 2.3 General ideas of approximation by networks . . . . . | 20        |
| 2.3.1 ReLU networks . . . . .                            | 20        |
| 2.3.2 Floor-ReLU networks . . . . .                      | 23        |
| <b>3 Basic results of ReLU networks</b>                  | <b>25</b> |
| 3.1 Wide networks to deep ones . . . . .                 | 25        |
| 3.2 Power of networks to fit points . . . . .            | 27        |
| 3.2.1 Width power of networks to fit points . . . . .    | 28        |

|          |  |            |
|----------|--|------------|
| 3.2.2    | Depth power of networks to fit points . . . . .          | 40         |
| 3.3      | Approximation in the trifling region . . . . .           | 48         |
| 3.4      | Approximation of step functions . . . . .                | 59         |
| <b>4</b> | <b>Approximation by ReLU networks</b>                    | <b>65</b>  |
| 4.1      | Approximation of polynomials . . . . .                   | 65         |
| 4.1.1    | Main theorem . . . . .                                   | 65         |
| 4.1.2    | Approximation of $x^2$ . . . . .                         | 66         |
| 4.1.3    | Approximation of $x_1x_2\cdots x_k$ . . . . .            | 70         |
| 4.1.4    | Proof of main theorem . . . . .                          | 75         |
| 4.2      | Approximation of continuous functions . . . . .          | 76         |
| 4.2.1    | Main theorem and its proof . . . . .                     | 76         |
| 4.2.2    | Proof of auxiliary theorem . . . . .                     | 79         |
| 4.2.3    | Proof of key proposition for auxiliary theorem . . . . . | 86         |
| 4.3      | Approximation of smooth functions . . . . .              | 91         |
| 4.3.1    | Main theorem and its proof . . . . .                     | 91         |
| 4.3.2    | Ideas of proving auxiliary theorem . . . . .             | 94         |
| 4.3.3    | Proof of auxiliary theorem . . . . .                     | 96         |
| 4.3.4    | Proof of key proposition for auxiliary theorem . . . . . | 103        |
| 4.4      | Optimality of approximation by networks . . . . .        | 106        |
| 4.4.1    | Hölder continuous functions . . . . .                    | 107        |
| 4.4.2    | Smooth functions . . . . .                               | 111        |
| <b>5</b> | <b>Approximation by Floor-ReLU networks</b>              | <b>117</b> |
| 5.1      | Main theorem and its proof . . . . .                     | 117        |
| 5.2      | Proof of auxiliary theorem . . . . .                     | 122        |
| 5.3      | Proof of key proposition for auxiliary theorem . . . . . | 128        |
| <b>6</b> | <b>Conclusion</b>  | <b>135</b> |
|          | <b>Bibliography</b>                                      | <b>137</b> |



---

# Summary

---

This dissertation is a summary of our previous papers [38, 52, 53, 54], focusing on the approximation theory of neural networks. We provide (nearly optimal) approximation error estimates in terms of the width and depth when constructing ReLU ( $\max\{0, x\}$ ) networks, via the idea of function compositions, to uniformly approximate polynomials, (Hölder) continuous functions, and smooth functions on a hypercube  $[0, 1]^d$ . The optimality is discussed via studying the connection between the approximation error and VC-dimension.<sup>①</sup> To overcome the limitation of ReLU networks that (nearly) exponential approximation errors<sup>②</sup> hold only for polynomials among all function spaces considered, we introduce new networks built with either Floor ( $\lfloor x \rfloor$ ) or ReLU as the activation function in each neuron. We call such networks Floor-ReLU networks. It is proved by construction that nearly exponential approximation errors can be attained when using Floor-ReLU networks to approximate (Hölder) continuous functions on  $[0, 1]^d$ . See Table 1.1 for a summary.

Chapter 1 is the introduction of this dissertation, including our main contributions and the literature review.

---

<sup>①</sup>See the definition of VC-dimension in Section 4.4.

<sup>②</sup>Throughout this dissertation, “exponential (approximation) error(s)” means “(approximation) error(s) with exponential decay”, similar to [21, 46, 54].

Chapter 2 is the preliminary chapter. In this chapter, we present the notations used throughout this dissertation, discuss the architecture of neural networks, and provide the general ideas of constructing neural networks to approximate given functions.

In Chapter 3, we prove several basic results of ReLU networks, which will be used in later chapters. The chapter consists of four parts. The first part investigates representing shallow ReLU networks by deep ones with a similar number of neurons. Part 2 discusses the width and depth power of ReLU networks to fit points. The third part looks at the approximation in a small region if a ReLU network approximates the target function well except for this small region. That is, we modify this network to let it approximate the target function uniformly well on the whole region. The final part deals with the approximation of step functions by ReLU networks.

Chapter 4 focuses on the approximation by ReLU networks and is divided into four parts. The first part aims to construct ReLU networks to approximate general polynomials on  $[0, 1]^d$  with exponential approximation errors. The second and third parts provide the detailed constructions of ReLU networks for approximating (Hölder) continuous functions and smooth functions on  $[0, 1]^d$  with (nearly optimal) approximation errors, respectively. The final part looks at the optimality of the approximation by ReLU networks via studying the connection between the approximation error and VC-dimension.

Chapter 5 aims to reveal the approximation power of Floor-ReLU networks. We provide nearly exponential approximation error estimates when constructing Floor-ReLU networks with fixed architectures<sup>③</sup> to uniformly approximate (Hölder) continuous functions on  $[0, 1]^d$ . In other words, the approximation errors are improved from polynomial ones to nearly exponential ones by adding a simple activation function (Floor) to ReLU networks.

Chapter 6 concludes this dissertation with a short discussion.

---

<sup>③</sup>A Floor-ReLU network with a fixed architecture means all the components of this network architecture is determined except for the values of the parameters. In particular, the choice of activation functions (Floor or ReLU) in each neuron is independent of the target function.

# List of Tables

|   |    |
|---|----|
| Table 1.1: A summary of the main results in this dissertation, aiming to<br>design neural networks to approximate functions in several func-<br>tion spaces . . . . . | 2  |
| Table 4.1: Key ideas of approximating a smooth function . . . . .   | 97 |

This page is intentionally left blank.

# List of Figures

|  |    |
|--|----|
| Figure 1.1: A sketch of most existing results and new results in this dissertation . . . . .   | 3  |
| Figure 2.1: Two examples of trifling regions . . . . .   | 13 |
| Figure 2.2: An example of a ReLU network with width 4 and depth 3 . . .  | 16 |
| Figure 2.3: A detailed example of a ReLU network with two inputs $x_1, x_2$ and an output $\phi(x_1, x_2)$ . . . . .   | 17 |
| Figure 2.4: An illustration of the implementation of the identity map by a ReLU network . . . . .  | 19 |
| Figure 2.5: Illustrations of $\Omega([0, 1]^d, K, \delta)$ , $Q_\beta$ , and $\mathbf{x}_\beta$ for any $\beta \in \{0, 1, \dots, K-1\}^d$ . . . . .   | 21 |
| Figure 2.6: An illustration of the main ideas of constructing $\phi = \phi_2 \circ \Phi_1$ .   | 21 |
| Figure 2.7: An example of a step function for the case $K = 4$ and $d = 1$ .   | 22 |
| Figure 2.8: Illustrations of $Q_\beta$ and $\mathbf{x}_\beta$ for any $\beta \in \{0, 1, \dots, K-1\}^d$ . . .   | 24 |
| Figure 3.1: An illustration of the main idea of proving Theorem 3.1 . . . .  | 26 |
| Figure 3.2: An illustration of the desired network architecture for proving Theorem 3.1 . . . . .  | 27 |
| Figure 3.3: An illustration of $\mathbf{g} = (g_0, g_1^+, g_1^-, \dots, g_n^+, g_n^-)$ and $\tilde{\mathbf{g}} = \sigma(\mathbf{g}) = (\tilde{g}_0, \tilde{g}_1^+, \tilde{g}_1^-, \dots, \tilde{g}_n^+, \tilde{g}_n^-)$ for the case $m = n = 2$ . . . . . | 30 |
| Figure 3.4: Illustrations of the proof of Theorem 3.2 . . . . .  | 34 |
| Figure 3.5: An illustration of $g_k^+$ and $\tilde{g}_k^+ = \sigma(g_k^+)$ . . . . .   | 36 |
| Figure 3.6: An illustration of the network architecture implementing the desired function $\phi$ based on Equation (3.4) . . . . .   | 43 |

|   |    |
|---|----|
| Figure 3.7: An illustration of the ReLU network architecture for proving<br>Lemma 3.6 . . . . .   | 47 |
| Figure 3.8: An illustration of the desired network architecture for proving<br>Theorem 3.5 . . . . .  | 48 |
| Figure 3.9: An illustration of the desired network architecture implement-<br>ing $\max(x_1, x_2, x_3)$ . . . . .   | 50 |
| Figure 3.10: An illustration of $Q_{k,i}$ for $i = 1, 2, 3, 4$ . . . . .  | 52 |
| Figure 3.11: Illustrations of $E_\ell$ for $\ell = 0, 1, 2$ when $K = 4$ and $d = 2$ . . . . .  | 56 |
| Figure 3.12: An example of a step function for the case $K = 4$ and $d = 1$ . . . . .   | 60 |
| Figure 3.13: An illustration of the desired network architecture for proving<br>Theorem 3.12 . . . . .  | 62 |
| Figure 4.1: Examples of “sawtooth” functions $T_1, T_2, T_3$ , and $T_4$ . . . . .  | 67 |
| Figure 4.2: Illustrations of $f_1, f_2$ , and $f_3$ for approximating $x^2$ . . . . .   | 68 |
| Figure 4.3: Illustrations of $f_1 - f_2$ and $f_2 - f_3$ . . . . .  | 68 |
| Figure 4.4: An illustration of the target network architecture for approxi-<br>mating $x^2$ on $x \in [0, 1]$ . . . . .   | 69 |
| Figure 4.5: An illustration of the network architecture implementing $\phi$ for<br>approximating $xy$ on $[0, 1]^2$ . . . . .                                       | 71 |
| Figure 4.6: An illustration of the network architecture implementing $\phi$ for<br>approximating $xy$ on $[a, b]^2$ . . . . .                                       | 72 |
| Figure 4.7: Illustrations of $\Omega([0, 1]^d, K, \delta)$ , $Q_\beta$ , and $\mathbf{x}_\beta$ for any $\beta \in \{0, 1, \dots, K-1\}^d$ . . . . .                | 81 |
| Figure 4.8: An illustration of $\mathcal{A}_1, \mathcal{A}_2, \{1\}$ , and $g$ for the case $d = 2$ and<br>$K = 4$ . . . . .  | 84 |
| Figure 4.9: Illustrations of two sub-network architectures for implementing<br>the desired function $\phi = \phi_2 \circ \phi_1$ based on Equation (4.15) . . . . . | 88 |
| Figure 4.10: An illustration of the network architecture implementing the<br>desired function $\phi$ based Equation (4.17) . . . . .                                | 90 |

|  |     |
|--|-----|
| Figure 4.11: Illustrations of $\Omega([0, 1]^d, K, \delta)$ , $Q_\beta$ , and $\mathbf{x}_\beta$ for any $\beta \in \{0, 1, \dots, K-1\}^d$ . . . . .  | 95  |
| Figure 4.12: An illustration of the sub-network architecture implementing $\varphi\left(\frac{\phi_\alpha(\Psi(\mathbf{x}))}{\alpha!}, P_\alpha(\mathbf{x} - \Psi(\mathbf{x}))\right)$ for each $\alpha \in \mathbb{N}^d$ with $\ \alpha\  \leq s-1$ . . . . . | 103 |
| Figure 4.13: An illustration of the network architecture implementing $\tilde{\phi}(i) = \sum_{j=1}^J 2^{-j} \phi_j(i)$ . . . . .  | 105 |
| Figure 4.14: An illustration of the network architecture for proving Proposition 4.14 . . . . .  | 106 |
| Figure 5.1: An example of a Floor-ReLU network with width 5 and depth 2 . . . . .  | 118 |
| Figure 5.2: Illustrations of $Q_\beta$ and $\mathbf{x}_\beta$ for any $\beta \in \{0, 1, \dots, K-1\}^d$ . . . . .   | 124 |
| Figure 5.3: An illustration of $\phi_1$ on $[0, 1]$ for the case $K = 4$ . . . . .   | 125 |
| Figure 5.4: An illustration of the desired network architecture implementing $\phi_2$ . . . . .  | 127 |
| Figure 5.5: An illustration of the network architecture implementing $\tilde{\phi} = \phi_2 \circ \Phi_1$ . . . . .  | 128 |
| Figure 5.6: An illustration of $g(x) = \sigma(\sigma(x) - \sigma(\frac{x+\delta-1}{\delta}))$ . . . . .  | 129 |
| Figure 5.7: An illustration of the desired network architecture for proving Lemma 5.7 . . . . .  | 130 |
| Figure 5.8: An illustration of the Floor-ReLU network architecture implementing $\phi_{k+1}$ based on Equation (5.6), (5.7), and (5.8). . . . .  | 132 |

This page is intentionally left blank.



# Introduction

Deep neural networks have made significant impacts in many fields of computer science and engineering, especially for large-scale and high-dimensional learning problems. Well-designed neural network architectures, efficient training algorithms, and high-performance computing technologies have made neural-network-based methods very successful in a great number of real applications. Especially in supervised learning, *e.g.*, image classification and objective detection, the great advantages of neural-network-based methods have been demonstrated over traditional learning methods. Understanding the approximation capacity of deep neural networks has become a key question for revealing the power of deep learning. A large number of experiments in real applications have shown the large capacity of deep neural networks from many empirical perspectives, drawing a great deal of attention to the theoretical foundation of the approximation theory of deep neural networks.

In particular, there are three main directions in the error analysis of the approximation theory of neural networks: the **approximation error** estimate, the **optimization error** estimate, and the **generalization error** estimate. See [38, 54] for the introduction of these three error estimates. This dissertation concentrates on the approximation error estimate for neural networks. To this end, we need to solve three fundamental problems listed below.

**Problem 1:** How do we construct a neural network to approximate a function in a

given space?

**Problem 2:** Is there an error estimate for the approximation in Problem 1 in terms of the size of networks, characterized by either the number of parameters or the width and depth simultaneously?

**Problem 3:** If an error estimate exists in Problem 2, is this error estimate (nearly) optimal for the given function space?

This dissertation solves these three problems for several function spaces. See Table 1.1 for a summary of the main results in this dissertation, focusing on designing neural networks to approximate functions in several given function spaces.

## 1.1 Contributions

The main contribution of this dissertation is to provide (nearly optimal) approximation error estimates in terms of the width and depth when constructing neural networks to uniformly approximate polynomials, Hölder continuous functions of order  $\alpha \in (0, 1]$  with a Hölder constant  $\lambda > 0$  ( $\text{Hölder}([0, 1]^d, \alpha, \lambda)$ ), continuous functions ( $C([0, 1]^d)$ ), and smooth functions ( $C^s([0, 1]^d)$ ) on  $[0, 1]^d$ . See Table 1.1 for a summary. Note that all approximation errors in Table 1.1 hold for **arbitrary**  $N, L \in \mathbb{N}^+$  and on  $[0, 1]^d$  **uniformly**. All constants in  $\mathcal{O}(\cdot)$  are **explicitly** estimated in this dissertation, and  $\omega_f(\cdot)$  is the modulus of continuity of  $f$  defined by  $\omega_f(r) = \sup\{|f(\mathbf{x}) - f(\mathbf{y})| : \|\mathbf{x} - \mathbf{y}\|_2 \leq r, \mathbf{x}, \mathbf{y} \in [0, 1]^d\}$ .

Table 1.1: A summary of the main results in this dissertation, aiming to design neural networks to approximate functions in several function spaces.

|               | target function  | activation function | width                     | depth (#hidden-layer)     | approximation error  | optimality  |
|---------------|--|---------------------|---------------------------|---------------------------|--|---|
| Lemma 4.2     | $f(x) = x^2$   | ReLU                | $3N$                      | $L$                       | $N^{-L}$   |   |
| Theorem 4.1   | polynomial $f(\mathbf{x}) = x_1^{\alpha_1} x_2^{\alpha_2} \cdots x_d^{\alpha_d}$ | ReLU                | $\mathcal{O}(N)$          | $\mathcal{O}(L)$          | $\mathcal{O}(N^{-L})$  |   |
| Corollary 4.7 | $f \in \text{Hölder}([0, 1]^d, \alpha, \lambda)$                                 | ReLU                | $\mathcal{O}(N)$          | $\mathcal{O}(L)$          | $\mathcal{O}(\lambda N^{-2\alpha/d} L^{-2\alpha/d})$                   | nearly optimal in $N$ and $L$ , see Section 4.4.1 |
| Theorem 4.6   | $f \in C([0, 1]^d)$  | ReLU                | $\mathcal{O}(N)$          | $\mathcal{O}(L)$          | $\mathcal{O}(\omega_f(N^{-2/d} L^{-2/d}))$                             |   |
| Theorem 4.11  | $f \in C^s([0, 1]^d)$ , $s \in \mathbb{N}^+$                                     | ReLU                | $\mathcal{O}(N \ln(N+1))$ | $\mathcal{O}(L \ln(L+1))$ | $\mathcal{O}(\ f\ _{C^s} N^{-2s/d} L^{-2s/d})$                         | nearly optimal in $N$ and $L$ , see Section 4.4.2 |
| Corollary 5.3 | $f \in \text{Hölder}([0, 1]^d, \alpha, \lambda)$                                 | Floor and ReLU      | $\max\{d, 5N+13\}$        | $64dL+3$                  | $3\lambda d^{d/2} N^{-\alpha\sqrt{L}}$                                 |   |
| Theorem 5.1   | $f \in C([0, 1]^d)$  | Floor and ReLU      | $\max\{d, 5N+13\}$        | $64dL+3$                  | $\omega_f(\sqrt{d} N^{-\sqrt{L}}) + 2\omega_f(\sqrt{d}) N^{-\sqrt{L}}$ |   |

We would like to point out that most results in Table 1.1 can be generalized from  $[0, 1]^d$  to any compact set  $E \subseteq \mathbb{R}^d$ . Such a generalization is mainly based on two key ideas: 1) an affine linear map  $\mathcal{L}_{a,b}(\mathbf{x}) = (b - a)\mathbf{x} + a$  with proper  $a, b \in \mathbb{R}$  satisfying  $E \subseteq [a, b]^d$ ; 2) the function extension (*e.g.*, see Lemma 4.2 of [53] for the extension of continuous functions).

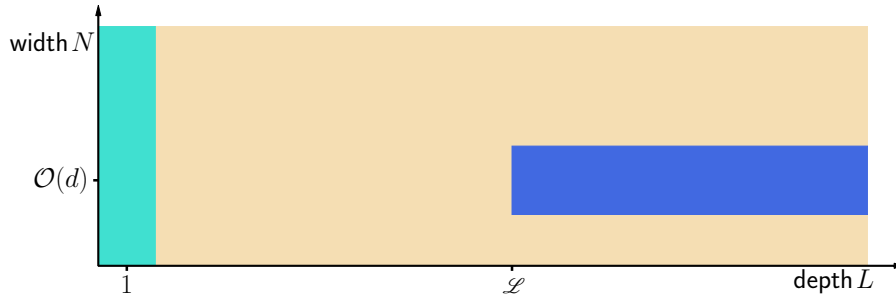


Figure 1.1: A sketch of most existing results and new results in this dissertation.  $\mathcal{L}$  represents a sufficiently large unknown number. Most existing results (*e.g.*, [18, 23, 35, 39, 57, 59, 60]) are applicable in the areas in  or , while our results are suitable for almost all areas characterized by .

As far as we know, most existing works focus on either one-hidden-layer networks (visualized by the region in  in Figure 1.1), or very deep networks with a constant width (visualized by the region in  in Figure 1.1). Meanwhile, these works only provide asymptotic<sup>①</sup> approximation errors in terms of the number of parameters, which are valid for particular network architectures. They are unable to give approximation error estimates for other network architectures with the same number of parameters. To overcome this, we provide general approximation error characterizations with explicit formulas for the prefactors, in terms of the width and depth simultaneously (visualized by the region in  in Figure 1.1), which is of more practical interest in real applications and requires innovative proofs. This gives us much more freedom to design neural networks for a good approximation and we can always give an error estimate via the width and depth no matter what network architecture is given, though the error estimate may not be optimal for

<sup>①</sup>“Asymptotic” means the approximation error is described via big O notation  $\mathcal{O}(\cdot)$  without an explicit formula for the prefactor.

unusual architectures. In fact, many results in previous works can be regarded as the corollaries of this dissertation.

Problem 1 and 2 are completely solved by providing approximation error estimates in terms of the width and depth simultaneously. For Problem 3, we use VC-dimension to show our approximation error estimates are nearly optimal for the Hölder continuous function space ( $\text{Hölder}([0, 1]^d, \alpha, \lambda)$ ) and the smooth function space ( $C^s([0, 1]^d)$ ). The optimality becomes insignificant if (nearly) exponential approximation errors are attained.

Table 1.1 (Theorem 4.1) shows that ReLU networks with width  $\mathcal{O}(N)$  and depth  $\mathcal{O}(L)$  are able to approximate  $d$ -dimensional polynomials on  $[0, 1]^d$  within an error  $\mathcal{O}(N^{-L})$ . This reveals the power of depth in ReLU networks for approximating polynomials, from function compositions. Generally speaking, such an approximation error is the best (up to constants) what we can expect since ReLU networks with width  $\mathcal{O}(N)$  and depth  $\mathcal{O}(L)$  are continuous piecewise linear functions with at most  $\mathcal{O}(N)^{\mathcal{O}(L)}$  linear pieces. The starting point of a good approximation of functions is to approximate polynomials with high accuracy. In classical approximation theory, the approximation power of a lot of numerical schemes depends on the degree of polynomials that can be locally reproduced. Being able to approximate polynomials with an exponential error plays a vital role in the approximation power of deep ReLU networks. It is interesting to study whether there are any other function spaces with a reasonable size, besides the polynomial space, having an exponential error when approximated by neural networks.

In particular, we introduce new networks built with either Floor ( $\lfloor x \rfloor$ ) or ReLU ( $\max\{0, x\}$ ) as the activation function in each neuron. We call such networks Floor-ReLU networks. It is proved by construction that nearly exponential approximation errors can be attained when using Floor-ReLU networks with fixed architectures to approximate Hölder continuous functions and general continuous functions on  $[0, 1]^d$ . As shown in Table 1.1, approximation errors are improved from polynomial ones to nearly exponential ones by adding a simple activation function (Floor) to ReLU

networks. This reveals the power of deep Floor-ReLU network architectures. As we shall see later, the idea of function compositions is the most significant cornerstone of the proofs for the results listed in Table 1.1. Finally, we would like to remark that the architecture of the final Floor-ReLU network is independent of the target function  $f$ . That is, only the values of the parameters rely on the target function  $f$ . In particular, the choice of activation functions (Floor or ReLU) in each neuron is independent of the target function  $f$ .

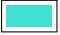

## 1.2 Related work

This dissertation is a summary of our previous papers [38, 52, 53, 54], focusing on the approximation error estimate for neural networks. Thus, all the contents of this dissertation focus on three main problems, Problem 1, 2, and 3. In the following, only the previous works related to them are reviewed.

The approximation theory of neural networks has been an active research topic in the past few decades. Previously, as a special kind of ridge function approximation, shallow neural networks with one hidden layer and various activation functions (*e.g.*, wavelets pursuits [12, 41], adaptive splines [19, 49], radial basis functions [10, 18, 23, 47, 57], sigmoid functions [8, 13, 14, 15, 24, 32, 33, 37, 40]) were widely discussed and admit good approximation properties, *e.g.*, the universal approximation property [16, 24, 25], overcoming the curse of dimensionality [3], and providing attractive approximation errors in nonlinear approximation [12, 18, 19, 23, 41, 49, 57].

The introduction of deep neural networks with more than one hidden layers has made significant impacts in many fields in computer science and engineering including computer vision [31] and natural language processing [1]. New scientific computing tools based on deep networks have also emerged and facilitated large-scale and high-dimensional problems that were impractical previously [20, 22]. The design of deep ReLU networks and high-performance computing technologies are the key of such a revolution. These breakthroughs have stimulated broad research topics from

different points of views to study the power of deep ReLU networks, *e.g.*, in terms of combinatorics [44], topology [7], Vapnik-Chervonenkis (VC) dimension [4, 5, 51], fat-shattering dimension [2, 29], information theory [48], classical approximation theory [3, 16, 25, 38, 52, 53, 54, 59], optimization [27, 28, 45], *etc.*

Particularly in approximation theory, non-quantitative and asymptotic approximation errors of ReLU networks have been proposed for various types of functions. For example, smooth functions [21, 34, 38, 39, 58], piecewise smooth functions [48], band-limited functions [43], continuous functions [53, 59]. However, to the best of our knowledge, existing theories [17, 21, 34, 39, 42, 43, 48, 55, 58, 59] can only provide implicit formulas. In other words, the approximation error contains an unknown prefactor, or they work only for sufficiently large network size. For example, an approximation error  $c_d L^{-2\alpha/d}$  for Lipschitz continuous functions on  $[0, 1]^d$  is estimated in [59] via a narrow and deep ReLU network with  $L$  hidden layers, where  $c_d$  is an unknown number depending on  $d$ . For another example, the existence of a ReLU network with a constant width and  $W$  parameters is constructed in [60] to approximate a smooth function in  $C^s([0, 1]^d)$  within an error  $c_{s,d} W^{-2s/d} (\ln W)^{2s/d}$ , where  $c_{s,d}$  is still an unknown number depending on  $s$  and  $d$ . Generally, most of these works can be divided into two cases: 1) networks with varying width and only one hidden layer [18, 23, 35, 57] (visualized by the region in  in Figure 1.1); 2) networks with a fixed width of  $\mathcal{O}(d)$  and a varying depth larger than an unknown number  $\mathcal{L}$  [39, 59, 60] (represented by the region in  in Figure 1.1).

Almost all works listed above answer Problem 1 and 2 for given activation functions and special function spaces. Most of them estimate the approximation error in terms of the number of parameters. In other words, their approximation errors are only valid for very special network architectures, such as very deep but very narrow networks, complicated networks generated by compositing shallow-wide sub-networks and deep-narrow sub-networks, *etc.*, while our approximation error estimates in this dissertation are valid for arbitrary width and depth up to absolute constants. It means the shape of our network architectures is a rectangle

with free choice of width (the maximum width of networks) and length (the depth of networks), which is of more practical interest in real applications and requires innovative proofs.

Finally, let us turn to Problem 3. A typical method characterizing optimality in the approximation theory of neural networks is essentially to study the connection between the approximation error and VC-dimension [38, 52, 53, 58, 59, 60]. Of course, this method relies on the VC-dimension upper bound given in [4]. In this dissertation, we adopt this method with several modifications to simplify the proof. As we shall see later in Section 4.4, the optimality is discussed for two function spaces: 1) the Hölder function space (see Section 4.4.1); 2) the smooth function space (see Section 4.4.2).

This page is intentionally left blank.



# Preliminaries

Before moving to the main body of this dissertation, we first introduce the preliminaries related to this dissertation including notations used throughout this dissertation, the architecture of neural networks, and the general ideas of the approximation by neural networks.

## 2.1 Notations

For convenience, we present all notations used throughout this dissertation in this section. Several notations used only in a particular section are not presented here.

### 2.1.1 Basic notations

Basic notations are listed below.

- Let  $\mathbb{Z}$  and  $\mathbb{R}$  denote the set of integers and real numbers, respectively.
- Let  $\mathbb{N}$  denote the set of natural numbers and  $\mathbb{N}^+$  denote the set containing all positive integers, *i.e.*,  $\mathbb{N} = \{0, 1, 2, \dots\}$  and  $\mathbb{N}^+ = \{1, 2, 3, \dots\}$ .
- Matrices are denoted by bold uppercase letters, *e.g.*,  $\mathbf{A} \in \mathbb{R}^{m \times n}$  is a real matrix of size  $m \times n$ , and  $\mathbf{A}^T$  denotes the transpose of  $\mathbf{A}$ . Correspondingly,  $\mathbf{A}(i, j)$

is the  $(i, j)$ -th entry of  $\mathbf{A}$ ;  $\mathbf{A}(:, j)$  is the  $j$ -th column of  $\mathbf{A}$ ;  $\mathbf{A}(i, :)$  is the  $i$ -th row of  $\mathbf{A}$ .

- Vectors are denoted as bold lowercase letters, *e.g.*,

$$\mathbf{v} = [v_1, \dots, v_d]^T = \begin{bmatrix} v_1 \\ \vdots \\ v_d \end{bmatrix}$$

is a column vector of size  $d$  and  $\mathbf{v}(i)$  is the  $i$ -th element of  $\mathbf{v}$ . For simplicity, a vector  $\mathbf{v} \in \mathbb{R}^d$  can also be denoted by  $\mathbf{v} = (v_1, \dots, v_d)$ .

- By convention, “[” and “]” are used to partition matrices (vectors) into blocks, *e.g.*, a matrix  $\mathbf{A}$  can be partitioned into  $\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}$  and a row vector  $\mathbf{v}$  can be denoted by  $\mathbf{v} = [v_1, v_2, \dots, v_d] \in \mathbb{R}^d$ .
- We say a map (transform)  $\mathcal{L} : \mathbb{R}^m \rightarrow \mathbb{R}^n$  is **affine linear** if there exist  $\mathbf{W} \in \mathbb{R}^{n \times m}$  and  $\mathbf{b} \in \mathbb{R}^n$  such that  $\mathcal{L}(\mathbf{x}) = \mathbf{W} \cdot \mathbf{x} + \mathbf{b}$  for any  $\mathbf{x} \in \mathbb{R}^m$ . In particular, an affine linear map is also called a **linear** function in the case  $n = 1$ .
- For a real number  $p \in [1, \infty)$ , the  $p$ -norm (or  $\ell^p$ -norm) of a vector  $\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$  is defined by

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

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

$$* \|\boldsymbol{\alpha}\|_1 = |\alpha_1| + |\alpha_2| + \dots + |\alpha_d|;$$

$$* \mathbf{x}^{\boldsymbol{\alpha}} = x_1^{\alpha_1} x_2^{\alpha_2} \dots x_d^{\alpha_d}, \text{ where } \mathbf{x} = (x_1, x_2, \dots, x_d);$$

$$* \boldsymbol{\alpha}! = \alpha_1! \alpha_2! \dots \alpha_d!;$$

$$* \partial^{\boldsymbol{\alpha}} = \frac{\partial^{\alpha_1}}{\partial x_1^{\alpha_1}} \frac{\partial^{\alpha_2}}{\partial x_2^{\alpha_2}} \dots \frac{\partial^{\alpha_d}}{\partial x_d^{\alpha_d}}.$$

- Let  $\mathcal{O}(\cdot)$  denote the big O notation. That is, for any  $\mathbf{n} \in \mathbb{N}^d$  and functions  $f$  and  $g$ ,  $f(\mathbf{n}) = \mathcal{O}(g(\mathbf{n}))$  means that there exist  $C > 0$  and  $\mathbf{n}_0 \in \mathbb{N}^d$  independent of  $\mathbf{n}$ ,  $f$ , and  $g$  such that  $f(\mathbf{n}) \leq Cg(\mathbf{n})$  when  $\mathbf{n}(i) \geq \mathbf{n}_0(i)$  for all  $i$ .
- The floor function (Floor) is defined as  $\lfloor x \rfloor := \max\{n : n \leq x, n \in \mathbb{Z}\}$  for any  $x \in \mathbb{R}$ .  $\lfloor \mathbf{x} \rfloor$  means applying  $\lfloor \cdot \rfloor$  elementwisely to  $\mathbf{x}$ . Similarly, the ceiling function (Ceiling) is defined as  $\lceil x \rceil := \min\{n : n \geq x, n \in \mathbb{Z}\}$  for any  $x \in \mathbb{R}$ .
- Similar to “min” and “max”, let  $\text{mid}(x_1, x_2, x_3)$  denote the middle value of **three** inputs  $x_1$ ,  $x_2$ , and  $x_3$ .<sup>①</sup>
- 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  $\theta$ , *i.e.*,  $\text{bin}0.\theta_1\theta_2 \cdots \theta_L := \sum_{\ell=1}^L \theta_{\ell} 2^{-\ell}$ .
- Let  $\text{H\"older}([0, 1]^d, \alpha, \lambda)$  denote the space of H\"older continuous functions of order  $\alpha \in (0, 1]$  on  $[0, 1]^d$  with a H\"older constant  $\lambda > 0$ . To be precise, each function  $f$  of  $\text{H\"older}([0, 1]^d, \alpha, \lambda)$  satisfies

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

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

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

Generally,  $E$  is assigned as  $[0, 1]^d$  in this dissertation. In particular, the closed

---

<sup>①</sup>Note that “mid” can be defined via  $\text{mid}(x_1, x_2, x_3) = x_1 + x_2 + x_3 - \max(x_1, x_2, x_3) - \min(x_1, x_2, x_3)$ , which can be implemented by a ReLU network with width 14 and depth 2, as shown in Lemma 3.8.

unit ball of  $C^s([0, 1]^d)$  is denoted by

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

- The modulus of continuity of a continuous function  $f \in C([0, 1]^d)$  is defined by

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

Clearly,  $\omega_f(nr) \leq n\omega_f(r)$  for any  $n \in \mathbb{N}^+$  and  $r \geq 0$ .

### 2.1.2 Set notations

All set notations used in this dissertation can be found below.

- The Lebesgue measure of a measurable set  $S \in \mathbb{R}^d$  is denoted by  $\mu(S)$ .
- Let  $|S|$  denote the size of a finite set  $S$ , *i.e.*, the number of all elements in  $S$ .
- The set difference of two sets  $A$  and  $B$  is denoted by  $A \setminus B := \{x : x \in A, x \notin B\}$ .
- For a set of numbers  $A$  and a real number  $x$ ,  $A - x := \{y - x : y \in A\}$ .
- Let  $1_S$  be the characteristic function on a set  $S$ , *i.e.*,  $1_S$  is equal to 1 on  $S$  and 0 outside  $S$ .  $S$  can be simply described by one or more conditions, *e.g.*,  $1_{\{n \leq m\}}$  is equal to 1 if  $n \leq m$  and 0 if  $n > m$ .
- Let  $\mathcal{B}(\mathbf{x}, r) \subseteq \mathbb{R}^d$  denote the closed ball, in  $\ell^2$ -norm, with a center  $\mathbf{x} \subseteq \mathbb{R}^d$  and a radius  $r$ , *i.e.*,

$$\mathcal{B}(\mathbf{x}, r) := \{\mathbf{y} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{y}\|_2 \leq r\}.$$

- 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

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

In particular,  $\Omega([0, 1]^d, K, \delta) = \emptyset$  if  $K = 1$ . See Figure 2.1 for two examples of trifling regions.

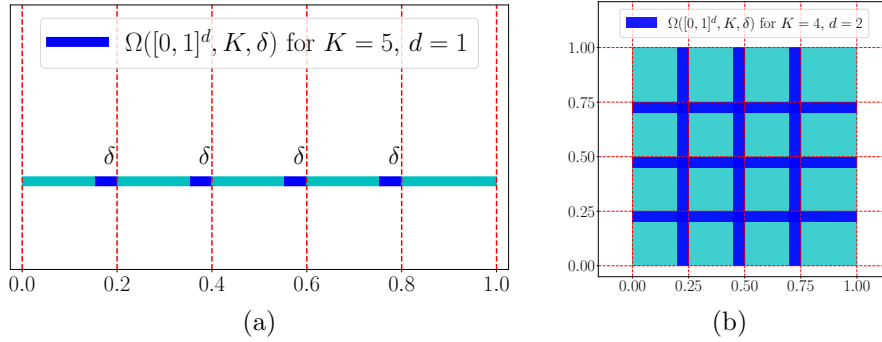


Figure 2.1: Two examples of trifling regions. (a)  $K = 5, d = 1$ . (b)  $K = 4, d = 2$ .

### 2.1.3 Neural network notations

We list neural network notations as follows.

- 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  $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$ .
- The expression “a network (architecture) with width  $N$  and depth  $L$ ” means
  - \* The maximum width of this network (architecture) for all **hidden** layers is no more than  $N$ .
  - \* The number of **hidden** layers of this network (architecture) is no more than  $L$ .
- The expression “a (vector-valued) function is implemented by a network (architecture)” means, by specifying the parameters as proper real numbers, this

network (architecture) has the same output as this function for each input.

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

- \*  $\phi$  is a two-dimensional function that maps from  $\mathbb{R}^2$  to  $\mathbb{R}$ .
- \*  $\phi$  can be implemented by a two-hidden-layer ReLU network that the number of neurons in each hidden layer is 100.

We would like to point out that each element of  $\mathcal{NN}(c_1; c_2; \dots; c_m)$  is a continuous piecewise linear function.

## 2.2 Architecture of neural networks

There are a large number of types of neural network architectures, *e.g.*, convolution neural networks (CNN), recurrent neural networks (RNN), and generative adversarial networks (GAN), variational auto encoders (VAE), residual networks (ResNet), *etc.* This dissertation focuses on feed-forward fully connected neural networks. If there are no special instructions, “feed-forward fully connected neural network(s)” is abbreviated to “network(s)” throughout this dissertation. In this section, we will describe the architecture of networks mathematically and intuitively in

Section 2.2.1 and study the compositions and combinations of network architectures in Section 2.2.2.

### 2.2.1 Descriptions

First, we use mathematical formulas to describe network architectures. Assume  $\varrho_1, \dots, \varrho_r$  are one-dimensional functions. Let  $N_0 = d$ ,  $N_{L+1} \in \mathbb{N}^+$ , and  $N_\ell$  be the number of neurons in  $\ell$ -th hidden layer of a network with activation functions  $\varrho_1, \dots, \varrho_r$  for  $\ell = 1, 2, \dots, L$ , then the architecture of this network with input  $\mathbf{x}$  and output  $\phi(\mathbf{x})$  can be described as

$$\mathbf{x} = \tilde{\mathbf{h}}_0 \xrightarrow{\mathbf{W}_0, \mathbf{b}_0} \mathbf{h}_1 \xrightarrow{\varrho_1, \dots, \varrho_r} \tilde{\mathbf{h}}_1 \quad \dots \quad \xrightarrow{\mathbf{W}_{L-1}, \mathbf{b}_{L-1}} \mathbf{h}_L \xrightarrow{\varrho_1, \dots, \varrho_r} \tilde{\mathbf{h}}_L \xrightarrow{\mathbf{W}_L, \mathbf{b}_L} \mathbf{h}_{L+1} = \phi(\mathbf{x}), \quad (2.2)$$

where  $\mathbf{W}_\ell \in \mathbb{R}^{N_{\ell+1} \times N_\ell}$ ,  $\mathbf{b}_\ell \in \mathbb{R}^{N_{\ell+1}}$ ,

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

and

$$\tilde{\mathbf{h}}_{\ell,n} \in \{\varrho_1(\mathbf{h}_{\ell,n}), \dots, \varrho_r(\mathbf{h}_{\ell,n})\}, \quad \text{for } \ell = 1, 2, \dots, L \text{ and } n = 1, 2, \dots, N_\ell,$$

where  $\mathbf{h}_\ell = (\mathbf{h}_{\ell,1}, \dots, \mathbf{h}_{\ell,N_\ell})$ ,  $\tilde{\mathbf{h}}_\ell = (\tilde{\mathbf{h}}_{\ell,1}, \dots, \tilde{\mathbf{h}}_{\ell,N_\ell})$  for each  $\ell = 1, 2, \dots, L$ , and  $\mathcal{L}_\ell$  is an affine linear map given by  $\mathcal{L}_\ell(\mathbf{z}) := \mathbf{W}_\ell \cdot \mathbf{z} + \mathbf{b}_\ell$  for each  $\ell = 0, 1, 2, \dots, L$ .

The most common type of activation function is the rectified linear unit (ReLU), denoted by  $\sigma$  in this dissertation. We remark that using the ReLU activation function is not much different from using any other **continuous piecewise linear** activation function with finitely many linear pieces. In fact, if we let  $\tilde{\sigma}$  be a continuous piecewise linear activation function with finitely many linear pieces, then we can always replace a network, using one of  $\{\sigma, \tilde{\sigma}\}$  as activation function, by another network having the other activation function in  $\{\sigma, \tilde{\sigma}\}$  while only increasing the width and depth by absolute constant factors.

The networks with only ReLU activation function, *i.e.*,  $\varrho_1 = \cdots = \varrho_r = \sigma$  in Equation (2.2), are called ReLU networks. In this case, the set of functions implemented by the architecture in Equation (2.2) is exactly  $\mathcal{NN}(\#input = d; \text{widthvec} = [N_1, N_2, \dots, N_L])$ . Moreover, the (vector-valued) function  $\phi$  implemented by the network in the Equation (2.2) can also be represented in a compositive manner by

$$\phi = \mathcal{L}_L \circ \sigma \circ \mathcal{L}_{L-1} \circ \sigma \circ \cdots \circ \mathcal{L}_2 \circ \sigma \circ \mathcal{L}_1 \circ \sigma \circ \mathcal{L}_0.$$

In particular, if  $r = 2$ ,  $\varrho_1 = \sigma$ , and  $\varrho_2(x) = \lfloor x \rfloor$  for any  $x \in \mathbb{R}$ , the network described by Equation (2.2) is a Floor-ReLU network. We will discuss more details of Floor-ReLU networks in Chapter 5.

To visualize the network architecture, we take ReLU networks as examples. Figure 2.2 provides an example of a ReLU network with width 4 and depth 3. Note that the affine linear transform and the activation function are contained in a single neuron in Figure 2.2. To make the architecture of a ReLU network more clear, we put the affine linear transform and the activation function into different neurons in another example shown in Figure 2.3.

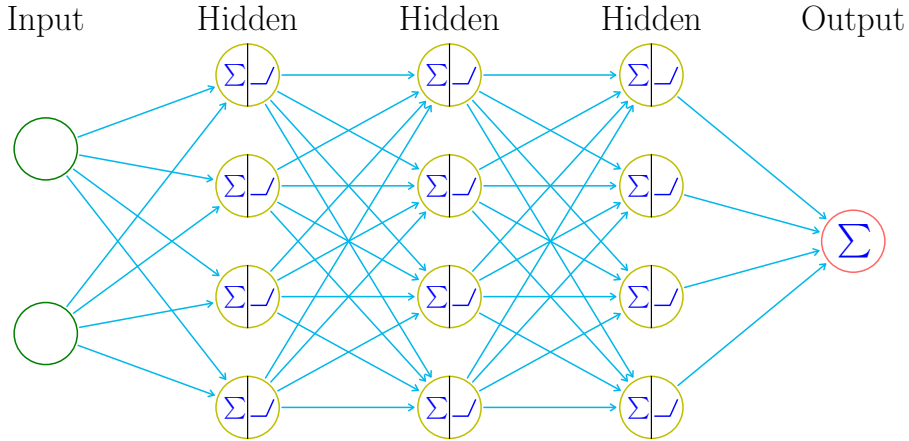


Figure 2.2: An example of a ReLU network with width 4 and depth 3. This network has two neurons in the input layer, one neuron in the output layer, and four neurons in each hidden layer.



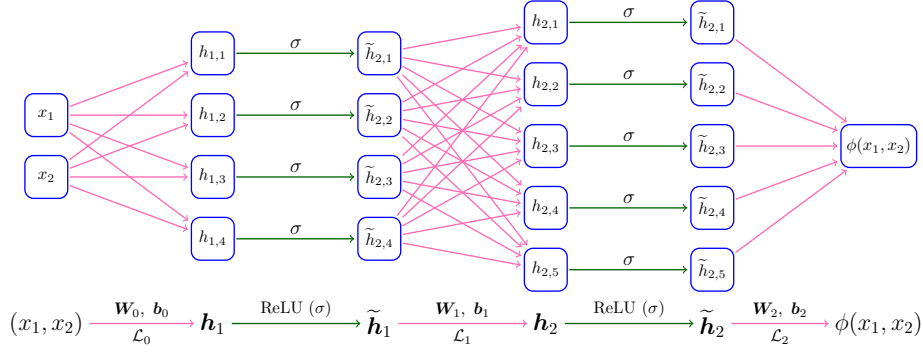


Figure 2.3: A detailed example of a ReLU network with two inputs  $x_1, x_2$  and an output  $\phi(x_1, x_2)$ . Here,  $\mathbf{h}_1 = (h_{1,1}, h_{1,2}, h_{1,3}, h_{1,4})$ ,  $\mathbf{h}_2 = (h_{2,1}, h_{2,2}, h_{2,3}, h_{2,4}, h_{2,5})$ ,  $\tilde{\mathbf{h}}_1 = \sigma(\mathbf{h}_1) = (\tilde{h}_{1,1}, \tilde{h}_{1,2}, \tilde{h}_{1,3}, \tilde{h}_{1,4})$ , and  $\tilde{\mathbf{h}}_2 = \sigma(\mathbf{h}_2) = (\tilde{h}_{2,1}, \tilde{h}_{2,2}, \tilde{h}_{2,3}, \tilde{h}_{2,4}, \tilde{h}_{2,5})$ .

### 2.2.2 Compositions and combinations

We use a lemma below to describe the compositions and combinations of ReLU network architectures.

**Lemma 2.1.** *The following three statements hold.*

- (i) For any  $N, L, d_1, d_2, d_3, d_4 \in \mathbb{N}^+$ , assume that  $\mathcal{L}_1 : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2}$  and  $\mathcal{L}_2 : \mathbb{R}^{d_3} \rightarrow \mathbb{R}^{d_4}$  are two affine linear maps, and  $\Phi \in \mathcal{NN}(\#input = d_2; \text{width} \leq N; \text{depth} \leq L; \#output = d_3)$ . Then

$$\Phi \circ \mathcal{L}_1 \in \mathcal{NN}(\#input = d_1; \text{width} \leq N; \text{depth} \leq L; \#output = d_3)$$

and

$$\mathcal{L}_2 \circ \Phi \in \mathcal{NN}(\#input = d_2; \text{width} \leq N; \text{depth} \leq L; \#output = d_4).$$

- (ii) For any  $N_1, N_2, L_1, L_2, d_1, d_2, d_3 \in \mathbb{N}^+$ , if  $\Phi_1 \in \mathcal{NN}(\#input = d_1; \text{width} \leq N_1; \text{depth} \leq L_1; \#output = d_2)$  and  $\Phi_2 \in \mathcal{NN}(\#input = d_2; \text{width} \leq N_2; \text{depth} \leq L_2; \#output = d_3)$ , then  $\Phi_2 \circ \Phi_1$  is in

$$\mathcal{NN}(\#input = d_1; \text{width} \leq \max\{N_1, N_2\}; \text{depth} \leq L_1 + L_2; \#output = d_3).$$

(iii) For any  $N_1, N_2, L_1, L_2, d \in \mathbb{N}^+$  and  $a, b, c \in \mathbb{R}$  with  $N_1 \geq 2$  and  $N_2 \geq 2$ , if  $\phi_1 \in \mathcal{NN}(\#input = d; \text{width} \leq N_1; \text{depth} \leq L_1; \#output = 1)$  and  $\phi_2 \in \mathcal{NN}(\#input = d; \text{width} \leq N_2; \text{depth} \leq L_2; \#output = 1)$ , then  $a\phi_1 + b\phi_2 + c$  is in

$$\mathcal{NN}(\#input = d; \text{width} \leq N_1 + N_2; \text{depth} \leq \max\{L_1, L_2\}; \#output = 1)$$

*Proof.* Let first prove Part (i). The case  $L = 1$  is trivial, we consider  $L \geq 2$  below. Since  $\Phi \in \mathcal{NN}(\text{width} \leq N; \text{depth} \leq L)$ , there exist two affine linear maps  $\widehat{\mathcal{L}}_1, \widehat{\mathcal{L}}_2$  and  $\Psi_1, \Psi_2 \in \mathcal{NN}(\text{width} \leq N; \text{depth} \leq L - 1)$  such that

$$\Phi = \Psi_1 \circ \sigma \circ \widehat{\mathcal{L}}_1, \quad \text{and} \quad \Phi = \widehat{\mathcal{L}}_2 \circ \sigma \circ \Psi_2.$$

Therefore,

$$\Phi \circ \mathcal{L}_1 = \Psi_1 \circ \sigma \circ \widehat{\mathcal{L}}_1 \circ \mathcal{L}_1 = \Psi_1 \circ \sigma \circ \widetilde{\mathcal{L}}_1, \quad \text{and} \quad \mathcal{L}_2 \circ \Phi = \mathcal{L}_2 \circ \widehat{\mathcal{L}}_2 \circ \sigma \circ \Psi_2 = \widetilde{\mathcal{L}}_2 \circ \sigma \circ \Psi_2,$$

where  $\widetilde{\mathcal{L}}_1 = \widehat{\mathcal{L}}_1 \circ \mathcal{L}_1$  and  $\widetilde{\mathcal{L}}_2 = \mathcal{L}_2 \circ \widehat{\mathcal{L}}_2$  are two new affine linear maps, implying

$$\Phi \circ \mathcal{L}_1 \in \mathcal{NN}(\#input = d_1; \text{width} \leq N; \text{depth} \leq L; \#output = d_3)$$

and

$$\mathcal{L}_2 \circ \Phi \in \mathcal{NN}(\#input = d_2; \text{width} \leq N; \text{depth} \leq L; \#output = d_4).$$

Next, let us focus on Part (ii). The case  $L_1 = 1$  or  $L_2 = 1$  is trivial, so we assume  $L_1 \geq 2$  and  $L_2 \geq 2$  below. Since  $\Phi_1 \in \mathcal{NN}(\text{width} \leq N_1; \text{depth} \leq L_1)$  and  $\Phi_2 \in \mathcal{NN}(\text{width} \leq N_2; \text{depth} \leq L_2)$ , there exist  $\Psi_1 \in \mathcal{NN}(\text{width} \leq N_1; \text{depth} \leq L_1 - 1)$  and  $\Psi_2 \in \mathcal{NN}(\text{width} \leq N_2; \text{depth} \leq L_2 - 1)$  such that

$$\Phi_1 = \mathcal{L}_1 \circ \sigma \circ \Psi_1, \quad \text{and} \quad \Phi_2 = \Psi_2 \circ \sigma \circ \mathcal{L}_2,$$

where  $\mathcal{L}_1, \mathcal{L}_2$  are two affine linear maps. Therefore,

$$\Phi_2 \circ \Phi_1 = \Psi_2 \circ \sigma \circ \mathcal{L}_2 \circ \mathcal{L}_1 \circ \sigma \circ \Psi_1 = \Psi_2 \circ \sigma \circ \mathcal{L} \circ \sigma \circ \Psi_1,$$

where  $\mathcal{L} = \mathcal{L}_2 \circ \mathcal{L}_1$  is a new affine linear map. Thus,  $\Phi_2 \circ \Phi_1$  can be implemented by a ReLU network with width  $\max\{N_1, N_2\}$  and depth  $(L_1 - 1) + 1 + 1 + (L_2 - 1) = L_1 + L_2$ , implying  $\Phi_2 \circ \Phi_1$  is in

$$\mathcal{NN}(\#input = d_1; \text{ width } \leq \max\{N_1, N_2\}; \text{ depth } \leq L_1 + L_2; \#output = d_3).$$

Finally, let us consider Part (iii). Let  $\mathfrak{t}$  denote the one-dimensional identity map. As shown in Figure 2.4,  $\mathfrak{t}$  can be understood as an implementation of a ReLU network with an arbitrary number of hidden layers and width 2. Thus, for  $j = 1, 2$ ,  $\mathfrak{t} \circ \phi_j$  can

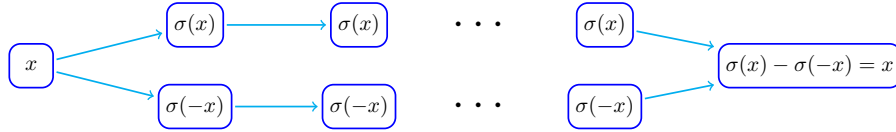


Figure 2.4: An illustration of the implementation of the identity map by a ReLU network based on the fact  $\sigma \circ \sigma = \text{id}$ .

be regarded as an implementation of ReLU network with  $\max\{L_1, L_2\}$  hidden layers and width  $\max\{N_j, 2\} = N_j$ . By placing the two networks implementing  $\mathfrak{t} \circ \phi_1$  and  $\mathfrak{t} \circ \phi_2$  in parallel (share the inputs), we have

$$\Phi \in \mathcal{NN}(\#input = d; \text{ width } \leq N_1 + N_2; \text{ depth } \leq \max\{L_1, L_2\}; \#output = 2),$$

where  $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^2$  is defined by  $\Phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}))$  for any  $\mathbf{x} \in \mathbb{R}^d$ . Define an affine linear map  $\mathcal{L} : \mathbb{R}^2 \rightarrow \mathbb{R}$  via  $\mathcal{L}(x, y) = ax + by + c$ . By Part (i), we have  $a\phi_1 + b\phi_2 + c = \mathcal{L} \circ \Phi \in \mathcal{NN}(\#input = d; \text{ width } \leq N_1 + N_2; \text{ depth } \leq \max\{L_1, L_2\}; \#output = 1)$ . So we finish the proof.  $\square$

## 2.3 General ideas of approximation by networks

In this section, we discuss the general ideas of approximation by networks. Universal approximation theorem shows that one-hidden-layer networks can approximate continuous functions arbitrarily well on  $[0, 1]^d$  as long as the network size is large enough. However, it is non-trivial to characterize the approximation error in terms of the width and depth simultaneously as we will do in later chapters. Thus, let us discuss the general ideas to warm up the later constructions and proofs.

### 2.3.1 ReLU networks

First, let us consider the approximation by ReLU networks. To illustrate the general ideas, we take continuous functions as examples. The ideas of smooth functions are similar by applying Taylor expansion, as we shall see later in Section 4.3. To approximate a continuous function  $f$  on  $[0, 1]^d$ , we essentially construct a piecewise constant function via function compositions. However, piecewise constant functions cannot be implemented by ReLU networks because of their discontinuity. To overcome this, we introduce the trifling region  $\Omega([0, 1]^d, K, \delta)$ , defined in Equation (2.1), and construct ReLU networks to implement **almost** piecewise linear functions to approximate the target functions outside the trifling region. For the sake of clarity, we divide the main ideas into four steps. See Figure 2.6 for an illustration.

1. Normalize  $f$  as  $\tilde{f}$ , partition  $[0, 1]^d$  into a union of sub-cubes<sup>②</sup>  $\{Q_\beta\}_{\beta \in \{0, 1, \dots, K-1\}^d}$  and the trifling region  $\Omega([0, 1]^d, K, \delta)$ , and let  $\mathbf{x}_\beta$  denote the vertex of  $Q_\beta$  with minimum  $\|\cdot\|_1$  norm, where  $K \in \mathbb{N}^+$  and  $\delta \in (0, \frac{1}{3K}]$  are two numbers determined later. See Figure 2.5 for illustrations of  $\Omega([0, 1]^d, K, \delta)$ ,  $Q_\beta$ , and  $\mathbf{x}_\beta$  for any  $\beta \in \{0, 1, \dots, K-1\}^d$ .
2. Construct a sub-network to implement a vector-valued function  $\Phi_1 : \mathbb{R}^d \rightarrow \mathbb{R}^d$  projecting the whole cube  $Q_\beta$  to the  $d$ -dimensional index  $\beta$  for each  $\beta$ , *i.e.*,  $\Phi_1(\mathbf{x}) = \beta$  for all  $\mathbf{x} \in Q_\beta$  and each  $\beta \in \{0, 1, \dots, K-1\}^d$ .

---

<sup>②</sup>For simplicity, we abbreviate ( $d$ -dimensional) hypercube to cube.

3. Construct a sub-network to implement a function  $\phi_2 : \mathbb{R}^d \rightarrow \mathbb{R}$  mapping the index  $\beta$  approximately to  $f(\mathbf{x}_\beta)$  for each  $\beta$ . Then  $\phi_2 \circ \Phi_1(\mathbf{x}) = \phi_2(\beta) \approx \tilde{f}(\mathbf{x}_\beta)$  for any  $\mathbf{x} \in Q_\beta$  and each  $\beta \in \{0, 1, \dots, K-1\}^d$ , implying  $\tilde{\phi} := \phi_2 \circ \Phi_1$  approximate  $\tilde{f}$  within an error  $\mathcal{O}(\omega_f(1/K))$  outside the trifling region.
4. Re-scale and shift  $\tilde{\phi}$  to obtain a function  $\phi$  approximating  $f$  well outside the trifling region. Then modify  $\phi$  to let it approximate  $f$  uniformly well on  $[0, 1]^d$  and determine the network architecture implementing the modified function  $\phi$ .

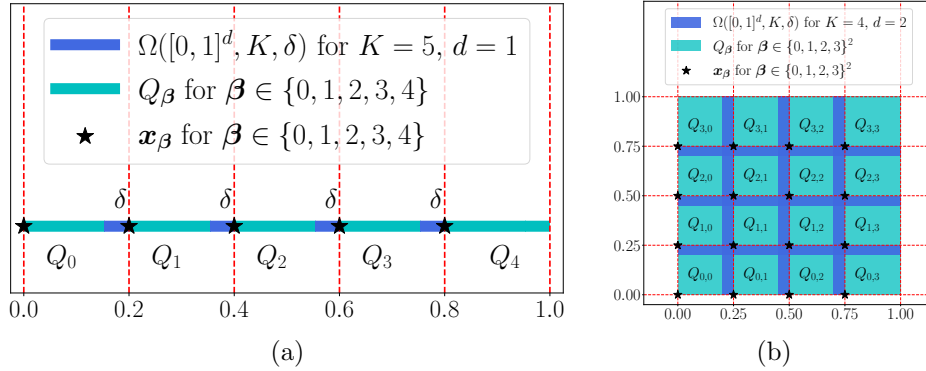


Figure 2.5: Illustrations of  $\Omega([0, 1]^d, K, \delta)$ ,  $Q_\beta$ , and  $\mathbf{x}_\beta$  for any  $\beta \in \{0, 1, \dots, K-1\}^d$ . (a)  $K = 5$  and  $d = 1$ . (b)  $K = 4$  and  $d = 2$ .

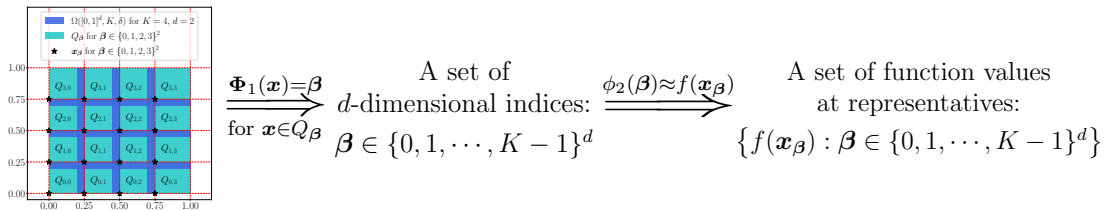


Figure 2.6: An illustration of the main ideas of constructing  $\phi = \phi_2 \circ \Phi_1$ . Note that  $\phi \approx f$  on  $[0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)$ , since  $\phi(\mathbf{x}) = \phi_2 \circ \Phi_1(\mathbf{x}) = \phi_2(\beta) \approx f(\mathbf{x}_\beta)$  for any  $\mathbf{x} \in Q_\beta$  and each  $\beta \in \{0, 1, \dots, K-1\}^d$ .

The first step is straightforward. The construction of  $\Phi_1$  in Step 2 is trivial if the network size is large enough. To control the width and depth of the network implementing  $\Phi_1$ , we establish a theorem, Theorem 3.12 in Section 3.4, to help construct  $\Phi_1$ . Assume  $\phi_1$  is the one-dimensional step function attained by Theorem 3.12, then

we can attain  $\Phi_1$  via defining

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

See Figure 2.7 for an illustration.

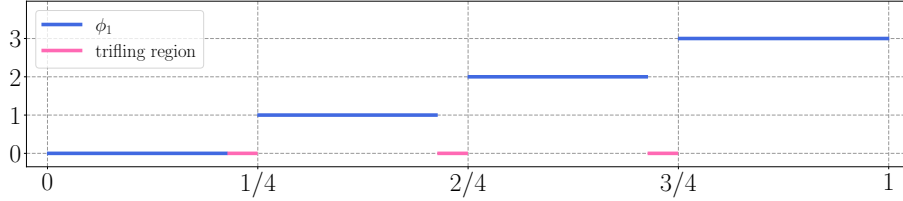


Figure 2.7: An example of a step function for the case  $K = 4$  and  $d = 1$ . We do not need to care about the values of  $\phi_1$  in the trifling region while constructing a ReLU network to implement  $\phi_1$ .

Step 3 is the core step. We would like to point out that we only need to let  $\phi_2$  map  $\beta$  approximately to  $\tilde{f}(\mathbf{x}_\beta)$  within an error  $\mathcal{O}(\omega_f(1/K))$  for each  $\beta \in \{0, 1, \dots, K-1\}^d$  when constructing  $\phi_2$  in Step 3. In other words, it is not necessary to care about the values of  $\phi_2$  outside the set of points  $\{0, 1, \dots, K-1\}^d$ , which plays a key role in constructing a ReLU network to implement  $\phi_2$  in Step 3. Thus, with  $\Phi_1$  in hand, a function approximation problem is converted to a point fitting problem for  $\phi_2$  via the idea of function compositions  $(\phi_2 \circ \Phi_1)$ ,<sup>③</sup> which reveals the power of function compositions in some sense. However, designing a network to solve such a point fitting problem is still a challenging task due to the limitation of the width and depth of the target network. To simplify the construction of a ReLU network solving this point fitting problem, we investigate the width power (Theorem 3.2) and the depth power (Theorem 3.4) of ReLU networks to fit a collection of points in Section 3.2.1 and 3.2.2, respectively. Then we can construct the desired ReLU network by combining these two properties together.

The final step is pretty technical, since  $\phi$  may oscillate greatly in the trifling region. To overcome this, we use two main ideas: “horizontal shift” and “middle

<sup>③</sup>Solving a point fitting problem is to design a function to fit a collection of points  $\{(\mathbf{x}_i, y_i)\}_i$  in  $\mathbb{R}^{d+1}$ , namely, the target function takes the value close to  $y_i$  at the location  $\mathbf{x}_i$ .

value". For example, if  $g$  approximates a one-dimensional continuous function  $f$  well except for an interval in  $\mathbb{R}$  with a small length  $\delta$ , then

$$\text{mid}(g(x - \delta), g(x), g(x + \delta))$$

can approximate  $f$  well on the whole domain  $\mathbb{R}$ , where  $\text{mid}(\cdot, \cdot, \cdot)$  is a function returning the middle value of three inputs. See Section 3.3 for more details.

### 2.3.2 Floor-ReLU networks

Next, let us discuss the general ideas of the approximation by Floor-ReLU networks, which are similar to those of ReLU networks except for the trifling region. Since Floor-ReLU networks can approximate step functions uniformly well on  $[0, 1]^d$ , we do not need to introduce the trifling region again. The main ideas can be divided into four steps as follows.

1. Normalize  $f$  as  $\tilde{f}$  satisfying  $\tilde{f}(\mathbf{x}) \in [0, 1]$  for any  $\mathbf{x} \in [0, 1]^d$ , partition  $[0, 1]^d$  into a set of non-overlapping cubes  $\{Q_\beta\}_{\beta \in \{0, 1, \dots, K-1\}^d}$ , and denote  $\mathbf{x}_\beta$  as the vertex of  $Q_\beta$  with minimum  $\|\cdot\|_1$  norm, where  $K$  is an integer determined later. See Figure 2.8 for the illustrations of  $Q_\beta$  and  $\mathbf{x}_\beta$  for any  $\beta \in \{0, 1, \dots, K-1\}^d$ .
2. Construct a Floor-ReLU sub-network to implement a vector-valued function  $\Phi_1 : \mathbb{R}^d \rightarrow \mathbb{R}^d$  projecting the whole cube  $Q_\beta$  to the index  $\beta$  for each  $\beta$ , *i.e.*,  $\Phi_1(\mathbf{x}) = \beta$  for all  $\mathbf{x} \in Q_\beta$  and each  $\beta \in \{0, 1, \dots, K-1\}^d$ .
3. Construct a Floor-ReLU sub-network to implement a function  $\phi_2 : \mathbb{R}^d \rightarrow \mathbb{R}$  mapping  $\beta \in \{0, 1, \dots, K-1\}^d$  approximately to  $\tilde{f}(\mathbf{x}_\beta)$ , *i.e.*,  $\phi_2(\beta) \approx \tilde{f}(\mathbf{x}_\beta)$  for each  $\beta$ . Then  $\phi_2 \circ \Phi_1(\mathbf{x}) = \phi_2(\beta) \approx \tilde{f}(\mathbf{x}_\beta)$  for any  $\mathbf{x} \in Q_\beta$  and each  $\beta \in \{0, 1, \dots, K-1\}^d$ , implying  $\tilde{\phi} := \phi_2 \circ \Phi_1$  approximates  $\tilde{f}$  within an error  $\mathcal{O}(\omega_f(1/K))$  on  $[0, 1]^d$ .
4. Re-scale and shift  $\tilde{\phi}$  to obtain the desired function  $\phi$  approximating  $f$  well and determine the final Floor-ReLU network to implement  $\phi$ .

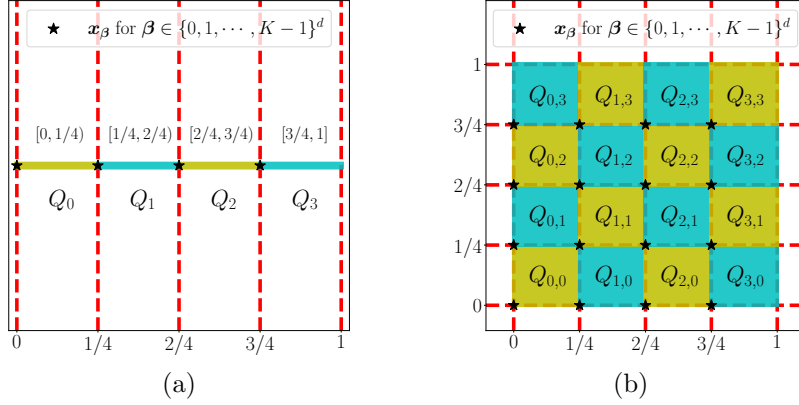


Figure 2.8: Illustrations of  $Q_\beta$  and  $x_\beta$  for any  $\beta \in \{0, 1, \dots, K-1\}^d$ . (a)  $K = 4, d = 1$ . (b)  $K = 4, d = 2$ .

The implementations of Step 1, 2, and 4 are straightforward. Step 3 is the core step. Similar to ReLU networks, we only need to solve a point fitting problem due to the power of function compositions. It is still a highly technical problem. Thus, we introduce a proposition, Proposition 5.6, to help implement this step. As we shall see later in Section 5.3, the key idea of proving Proposition 5.6 is the modified “bit extraction” technique derived from [5].

Finally, we would like to point out that the key reason Floor-ReLU networks can attain much better approximation errors than those of ReLU networks is that Floor ( $\lfloor x \rfloor$ ) has infinite (constant) pieces, while ReLU ( $\max\{0, x\}$ ) has only two (linear) pieces. Thus, roughly speaking, one Floor activation function can do what many ReLU activation functions do in our construction. For this reason, compared to ReLU networks, Floor-ReLU networks attain significantly better approximation errors.



## Basic results of ReLU networks

In this chapter, we introduce several basic results of ReLU networks, which will be used in the later chapters.

### 3.1 Wide networks to deep ones

Generally, it is easier to construct shallow and wide sub-networks to meet the requirements during designing the final network. To control the width of the final network, we consider representing wide and shallow networks by deep and narrow ones. To this end, we establish a theorem, Theorem 3.1 below, to convert wide networks with two hidden layers to deep and narrow ones.

**Theorem 3.1.** *For any  $N, L, d \in \mathbb{N}^+$ , it holds that*

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

This theorem shows that if a function  $\phi$  can be implemented by a two-hidden-layer ReLU network that the first and second hidden layers have  $N$  and  $NL$  neurons, respectively, then there exists a new ReLU network with width  $2N + 2$  and depth  $L + 1$  to implement  $\phi$ .

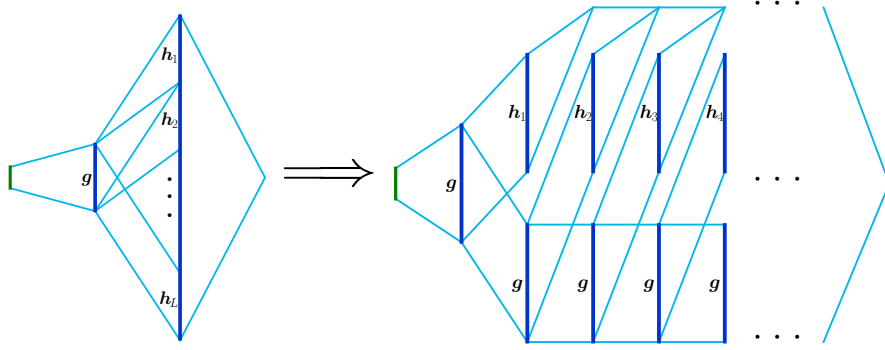


Figure 3.1: An illustration of the main idea of proving Theorem 3.1

The key idea to prove Theorem 3.1 is to re-assemble sub-networks in the shallower network in the left of Figure 3.1 to form a deeper one with width  $\mathcal{O}(N)$  and depth  $\mathcal{O}(L)$  on the right of Figure 3.1.

*Proof of Theorem 3.1.* For any  $\phi \in \mathcal{NN}(\#input = d; \text{widthvec} = [N, NL]; \#output = 1)$ ,  $\phi$  can be implemented by a ReLU network described as

$$\mathbf{x} \xrightarrow[\sigma]{\mathbf{W}_0, \mathbf{b}_0} \mathbf{g} \xrightarrow[\sigma]{\mathbf{W}_1, \mathbf{b}_1} \mathbf{h} \xrightarrow{\mathbf{W}_2, \mathbf{b}_2} \phi(\mathbf{x}),$$

where  $\mathbf{g}$  and  $\mathbf{h}$  are the output of the first and second hidden layers, respectively. That is,

$$\mathbf{g} = \sigma(\mathbf{W}_0 \cdot \mathbf{x} + \mathbf{b}_0), \quad \mathbf{h} = \sigma(\mathbf{W}_1 \cdot \mathbf{g} + \mathbf{b}_1), \quad \text{and} \quad \phi(\mathbf{x}) = \mathbf{W}_2 \cdot \mathbf{h} + \mathbf{b}_2.$$

We can evenly divide  $\mathbf{h} \in \mathbb{R}^{NL}$ ,  $\mathbf{b}_1 \in \mathbb{R}^{NL}$ ,  $\mathbf{W}_1 \in \mathbb{R}^{NL \times N}$ , and  $\mathbf{W}_2 \in \mathbb{R}^{1 \times NL}$  into  $L$  parts as follows.

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

and  $\mathbf{W}_2 = [\mathbf{W}_{2,1}, \mathbf{W}_{2,2}, \dots, \mathbf{W}_{2,L}]$ , where  $\mathbf{h}_\ell \in \mathbb{R}^N$ ,  $\mathbf{b}_{1,\ell} \in \mathbb{R}^N$ ,  $\mathbf{W}_{1,\ell} \in \mathbb{R}^{N \times N}$ , and

$\mathbf{W}_{2,\ell} \in \mathbb{R}^{1 \times N}$  for  $\ell = 1, 2, \dots, L$ . Then, for  $\ell = 1, 2, \dots, L$ , we have

$$\mathbf{h}_\ell = \sigma(\mathbf{W}_{1,\ell} \cdot \mathbf{g} + \mathbf{b}_{1,\ell}) \quad \text{and} \quad \phi(\mathbf{x}) = \mathbf{W}_2 \cdot \mathbf{h} + \mathbf{b}_2 = \sum_{j=1}^L \mathbf{W}_{2,j} \cdot \mathbf{h}_j + \mathbf{b}_2. \quad (3.1)$$

Define

$$s_0 := 0 \quad \text{and} \quad s_\ell := \sum_{j=1}^{\ell} \mathbf{W}_{2,j} \cdot \mathbf{h}_j, \quad \text{for } \ell = 1, 2, \dots, L.$$

Then  $\phi(\mathbf{x}) = \mathbf{W}_2 \cdot \mathbf{h} + \mathbf{b}_2 = s_L + \mathbf{b}_2$  and

$$s_\ell = s_{\ell-1} + \mathbf{W}_{2,\ell} \cdot \mathbf{h}_\ell, \quad \text{for } \ell = 1, 2, \dots, L. \quad (3.2)$$

Hence, it is easy to check that  $\phi$  can be also implemented by the deep network shown in Figure 3.2. Clearly, the network architecture in Figure 3.2 is with width  $2N + 2$

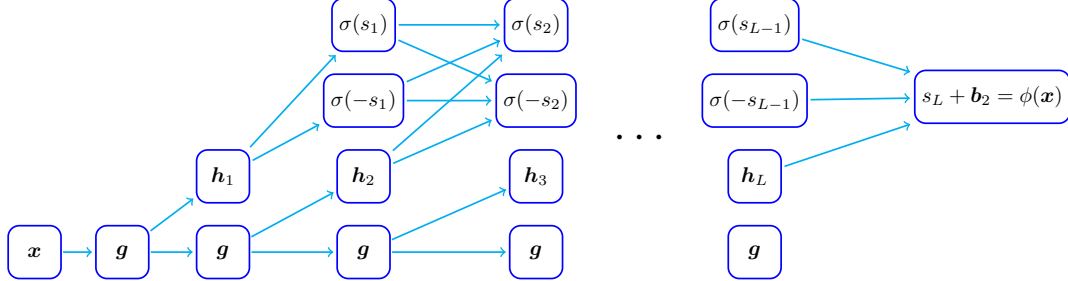


Figure 3.2: An illustration of the desired network implementing  $\phi$  based on Equation (3.1) and (3.2), and the fact  $x = \sigma(x) - \sigma(-x)$  for any  $x \in \mathbb{R}$ .<sup>①</sup>

and depth  $L + 1$ . So we finish the proof.  $\square$

## 3.2 Power of networks to fit points

As mentioned earlier in Section 2.3, we need to construct a ReLU sub-networks with the desired width and depth to solve a point fitting problem. To this end,

<sup>①</sup>In this figure, we omit ReLU ( $\sigma$ ) for a neuron if its output is non-negative without ReLU. Such a simplification will be applied to similar figures in the rest of this dissertation.

we discuss the power of ReLU networks to fit points from two perspectives: 1) the width power of ReLU networks to fit points in Section 3.2.1; 2) the depth power of ReLU networks to fit points in Section 3.2.2.

### 3.2.1 Width power of networks to fit points

Let us first discuss the width power of ReLU network to fit points. Roughly speaking, we would like to minimize the width by fixing the depth when constructing ReLU networks to fit a given number of points. In fact, we prove in Theorem 3.2 that a function  $\phi \in \mathcal{NN}(\#input = 1; \text{widthvec} = [2m, 2n + 1]; \#output = 1)$  can fit  $m(n + 1) + 1$  points in  $\mathbb{R}^2$  with several conditions.

**Theorem 3.2.** *For any  $m, n \in \mathbb{N}^+$ , given any  $m(n + 1) + 1$  samples  $(x_i, y_i) \in \mathbb{R}^2$  with  $x_0 < x_1 < x_2 < \dots < x_{m(n+1)}$  and  $y_i \geq 0$  for  $i = 0, 1, \dots, m(n + 1)$ , there exists  $\phi \in \mathcal{NN}(\#input = 1; \text{widthvec} = [2m, 2n + 1]; \#output = 1)$  satisfying the following three conditions.*

- (i)  $\phi(x_i) = y_i$  for  $i = 0, 1, \dots, m(n + 1)$ .
- (ii)  $\phi$  is linear on each interval  $[x_{i-1}, x_i]$  for all  $i \notin \{j(n + 1) : j = 1, 2, \dots, m\}$ .
- (iii)  $\phi$  is bounded by a constant determined by  $m, n, x_i, y_i$  for  $i = 0, 1, \dots, m(n + 1)$ .

To be exact,

$$\sup_{x \in [x_0, x_{m(n+1)}]} |\phi(x)| \leq C \max_{i \in \{0, 1, \dots, m(n+1)\}} y_i,$$

where

$$C = 1 + \prod_{k=1}^n \left( 1 + \max \left\{ \frac{x_{j(n+1)+n} - x_{j(n+1)+k-1}}{x_{j(n+1)+k} - x_{j(n+1)+k-1}} : j = 0, 1, \dots, m - 1 \right\} \right).$$

We would like to point out that  $\phi$  may not be linear on an interval  $[x_{i-1}, x_i]$  for some  $i \in \{j(n + 1) : j = 1, 2, \dots, m\}$ . So  $\phi$  may oscillate greatly in the region

$$\bigcup_{i \in \{j(n+1) : j=1, 2, \dots, m\}} [x_{i-1}, x_i],$$

which is called the “don’t-care” region in the proof of Theorem 3.2. However, we are able to choose the values of  $x_0, x_1, \dots, x_{m(n+1)}$  properly to make the “don’t-care” region small enough, the idea of which is similar to that of the trifling region defined in Equation (2.1).

Before proving Theorem 3.2, let us first study the properties of ReLU networks with only one hidden layer to warm up in Lemma 3.3 below. Recall that for a continuous piecewise linear function  $f(x)$ , the  $x$  values where the slope changes are typically called **breakpoints**.

**Lemma 3.3.** *Suppose  $\phi \in \text{NN}(\#input = 1; \text{widthvec} = [N]; \#output = 1)$  can be implemented by a ReLU network architecture*

$$x \xrightarrow{\mathbf{W}_0, \mathbf{b}_0} \mathbf{h} \xrightarrow{\sigma} \tilde{\mathbf{h}} \xrightarrow{\mathbf{W}_1, \mathbf{b}_1} \phi(x).$$

That is,  $\phi(x)$  is a function determined by  $\mathbf{W}_0$ ,  $\mathbf{b}_0$ ,  $\mathbf{W}_1$ , and  $\mathbf{b}_1$ . Given a sequence of strictly increasing numbers  $x_0, x_1, \dots, x_N$ , set  $\mathbf{W}_0 = (1, 1, \dots, 1) \in \mathbb{R}^{N \times 1}$  and  $\mathbf{b}_0 = (-x_0, -x_1, \dots, -x_{N-1}) \in \mathbb{R}^N$ . Then we have

- (i) The breakpoints of  $\phi$  are exactly  $x_0, x_1, \dots, x_N$  on the interval  $[x_0, x_N]$ <sup>②</sup>;
- (ii) For any sequence  $(y_i)_{i=0}^N$ , we are able to choose  $\mathbf{W}_1$  and  $\mathbf{b}_1$  properly such that  $\phi(x_i) = y_i$  for  $i = 0, 1, \dots, N$  and  $\phi$  is linear on each interval  $[x_i, x_{i+1}]$  for  $i = 0, 1, \dots, N - 1$ .

Part (i) in Lemma 3.3 is straightforward. The existence in Part (ii) is equivalent to the existence of a solution for a non-singular system of linear equations, which is left for the reader.

With Lemma 3.3 in hand, we are ready to prove Theorem 3.2.

*Proof of Theorem 3.2.* For any  $\phi \in \mathcal{NN}(\#input = 1; \text{widthvec} = [2m, 2n+1]; \#output =$

---

<sup>②</sup>We only consider the interval  $[x_0, x_N]$  and hence  $x_0$  and  $x_N$  are treated as breakpoints.  $\phi(x)$  might not have a real breakpoint in a small open neighborhood of  $x_0$  or  $x_N$ .

1),  $\phi$  can be implemented by the following ReLU network architecture

$$x \xrightarrow{\mathbf{W}_0, \mathbf{b}_0} \mathbf{h} \xrightarrow{\sigma} \tilde{\mathbf{h}} \xrightarrow{\mathbf{W}_1, \mathbf{b}_1} \mathbf{g} \xrightarrow{\sigma} \tilde{\mathbf{g}} \xrightarrow{\mathbf{W}_2, \mathbf{b}_2} \phi(x). \quad (3.3)$$

Clearly,  $\phi(x)$  is a function determined by  $\mathbf{W}_0, \mathbf{b}_0, \mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2$ . So our goal is to choose  $\mathbf{W}_0, \mathbf{b}_0, \mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2$  properly in order to make Condition (i)-(iii) true.

Note that  $\mathbf{g} = \mathbf{g}(x)$  is a vector-valued function mapping  $x \in \mathbb{R}$  to  $\mathbf{g}(x) \in \mathbb{R}^{2n+1}$  and determined by  $\mathbf{W}_0, \mathbf{b}_0, \mathbf{W}_1, \mathbf{b}_1$ . Hence each entry of  $\mathbf{g}(x)$  itself is a function implemented by a sub-network with one hidden layer. Denote  $\mathbf{g} = (g_0, g_1^+, g_1^-, \dots, g_n^+, g_n^-)$ , then  $\{g_0, g_1^+, g_1^-, \dots, g_n^+, g_n^-\} \subseteq \mathcal{NN}(\#input = 1; \text{widthvec} = [2m]; \#output = 1)$ . See Figure 3.3 for an illustration of  $\mathbf{g} = (g_0, g_1^+, g_1^-, \dots, g_n^+, g_n^-)$  and  $\tilde{\mathbf{g}} = \sigma(\mathbf{g}) = (\tilde{g}_0, \tilde{g}_1^+, \tilde{g}_1^-, \dots, \tilde{g}_n^+, \tilde{g}_n^-)$  for the case  $m = n = 2$ . Our proof of Theorem 3.2 is mainly

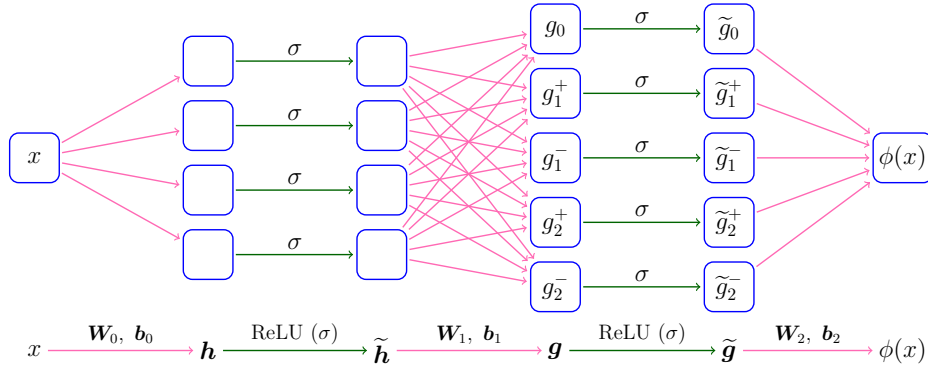


Figure 3.3: An illustration of  $\mathbf{g} = (g_0, g_1^+, g_1^-, \dots, g_n^+, g_n^-)$  and  $\tilde{\mathbf{g}} = \sigma(\mathbf{g}) = (\tilde{g}_0, \tilde{g}_1^+, \tilde{g}_1^-, \dots, \tilde{g}_n^+, \tilde{g}_n^-)$  for the case  $m = n = 2$ .

based on the repeated applications of Lemma 3.3 to determine  $\mathbf{W}_0, \mathbf{b}_0, \mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2$  such that Conditions (i)-(iii) hold.

To simplify the notations, we define two index sets  $\mathcal{I}_1(m, n)$  and  $\mathcal{I}_2(m, n)$  for any  $m, n \in \mathbb{N}^+$  as

$$\mathcal{I}_1(m, n) := \{j(n+1) : j = 1, 2, \dots, m\}$$

and

$$\mathcal{I}_2(m, n) := \mathcal{I}_1(m, n) \cup (\mathcal{I}_1(m, n) - 1) \cup \{0\},$$

where  $\mathcal{I}_1(m, n) - 1 = \{k - 1 : k \in \mathcal{I}_1(m, n)\}$ . For example,  $\mathcal{I}_1(4, 4) = \{5, 10, 15, 20\}$  and  $\mathcal{I}_2(4, 4) = \{0, 4, 5, 9, 10, 14, 15, 19, 20\}$ .

**Step 1:** Determine  $\mathbf{W}_0$  and  $\mathbf{b}_0$ .

Clearly, the index set  $\mathcal{I}_2(m, n)$  has  $2m + 1$  elements. Convert the point set  $\{x_i : i \in \mathcal{I}_2(m, n)\}$  in ascending order to a vector  $\boldsymbol{\xi} = (\xi_0, \xi_1, \dots, \xi_{2m}) \in \mathbb{R}^{2m+1}$ . Then set  $\mathbf{W}_0 = (1, 1, \dots, 1) \in \mathbb{R}^{2m \times 1}$  and  $\mathbf{b}_0 = (-\xi_0, -\xi_1, \dots, -\xi_{2m-1}) \in \mathbb{R}^{2m}$ . Note that  $\xi_{2m} = x_{m(n+1)}$  is the right endpoint of the interval  $[x_0, x_{m(n+1)}]$ . By Lemma 3.3 (set  $N = 2m$  therein), we have

- All functions in  $\{g_0, g_1^+, g_1^-, \dots, g_n^+, g_n^-\}$  have the same set of breakpoints

$$\{\xi_j : j = 0, 1, \dots, 2m\} = \{x_i : i \in \mathcal{I}_2(m, n)\},$$

that is, each function in  $\{g_0, g_1^+, g_1^-, \dots, g_n^+, g_n^-\}$  is linear between any two adjacent points of  $\{x_i : i \in \mathcal{I}_2(m, n)\}$ , no matter what  $\mathbf{W}_1$  and  $\mathbf{b}_1$  are.

- We are able to identify  $\mathbf{W}_1 \in \mathbb{R}^{(2n+1) \times 2m}$  and  $\mathbf{b}_1 \in \mathbb{R}^{2n+1}$  to freely determine the values of each function in  $\{g_0, g_1^+, g_1^-, \dots, g_n^+, g_n^-\}$  at all points of  $\{x_i : i \in \mathcal{I}_2(m, n)\}$ .

**Step 2:** Determine  $\mathbf{W}_1$  and  $\mathbf{b}_1$ .

This is the key step of the proof. Our ultimate goal is to set up

$$\mathbf{g} = (g_0, g_1^+, g_1^-, \dots, g_n^+, g_n^-)$$

by determining  $\mathbf{W}_1$  and  $\mathbf{b}_1$  such that, after a nonlinear activation function (ReLU), there exists a linear combination in the last step of our network (specified by  $\mathbf{W}_2$  and  $\mathbf{b}_2$  as shown in Equation (3.3)) that can generate a desired function  $\phi(x)$  matching the sample points  $\{(x_i, y_i)\}_{0 \leq i \leq m(n+1)}$ . In the previous step, we have determined the breakpoints of  $\{g_0, g_1^+, g_1^-, \dots, g_n^+, g_n^-\}$  by setting up  $\mathbf{W}_0$  and  $\mathbf{b}_0$ ; in this step, we will

identify  $\mathbf{W}_1 \in \mathbb{R}^{(2n+1) \times 2m}$  and  $\mathbf{b}_1 \in \mathbb{R}^{2n+1}$  to fully determine  $\{g_0, g_1^+, g_1^-, \dots, g_n^+, g_n^-\}$ . This will be conducted in two sub-steps.

**Step 2.1:** Set up.

Let  $f_0(x)$  be a continuous piecewise linear function defined on  $[0, 1]$  satisfying

- $f_0(x_i) = y_i$  for all  $i \in \{0, 1, \dots, m(n+1)\}$ .
- $f_0$  is linear between any two adjacent points of  $\{x_i : i \in \{0, 1, \dots, m(n+1)\}\}$ .

Note that  $\{x_i : i \in \mathcal{I}_2(m, n)\}$  is the set of breakpoints of  $g_0$ . By Lemma 3.3 and the setting of Step 1, we are able to choose  $\mathbf{W}_1(1, :)$  and  $\mathbf{b}_1(1)$  properly such that  $g_0(x_i) = f_0(x_i)$  for all  $i \in \mathcal{I}_2(m, n)$  and  $g_0$  is linear between any two adjacent points of  $\{x_i : i \in \mathcal{I}_2(m, n)\}$ .

We would like to inductively construct a sequence of  $f_k$  for all  $k \in \{1, 2, \dots, n+1\}$  satisfying

- $f_k(x_i) = 0$  for all  $i \in \cup_{\ell=0}^{k-1} (\mathcal{I}_1(m, n) - n - 1 + \ell) \cup \{m(n+1)\}$ .
- $f_k$  is linear on each interval  $[x_{i-1}, x_i]$  for all  $i \notin \mathcal{I}_1(m, n)$ .

As we shall see later in Step 3, the construction of the final function  $\phi$  is mainly based on  $f_{n+1}$ .

First, let us consider the case  $k = 1$ . Define  $f_1 := f_0 - \tilde{g}_0$ , where  $\tilde{g}_0 = \sigma(g_0) = g_0$  as shown in Equation (3.3), since  $g_0$  is positive by the construction of Lemma 3.3. Note that

$$(\mathcal{I}_1(m, n) - n - 1) \cup \{m(n+1)\} = \{j(n+1) : j = 0, 1, \dots, m\} \subseteq \mathcal{I}_2(m, n).$$

Then we have

- $f_1(x_i) = f_0(x_i) - \tilde{g}_0(x_i) = 0$  for all  $i \in (\mathcal{I}_1(m, n) - n - 1) \cup \{m(n+1)\}$ .
- $f_1$  is linear on each interval  $[x_{i-1}, x_i]$  for all  $i$ .



Thus, the desired  $f_1$  has been constructed. See Figure 3.4 (a) for an illustration of  $f_0$ ,  $f_1$ , and  $g_0$ .

**Step 2.2:** Mathematical induction.

The initialization of the mathematical induction,  $f_1$ , has been constructed in Step 2.1. Hence, it is enough to show how to proceed with an arbitrary  $k$ . See Figure 3.4 (b)-(d) for the illustration of the first two induction steps.

Now assume  $f_k$  is defined for some  $k \in \{1, 2, \dots, n\}$ , we need to construct  $f_{k+1}$  satisfying similar conditions as follows.

- $f_{k+1}(x_i) = 0$  for all  $i \in \cup_{\ell=0}^k (\mathcal{I}_1(m, n) - n - 1 + \ell) \cup \{m(n+1)\}$ .
- $f_{k+1}$  is linear on each interval  $[x_{i-1}, x_i]$  for all  $i \notin \mathcal{I}_1(m, n)$ .

Then we shall determine

$$\mathbf{W}_1(2k, :), \quad \mathbf{b}_1(2k), \quad \mathbf{W}_1(2k+1, :), \quad \text{and} \quad \mathbf{b}_1(2k+1)$$

to completely specify  $g_k^+$  and  $g_k^-$ , which in turn can determine  $f_{k+1}$ . This induction process can be further divided into four sub-steps.

**Step 2.2.1:** Define index sets.

Define

$$\Lambda_k^+(m, n) := \{j : f_k(x_{j(n+1)+k}) \geq 0, \ 0 \leq j < m\}$$

and

$$\Lambda_k^-(m, n) := \{j : f_k(x_{j(n+1)+k}) < 0, \ 0 \leq j < m\}.$$

Clearly,  $\Lambda_k^+(m, n) \cup \Lambda_k^-(m, n) = \{0, 1, \dots, m-1\}$ . Recall that  $g_k^+$  and  $g_k^-$  are two continuous piecewise linear functions with the same set of breakpoints  $\{x_i : i \in \mathcal{I}_2(m, n)\}$ . We will use  $\Lambda_k^+(m, n)$  and  $\Lambda_k^-(m, n)$  to generate  $2m+1$  samples in

$$\left\{ (x, y) \in \mathbb{R}^2 : x \in \{x_i : i \in \mathcal{I}_2(m, n)\} \right\}$$

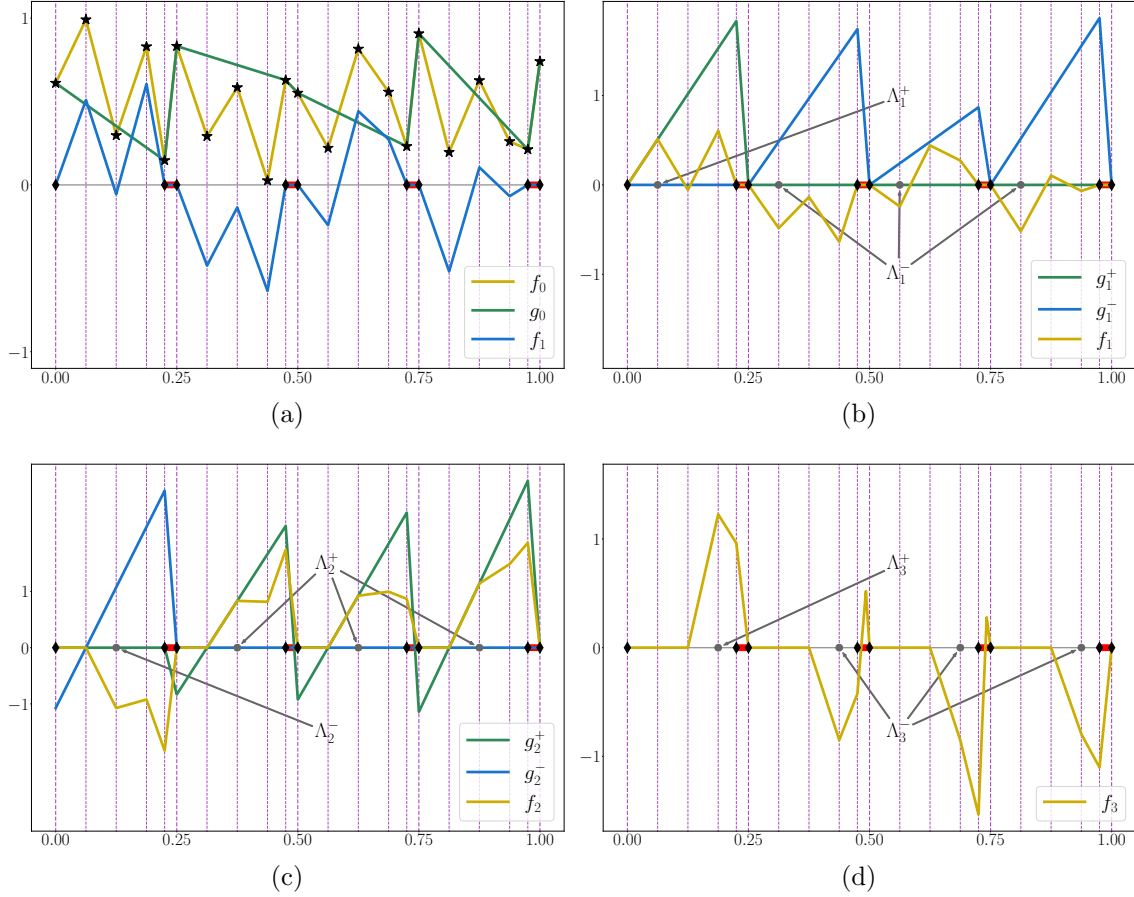


Figure 3.4: Illustrations of the proof of Theorem 3.2, especially Step 2 of the proof, when  $m = n = 4$ , with the “don’t-care” region  $\cup_{i \in \mathcal{I}_1(m,n)} [x_{i-1}, x_i]$  in red.  $g_0, g_1^+, g_1^-, \dots, g_n^+, g_n^-$  share the same set of breakpoints  $\{x_i : i \in \mathcal{I}_2(m,n)\}$  marked with black “diamonds”.  $\Lambda_k^+$  and  $\Lambda_k^-$  are short of  $\Lambda_k^+(m,n)$  and  $\Lambda_k^-(m,n)$ , respectively. (a) Given samples  $\{(x_i, y_i) : i = 0, 1, \dots, m(n+1)\}$  marked with black “stars”, let  $f_0(x)$  be the continuous piecewise linear function fitting these samples, construct  $g_0$  such that  $f_1 = f_0 - \sigma(g_0)$  is closer to 0 than  $f_0$  in a larger subset of the “important” region. (b) Construct  $g_1^+$  and  $g_1^-$  such that  $f_2 = f_1 - \sigma(g_1^+) + \sigma(g_1^-)$  is closer to 0 than  $f_1$  in a larger subset of the “important” region. (c) Construct  $g_2^+$  and  $g_2^-$  such that  $f_3 = f_2 - \sigma(g_2^+) + \sigma(g_2^-)$  is closer to 0 than  $f_2$  in a larger subset of the “important” region. (d) The visualization of  $f_3$ , which is 0 in the “important” areas that have been processed and may remain large near the “don’t-care” region.  $f_k$  will decay quickly to 0 outside the “don’t-care” region as  $k$  increases.

to **fully determine**  $g_k^+$  and  $g_k^-$  by identifying  $\mathbf{W}_1(2k, :)$ ,  $\mathbf{b}_1(2k)$ ,  $\mathbf{W}_1(2k+1, :)$ , and  $\mathbf{b}_1(2k+1)$  in the following steps.

**Step 2.2.2:** Determine  $\mathbf{W}_2(2k, :)$  and  $\mathbf{b}_2(2k)$ .

By Lemma 3.3 and the setting of Step 1, we can choose  $\mathbf{W}_2(2k, :)$  and  $\mathbf{b}_2(2k)$  to fully determine  $g_k^+$  such that each  $g_k^+(x_i)$  matches a specific value for all  $i \in \mathcal{I}_2(m, n)$ . Note that  $\mathcal{I}_2(m, n)$  is the union of three sets:  $\{m(n+1)\}$ ,

$$\{j(n+1) : j \in \Lambda_k^+(m, n) \cup \Lambda_k^-(m, n)\},$$

and

$$\{j(n+1) + n : j \in \Lambda_k^+(m, n) \cup \Lambda_k^-(m, n)\}.$$

The values of  $\{g_k^+(x_i) : i \in \mathcal{I}_2(m, n)\}$  are specified as follows.

- If  $j \in \Lambda_k^+(m, n)$ , specify the values of  $g_k^+(x_{j(n+1)})$  and  $g_k^+(x_{j(n+1)+n})$  such that

$$g_k^+(x_{j(n+1)+k-1}) = 0 \quad \text{and} \quad g_k^+(x_{j(n+1)+k}) = f_k(x_{j(n+1)+k}) \geq 0.$$

The existence of these values fulfilling the requirements above comes from the fact that  $g_k^+$  is linear on the interval  $[x_{j(n+1)}, x_{j(n+1)+n}]$  and  $g_k^+$  only depends on the values of  $g_k^+(x_{j(n+1)+k-1})$  and  $g_k^+(x_{j(n+1)+k})$  on  $[x_{j(n+1)}, x_{j(n+1)+n}]$ . See Figure 3.5 for an illustration. Now it is easy to verify that  $\tilde{g}_k^+ = \sigma(g_k^+)$  satisfies

$$* \quad \tilde{g}_k^+(x_{j(n+1)+\ell}) = 0 \text{ for } \ell = 0, 1, \dots, k-1 \text{ and}$$

$$\tilde{g}_k^+(x_{j(n+1)+k}) = f_k(x_{j(n+1)+k}) \geq 0.$$

$$* \quad \tilde{g}_k^+ \text{ is linear on each interval } [x_{j(n+1)+\ell}, x_{j(n+1)+\ell+1}] \text{ for all } \ell \in \{0, 1, \dots, n-1\}.$$

- If  $j \in \Lambda_k^-(m, n)$ , let  $g_k^+(x_{j(n+1)}) = g_k^+(x_{j(n+1)+n}) = 0$ . Then  $\tilde{g}_k^+ = \sigma(g_k^+) = 0$  on the interval  $[x_{j(n+1)}, x_{j(n+1)+n}]$ .

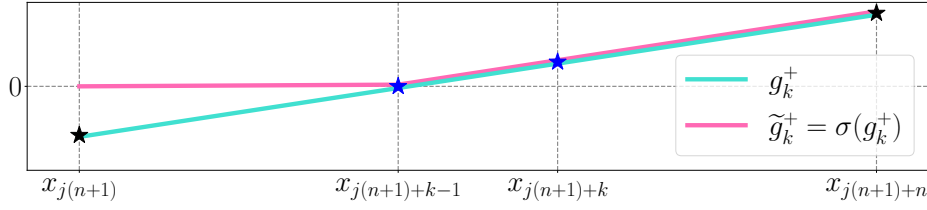


Figure 3.5: An illustration of  $g_k^+$  and  $\tilde{g}_k^+ = \sigma(g_k^+)$ . To design  $f_{k+1}$  with  $f_{k+1}(x_{j(n+1)+k}) = 0$ , we shall specify the  $y$ -coordinates of two blue “stars” as  $f_k(x_{j(n+1)+k-1}) = 0$  and  $f_k(x_{j(n+1)+k}) \geq 0$ , respectively. Four “stars” should be kept in a straight line. Thus, two blue “stars” determine two black “stars”, which in turn determine  $g_k^+$  on  $[x_{j(n+1)}, x_{j(n+1)+n}]$  since the  $x$ -coordinates of two black “stars” are two adjacent breakpoints of  $g_k^+$ . By doing so, we have  $f_k - \tilde{g}_k^+ = 0$  at  $x_{j(n+1)+\ell}$  for  $\ell = 0, 1, \dots, k$ , which is a big step forward in constructing  $f_{k+1}$ .

- Finally, specify the value of  $g_k^+$  at  $x_{m(n+1)}$  as 0.

**Step 2.2.3:** Determine  $\mathbf{W}_2(2k+1, :)$  and  $\mathbf{b}_2(2k+1)$ .

Similarly, we choose  $\mathbf{W}_2(2k+1, :)$  and  $\mathbf{b}_2(2k+1)$  such that  $g_k^-$  matches specific values as follows.

- If  $j \in \Lambda_k^-(m, n)$ , specify the values of  $g_k^-(x_{j(n+1)})$  and  $g_k^-(x_{j(n+1)+n})$  such that

$$g_k^-(x_{j(n+1)+k-1}) = 0 \quad \text{and} \quad g_k^-(x_{j(n+1)+k}) = -f_k(x_{j(n+1)+k}) > 0.$$

Then  $\tilde{g}_k^- = \sigma(g_k^-)$  satisfies

$$* \quad \tilde{g}_k^-(x_{j(n+1)+\ell}) = 0 \text{ for } \ell = 0, 1, \dots, k-1 \text{ and}$$

$$\tilde{g}_k^-(x_{j(n+1)+k}) = -f_k(x_{j(n+1)+k}) > 0.$$

$$* \quad \tilde{g}_k^- \text{ is linear on each interval } [x_{j(n+1)+\ell}, x_{j(n+1)+\ell+1}] \text{ for all } \ell \in \{0, 1, \dots, n-1\}.$$

- If  $j \in \Lambda_k^+(m, n)$ , let  $g_k^-(x_{j(n+1)}) = g_k^-(x_{j(n+1)+n}) = 0$ . Then  $\tilde{g}_k^- = \sigma(g_k^-) = 0$  on the interval  $[x_{j(n+1)}, x_{j(n+1)+n}]$ .

- Finally, specify the value of  $g_k^-$  at  $x_{m(n+1)}$  as 0.

**Step 2.2.4:** Construct  $f_{k+1}$  from  $g_k^+$  and  $g_k^-$ .

For the sake of clarity, the properties of  $g_k^+$  and  $g_k^-$  constructed in Step 2.2.3 are summarized below.

- $\tilde{g}_k^+(x_i) = \tilde{g}_k^-(x_i) = 0$  for all  $i \in \cup_{\ell=0}^{k-1} (\mathcal{I}_1(m, n) - n - 1 + \ell) \cup \{m(n+1)\}$ .
- If  $j \in \Lambda_k^+(m, n)$ ,  $\tilde{g}_k^+(x_{j(n+1)+k}) = f_k(x_{j(n+1)+k}) \geq 0$  and  $\tilde{g}_k^-(x_{j(n+1)+k}) = 0$ .
- If  $j \in \Lambda_k^-(m, n)$ ,  $\tilde{g}_k^-(x_{j(n+1)+k}) = -f_k(x_{j(n+1)+k}) > 0$  and  $\tilde{g}_k^+(x_{j(n+1)+k}) = 0$ .
- $\tilde{g}_k^+$  and  $\tilde{g}_k^-$  are linear on each interval  $[x_{j(n+1)+\ell}, x_{j(n+1)+\ell+1}]$  for each  $\ell \in \{0, 1, \dots, n-1\}$  and each  $j \in \Lambda_k^+(m, n) \cup \Lambda_k^-(m, n) = \{0, 1, \dots, m-1\}$ . In other words,  $\tilde{g}_k^+$  and  $\tilde{g}_k^-$  are linear on each interval  $[x_{i-1}, x_i]$  for all  $i \notin \{j(n+1) : j = 1, 2, \dots, m\} = \mathcal{I}_1(m, n)$ .

See Figure 3.4 (a)-(c) for the illustration of  $g_0$ ,  $g_1^+$ ,  $g_1^-$ ,  $g_2^+$ , and  $g_2^-$ , and to verify their properties as listed above. By the induction hypothesis, we have

- $f_k(x_i) = 0$  for all  $i \in \cup_{\ell=0}^{k-1} (\mathcal{I}_1(m, n) - n - 1 + \ell) \cup \{m(n+1)\}$ .
- $f_k$  is linear on each interval  $[x_{i-1}, x_i]$  for all  $i \notin \mathcal{I}_1(m, n)$ .

Thus,  $f_k(x_i) - \tilde{g}_k^+(x_i) + \tilde{g}_k^-(x_i) = 0$  for all  $i$  in

$$\begin{aligned} & \left( \cup_{\ell=0}^{k-1} (\mathcal{I}_1(m, n) - n - 1 + \ell) \cup \{m(n+1)\} \right) \cup \left\{ j(n+1) + k : j \in \Lambda_k^+(m, n) \cup \Lambda_k^-(m, n) \right\} \\ &= \cup_{\ell=0}^k (\mathcal{I}_1(m, n) - n - 1 + \ell) \cup \{m(n+1)\}, \end{aligned}$$

where the equality comes from the fact  $\Lambda_k^+(m, n) \cup \Lambda_k^-(m, n) = \{0, 1, \dots, m-1\}$ .

Therefore, by defining

$$f_{k+1} := f_k - \tilde{g}_k^+ + \tilde{g}_k^-,$$

we have

- $f_{k+1}(x_i) = 0$  for all  $i \in \cup_{\ell=0}^k (\mathcal{I}_1(m, n) - n - 1 + \ell) \cup \{m(n+1)\}$ .

- $f_{k+1}$  is linear on each interval  $[x_{i-1}, x_i]$  for all  $i \notin \mathcal{I}_1(m, n)$ .

See Figure 3.4 (b)-(d) for the illustration of  $f_1$ ,  $f_2$ , and  $f_3$ , and to verify their properties as listed just above. This finishes the mathematical induction process. As we can imagine based on Figure 3.4, when  $k$  increases, the support of  $f_k$  shrinks to the “don’t-care” region.

**Step 3:** Determine  $\mathbf{W}_2$  and  $\mathbf{b}_2$ .

With the special vector-valued function  $\mathbf{g} = (g_0, g_1^+, g_1^-, \dots, g_n^+, g_n^-)$  constructed in Step 2, we are able to specify  $\mathbf{W}_2$  and  $\mathbf{b}_2$  to generate a desired  $\phi(x)$  matching the samples  $\{(x_i, y_i)\}_{0 \leq i \leq m(n+1)}$ .

In fact, we can simply set  $\mathbf{W}_2 = (1, 1, -1, 1, -1, \dots, 1, -1) \in \mathbb{R}^{1 \times (2n+1)}$  and  $\mathbf{b}_2 = 0$ , which finishes the implementation of  $\phi = \tilde{g}_0 + \sum_{\ell=1}^n \tilde{g}_\ell^+ - \sum_{\ell=1}^n \tilde{g}_\ell^-$ . The rest of the proof is to verify the properties of  $\phi$ . By the principle of mathematical induction, we have

- $f_{n+1} = f_1 + \sum_{\ell=1}^n (\tilde{g}_\ell^- - \tilde{g}_\ell^+) = f_0 - \tilde{g}_0 - \sum_{\ell=1}^n \tilde{g}_\ell^+ + \sum_{\ell=1}^n \tilde{g}_\ell^- = f_0 - \phi$ .
- $f_{n+1}(x_i) = 0$  for all  $i$  in

$$\cup_{\ell=0}^n (\mathcal{I}_1(m, n) - n - 1 + \ell) \cup \{m(n+1)\} = \{0, 1, \dots, m(n+1)\}.$$

- $f_{n+1}$  is linear on each interval  $[x_{i-1}, x_i]$  for all  $i \notin \mathcal{I}_1(m, n)$ .

Hence,  $\phi = \tilde{g}_0 + \sum_{\ell=1}^n \tilde{g}_\ell^+ - \sum_{\ell=1}^n \tilde{g}_\ell^- = f_0 - f_{n+1}$ . Then

$$\phi(x_i) = f_0(x_i) - f_{n+1}(x_i) = y_i, \quad \text{for all } i \in \{0, 1, \dots, m(n+1)\},$$

which verifies Condition (i), and  $\phi = f_0 - f_{n+1}$  is linear on each interval  $[x_{i-1}, x_i]$  for  $i \notin \mathcal{I}_1(m, n)$ , which verifies Condition (ii). It remains to check that  $\phi$  satisfies Condition (iii).

By the definition of  $f_1 = f_0 - \tilde{g}_0$ , we have

$$-\max_{i \in \{0,1,\dots,m(n+1)\}} y_i \leq -\tilde{g}_0(x) \leq f_0(x) - \tilde{g}_0(x) \leq f_0(x) \leq \max_{i \in \{0,1,\dots,m(n+1)\}} y_i,$$

for any  $x \in [x_0, x_{m(n+1)}]$ , implying

$$\sup_{x \in [x_0, x_{m(n+1)}]} |f_1(x)| \leq \max_{i \in \{0,1,\dots,m(n+1)\}} y_i.$$

By the induction process in Step 2, for any  $k \in \{1, 2, \dots, n\}$ , it holds that

$$\sup_{x \in [x_0, x_{m(n+1)}]} |\tilde{g}_k^+(x)| \leq C_k(m, n) \sup_{x \in [x_0, x_{m(n+1)}]} |f_k(x)|$$

and

$$\sup_{x \in [x_0, x_{m(n+1)}]} |\tilde{g}_k^-(x)| \leq C_k(m, n) \sup_{x \in [x_0, x_{m(n+1)}]} |f_k(x)|,$$

where

$$C_k(m, n) := \max \left\{ \frac{x_{j(n+1)+n} - x_{j(n+1)+k-1}}{x_{j(n+1)+k} - x_{j(n+1)+k-1}} : j = 0, 1, \dots, m-1 \right\}.$$

Since either  $\tilde{g}_k^+(x)$  or  $\tilde{g}_k^-(x)$  is equal to 0 for any  $x \in [x_0, x_{m(n+1)}]$ , we have

$$\sup_{x \in [x_0, x_{m(n+1)}]} |\tilde{g}_k^+(x) - \tilde{g}_k^-(x)| \leq C_k(m, n) \sup_{x \in [x_0, x_{m(n+1)}]} |f_k(x)|.$$

It follows from  $f_{k+1} = f_k - \tilde{g}_k^+ + \tilde{g}_k^-$  that, for any  $k \in \{1, 2, \dots, n\}$ ,

$$\begin{aligned} \sup_{x \in [x_0, x_{m(n+1)}]} |f_{k+1}(x)| &\leq \sup_{x \in [x_0, x_{m(n+1)}]} |\tilde{g}_k^+(x) - \tilde{g}_k^-(x)| + \sup_{x \in [x_0, x_{m(n+1)}]} |f_k(x)| \\ &\leq (C_k(m, n) + 1) \sup_{x \in [x_0, x_{m(n+1)}]} |f_k(x)|. \end{aligned}$$

Hence,

$$\begin{aligned} \sup_{x \in [x_0, x_{m(n+1)}]} |f_{n+1}(x)| &\leq \left( \prod_{k=1}^n (C_k(m, n) + 1) \right) \sup_{x \in [x_0, x_{m(n+1)}]} |f_1(x)| \\ &\leq \left( \prod_{k=1}^n (C_k(m, n) + 1) \right) \max_{i \in \{0, 1, \dots, m(n+1)\}} y_i. \end{aligned}$$

Therefore,

$$\begin{aligned} \sup_{x \in [x_0, x_{m(n+1)}]} |\phi(x)| &\leq \sup_{x \in [x_0, x_{m(n+1)}]} |f_0(x)| + \sup_{x \in [x_0, x_{m(n+1)}]} |f_{n+1}(x)| \\ &\leq \left( 1 + \prod_{k=1}^n (C_k(m, n) + 1) \right) \max_{i \in \{0, 1, \dots, m(n+1)\}} y_i \\ &:= C \max_{i \in \{0, 1, \dots, m(n+1)\}} y_i, \end{aligned}$$

where

$$C = 1 + \prod_{k=1}^n \left( 1 + \max \left\{ \frac{x_{j(n+1)+n} - x_{j(n+1)+k-1}}{x_{j(n+1)+k} - x_{j(n+1)+k-1}} : j = 0, 1, \dots, m-1 \right\} \right).$$

So we finish the proof.  $\square$

### 3.2.2 Depth power of networks to fit points

Next, let us discuss the depth power of ReLU networks to fit points. Roughly speaking, we would like to minimize the depth by fixing the width while constructing ReLU networks to fit a given number of points. In fact, we prove in Theorem 3.4 that a function  $\phi \in \mathcal{NN}(\#input = 1; \text{width} \leq 8N + 6; \text{depth} \leq 5L + 7; \#output = 1)$  can fit  $N^2 L^2$  points in  $\mathbb{R}^2$  with several conditions.

**Theorem 3.4.** *For any  $N, L \in \mathbb{N}^+$  and any  $\theta_i \in \{0, 1\}$  for  $i = 0, 1, \dots, N^2 L^2 - 1$ , there exists a function  $\phi$  implemented by a ReLU network with width  $8N + 6$  and depth  $5L + 7$  such that*

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



We would like to remark that the key idea in the proof of Theorem 3.4 is the “bit extraction” technique in [5], which allows us to store  $L$  bits in a binary number  $\text{bin}0.\theta_1\theta_2\cdots\theta_L$  and extract each bit  $\theta_i$ . The extraction operator can be efficiently carried out via a deep ReLU network architecture, demonstrating the power of depth.

Next, we introduce Theorem 3.5, a variant of Theorem 3.4, which is easier to prove and can deduce 3.4 simply. Theorem 3.4 and 3.5 characterize the depth power of ReLU networks. Both of them will be used in the later chapters.

**Theorem 3.5.** *For any  $N, L \in \mathbb{N}^+$ , any  $\theta_{m,\ell} \in \{0, 1\}$  for  $m = 0, 1, \dots, M - 1$  and  $\ell = 0, 1, \dots, L - 1$ , where  $M = N^2L$ , there exists a function  $\phi$  implemented by a ReLU network with width  $4N + 3$  and depth  $3L + 3$  such that*

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

We denote  $M = N^2L$  in Theorem 3.5 because it is roughly the number of parameters. The choice of outputting  $\sum_{j=0}^{\ell} \theta_{m,j}$  rather than  $\theta_{m,\ell}$  not only guarantees the proof of Theorem 3.4 but also simplifies the construction of ReLU networks to approximate continuous functions in  $C([0, 1]^d)$  in Section 4.2.

Theorem 3.5 will be proven later in this section. Let us first prove Theorem 3.4 based on Theorem 3.5.

*Proof of Theorem 3.4.* The case  $L = 1$  is clear. We assume  $L \geq 2$  below.

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

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

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

By Theorem 3.5, there exist  $\phi_1, \phi_2 \in \mathcal{NN}(\text{width} \leq 4N + 3; \text{depth} \leq 3L + 3)$  such that

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

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

We consider the sample set

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

Its size is  $2M + 1 = N \cdot ((2NL - 1) + 1) + 1$ . By Theorem 3.2 (set  $m = N$  and  $n = 2NL - 1$  therein), there exists

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

such that

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

It follows that

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

implying

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

For  $i = 0, 1, \dots, N^2L^2 - 1$ , by representing  $i = mL + \ell$  for  $m = 0, 1, \dots, M - 1$

and  $\ell = 0, 1, \dots, L-1$ , we have  $\psi(i) = \psi(mL + \ell) = m$  and  $i - L\psi(i) = \ell$ , deducing

$$\begin{aligned}
 & \phi_1(\psi(i), i - L\psi(i)) - \phi_2(\psi(i), i - L\psi(i)) \\
 &= \phi_1(m, \ell) - \phi_2(m, \ell) = \sum_{j=0}^{\ell} a_{m,j} - \sum_{j=0}^{\ell} b_{m,j} \\
 &= \sum_{j=0}^{\ell} a_{m,j} - \sum_{j=1}^{\ell} a_{m,j-1} - b_0 = a_{m,\ell} = \theta_i.
 \end{aligned} \tag{3.4}$$

Therefore, the desired function  $\phi$  can be implemented by the network architecture described in Figure 3.6.

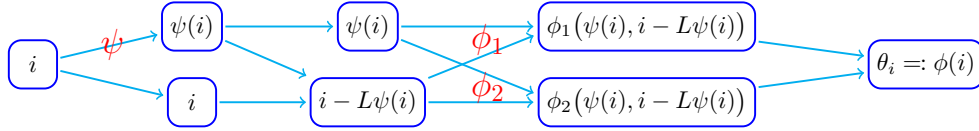


Figure 3.6: An illustration of the network architecture implementing the desired function  $\phi$  based on Equation (3.4) for  $i = 0, 1, \dots, N^2L^2 - 1$ .<sup>③</sup>

Note that

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

By Theorem 3.1,

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

Hence, the network architecture shown in Figure 3.6 is with width

$$\max\{4L + 2 + 1, 2(4L + 3)\} = 8N + 6$$

<sup>③</sup>In this figure, “ $\psi$ ”, “ $\phi_1$ ”, and “ $\phi_2$ ” and cyan arrows (“ $\longrightarrow$ ”) adjacent to them represent the ReLU networks implementing themselves. We use similar notations in the rest of this dissertation. For example, “ $\xrightarrow{\phi}$ ” means the network architecture that implements a function  $\phi: \mathbb{R}^2 \rightarrow \mathbb{R}$ .

and depth

$$(2L + 1) + 2 + (3L + 3) + 1 = 5L + 7,$$

implying  $\phi \in \mathcal{NN}(\text{width} \leq 8N + 6; \text{depth} \leq 5L + 7)$ . So we finish the proof.  $\square$

It remains to prove Theorem 3.5, which relies on the “bit extraction” technique introduced in [5]. We modify this technique to extract the sum of many bits rather than one bit and this modification can be summarized in Lemma 3.6 below.

**Lemma 3.6** (Bit extraction). *For any  $L \in \mathbb{N}^+$ , there exists a function  $\phi$  in*

$$\mathcal{NN}(\#input = 2; \text{width} \leq 7; \text{depth} \leq 2L + 1; \#output = 1)$$

*such that, for any  $\theta_1, \theta_2, \dots, \theta_L \in \{0, 1\}$ , we have*

$$\phi(\text{bin}0.\theta_1\theta_2\cdots\theta_L, \ell) = \sum_{j=1}^{\ell} \theta_j, \quad \text{for } \ell = 1, 2, \dots, L.$$

*Proof.* Given any  $\theta_1, \theta_2, \dots, \theta_L \in \{0, 1\}$ , define

$$\xi_j := \text{bin}0.\theta_j\theta_{j+1}\cdots\theta_L, \quad \text{for } j = 1, 2, \dots, L$$

and

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

Then we have

$$\theta_j = \mathcal{T}(\xi_j - 1/2), \quad \text{for } j = 1, 2, \dots, L,$$

and

$$\xi_{j+1} = 2\xi_j - \theta_j, \quad \text{for } j = 1, 2, \dots, L - 1.$$

We would like to point out that, by above two iteration equations, we can iteratively get  $\xi_1, \theta_1, \xi_2, \theta_2, \dots, \xi_L, \theta_L$  when  $\xi_1 = \text{bin}0.\theta_1\theta_2\cdots\theta_L$  is given. Based on this idea, the rest proof can be divided into three steps.

**Step 1:** Simplify two iteration equations.

Note that  $\mathcal{T}(x) = \sigma(x/\delta + 1) - \sigma(x/\delta)$  for any  $x \notin (-\delta, 0)$ . By setting  $\delta = 1/2 - \sum_{j=2}^L 2^{-j} = 2^{-L}$ , we have  $\xi_j - 1/2 \notin (-\delta, 0)$  for all  $j$ , implying

$$\begin{aligned}\theta_j &= \mathcal{T}(\xi_j - 1/2) = \sigma((\xi_j - 1/2)/\delta + 1) - \sigma((\xi_j - 1/2)/\delta) \\ &= \sigma(\mathcal{L}(\xi_j) + 1) - \sigma(\mathcal{L}(\xi_j)),\end{aligned}\tag{3.5}$$

for  $j = 1, 2, \dots, L$ , where  $\mathcal{L}$  is an affine linear map given by  $\mathcal{L}(x) = (x - 1/2)/\delta$ . It follows that, for  $j = 1, 2, \dots, L - 1$ ,

$$\xi_{j+1} = 2\xi_j - \theta_j = 2\xi_j - \sigma(\mathcal{L}(\xi_j) + 1) + \sigma(\mathcal{L}(\xi_j)).\tag{3.6}$$

**Step 2:** Design a ReLU network to output  $\sum_{j=1}^L \theta_j$ .

It is easy to design a ReLU network to output  $\theta_1, \theta_2, \dots, \theta_L$  by Equation (3.5) and (3.6) when using  $\xi_1 = \text{bin}0.\theta_1\theta_2 \dots \theta_L$  as the input. However, it is highly non-trivial to construct a ReLU network to output  $\sum_{j=1}^L \theta_j$  with another input  $\ell$ , since many operations like multiplication and comparison are not allowed in designing ReLU networks.

Now let us establish a formula to represent  $\sum_{j=1}^L \theta_j$  in a form of a ReLU network. Recall two facts: 1)  $x_1x_2 = \sigma(x_1 + x_2 - 1)$  for any  $x_1, x_2 \in \{0, 1\}$ ; 2)  $\mathcal{T}(n) = \sigma(n + 1) - \sigma(n)$  for any integer  $n$ . Thus, for  $\ell = 1, 2, \dots, L$ , we have

$$\begin{aligned}\sum_{j=1}^{\ell} \theta_j &= \sum_{j=1}^L \theta_j \mathcal{T}(\ell - j) = \sum_{j=1}^L \sigma(\theta_j + \mathcal{T}(\ell - j) - 1) \\ &= \sum_{j=1}^L \sigma(\theta_j + \sigma(\ell - j + 1) - \sigma(\ell - j) - 1).\end{aligned}$$

To simplify the notations, we define

$$z_{\ell,j} := \sigma(\theta_j + \sigma(\ell - j + 1) - \sigma(\ell - j) - 1),\tag{3.7}$$

for  $\ell = 1, 2, \dots, L$  and  $j = 1, 2, \dots, L$ . Then,

$$\sum_{j=1}^{\ell} \theta_j = \sum_{j=1}^L z_{\ell,j}, \quad \text{for } \ell = 1, 2, \dots, L. \quad (3.8)$$

With Equation (3.5), (3.6), (3.7), and (3.8) in hand, it is easy to construct a function  $\phi$  implemented by a ReLU network with the desired width and depth outputting  $\sum_{j=1}^{\ell} \theta_j = \sum_{j=1}^L z_{\ell,j}$  for the given input  $(\xi_1, \ell) = (\text{bin}0.\theta_1\theta_2 \cdots \theta_L, \ell)$  for  $\ell \in \{1, 2, \dots, L\}$  and  $\theta_1, \theta_2, \dots, \theta_L \in \{0, 1\}$ . The detailed construction is shown in Figure 3.7. It is easy to verify by Figure 3.7 that

$$\phi \in \mathcal{NN}(\#input = 2; \text{width} \leq 7; \text{depth} \leq 2L + 1; \#output = 1).$$

So we finish the proof.  $\square$

With Lemma 3.6 in hand, we are ready to prove Theorem 3.5.

*Proof of Theorem 3.5.* Define

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

Consider the sample set  $\{(m, y_m) : m = 0, 1, \dots, M\}$ , whose size is  $M + 1 = N((NL - 1) + 1) + 1$ . By Theorem 3.2 (set  $m = N$  and  $n = NL - 1$  therein), there exists

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

such that

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

By Lemma 3.6, there exists

$$\phi_2 \in \mathcal{NN}(\#input = 2; \text{width} \leq 7; \text{depth} \leq 2L + 1; \#output = 1)$$

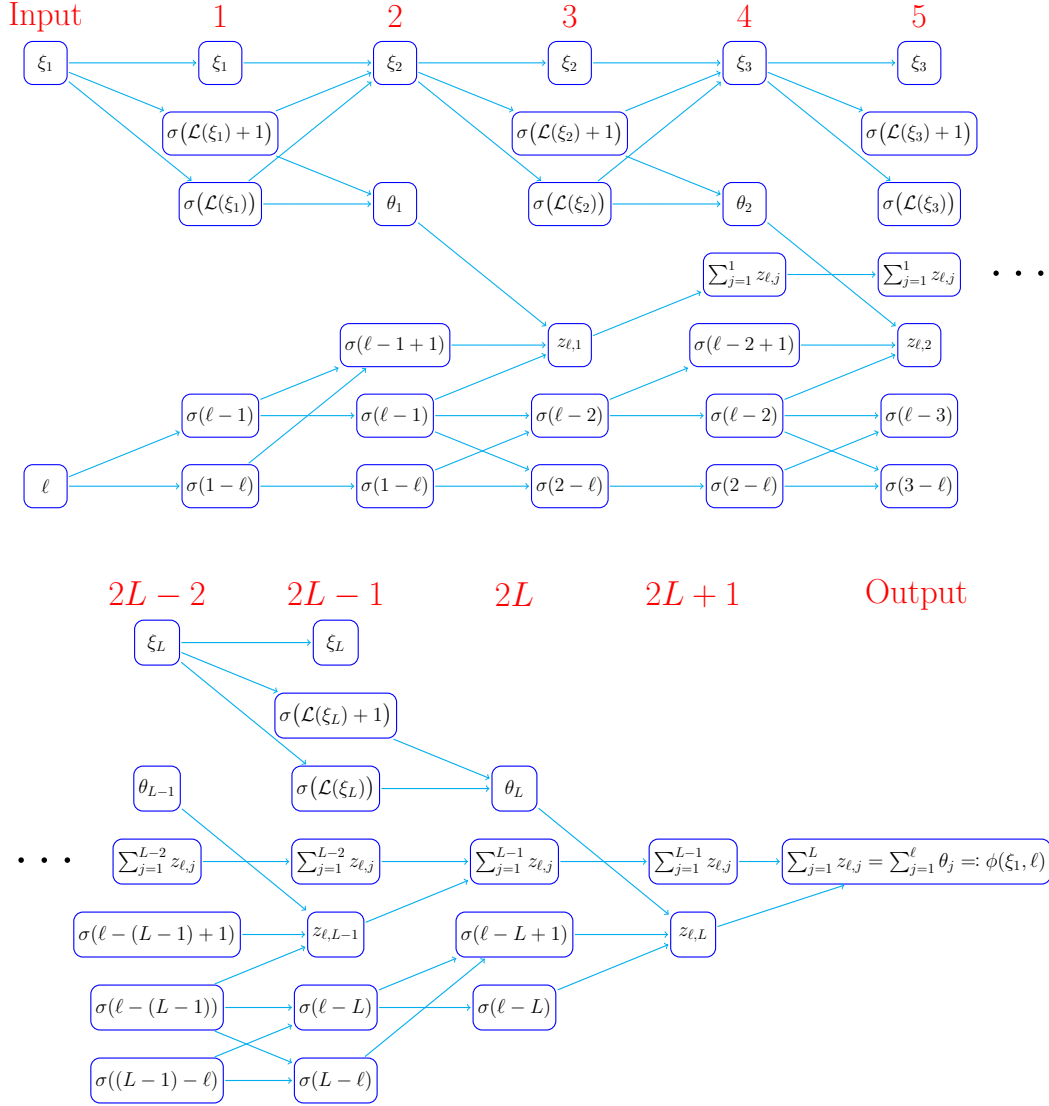


Figure 3.7: An illustration of the target ReLU network implementing  $\phi$  to output  $\sum_{j=1}^L z_{j,\ell} = \sum_{j=1}^L \theta_j = \phi(\xi_1, \ell)$  for the given input  $(\xi_1, \ell) = (\text{bin}0.\theta_1\theta_2\cdots\theta_L, \ell)$  for  $\ell \in \{1, 2, \dots, L\}$  and  $\theta_1, \theta_2, \dots, \theta_L \in \{0, 1\}$ . The construction is mainly based on Equation (3.5), (3.6), (3.7), and (3.8). The red numbers above the architecture indicate the order of hidden layers and every two adjacent layers builds a whole iteration step. We output both  $\sigma(\ell - j)$  and  $\sigma(j - \ell)$  in a hidden layer because we can get the value  $\ell - j$  in the next hidden layer because of the fact  $x = \sigma(x) - \sigma(-x)$  for any  $x \in \mathbb{R}$ . We omit ReLU ( $\sigma$ ) for a neuron if its output is non-negative without ReLU. Note that all parameters of this network are determined by Equation (3.5), (3.6), (3.7), and (3.8), which are valid no matter what  $\theta_1, \theta_2, \dots, \theta_L \in \{0, 1\}$  are. Thus, the desired function  $\phi$  implemented by this network is independent of  $\theta_1, \theta_2, \dots, \theta_L \in \{0, 1\}$ .

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

$$\phi_2(\text{bin}0.\xi_1\xi_2\cdots\xi_L, \ell) = \sum_{j=1}^{\ell} \xi_j, \quad \text{for } \ell = 1, 2, \dots, L.$$

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

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

Thus, for  $m = 0, 1, \dots, M - 1$  and  $\ell = 0, 1, \dots, L - 1$ , we have

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

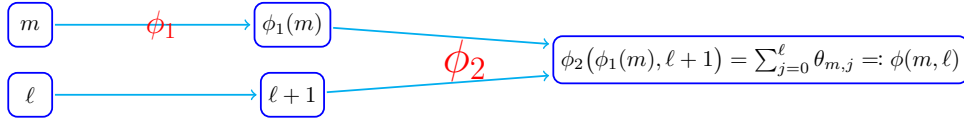


Figure 3.8: An illustration of the network architecture implementing the desired function  $\phi$  based on  $\phi_1$  and  $\phi_2$  for  $m = 0, 1, \dots, M - 1$  and  $\ell = 0, 1, \dots, L - 1$ .

Hence, the desired function  $\phi$  can be implemented by the network shown in Figure 3.8. By Theorem 3.1,  $\phi_1 \in \mathcal{NN}(\text{widthvec} = [2N, 2NL - 1]) \subseteq \mathcal{NN}(\text{width} \leq 4N + 2; \text{depth} \leq L + 1)$ . Recall that  $\phi_2 \in \mathcal{NN}(\text{width} \leq 7; \text{depth} \leq 2L + 1)$ . Therefore, the network in Figure 3.8 is with width  $\max\{(4N + 2) + 1, 7\} = 4N + 3$  and depth  $(L + 1) + 1 + (2L + 1) = 3L + 3$ . So we finish the proof.  $\square$

### 3.3 Approximation in the trifling region

As mentioned earlier in Section 2.3, we need to modify a ReLU network to let it approximate the target function  $f$  uniformly well on the whole region  $[0, 1]^d$ , if this ReLU network approximates  $f$  well outside the trifling region  $\Omega([0, 1]^d, K, \delta)$  defined in Equation (2.1).



**Theorem 3.7.** *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  $\tilde{\phi}$  is a function implemented by a ReLU network with width  $N$  and depth  $L$ . If*

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

*then there exists a new ReLU network with width  $3^d(N + 4)$  and depth  $L + 2d$  implementing a new function  $\phi$  such that*

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

Intuitively speaking, Theorem 3.7 shows that: If a function  $\tilde{\phi}$ , implemented by a ReLU network, approximates  $f$  well except for the trifling region, then we can modify  $\tilde{\phi}$  to get  $\phi$ , implemented by a new ReLU network with similar width and depth to the old one, to approximate  $f$  uniformly well on the whole domain. For example, if  $\tilde{\phi}$  approximates a one-dimensional continuous function  $f$  well except for an interval in  $\mathbb{R}$  with a small length  $\delta$ , then  $\text{mid}(\tilde{\phi}(x - \delta), \tilde{\phi}(x), \tilde{\phi}(x + \delta))$  can approximate  $f$  well on the whole domain, where  $\text{mid}(\cdot, \cdot, \cdot)$  is a function returning the middle value of three inputs and can be implemented via a ReLU network as shown in Lemma 3.8.

**Lemma 3.8.** *The middle value function  $\text{mid}(x_1, x_2, x_3)$  can be implemented by a ReLU network with width 14 and depth 2.*

*Proof.* Recall the fact

$$x = \sigma(x) - \sigma(-x) \quad \text{and} \quad |x| = \sigma(x) + \sigma(-x), \quad \text{for any } x \in \mathbb{R}. \quad (3.9)$$

Therefore,

$$\begin{aligned} \max(x, y) &= \frac{x + y + |x - y|}{2} \\ &= \frac{1}{2}\sigma(x + y) - \frac{1}{2}\sigma(-x - y) + \frac{1}{2}\sigma(x - y) + \frac{1}{2}\sigma(-x + y), \end{aligned} \quad (3.10)$$

for any  $x, y \in \mathbb{R}$ . Thus,  $\max(x_1, x_2, x_3)$  can be implemented by the network shown in Figure 3.9.

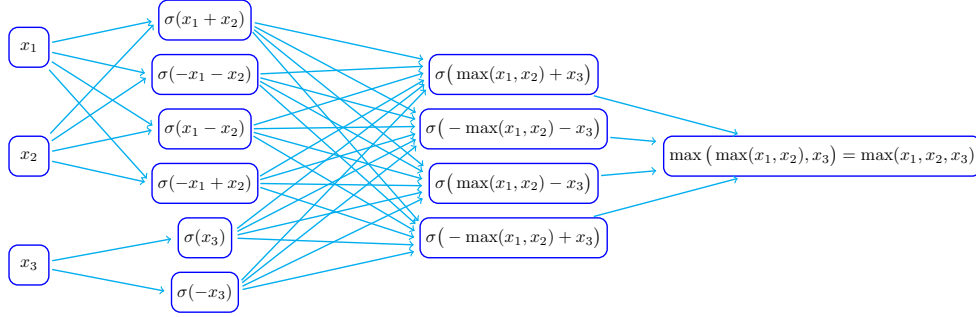


Figure 3.9: An illustration of the network architecture implementing  $\max(x_1, x_2, x_3)$  based on Equation (3.9) and (3.10).

Clearly,

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

Similarly, we have

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

It is easy to check that

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

Hence,

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

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

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

**Lemma 3.9.** *For any  $\varepsilon > 0$ , if at least two of  $\{x_1, x_2, x_3\}$  are in  $\mathcal{B}(y, \varepsilon)$ , then  $\text{mid}(x_1, x_2, x_3) \in \mathcal{B}(y, \varepsilon)$ .*

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

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

$$y - \varepsilon \leq x_1 \leq x_3 \leq x_2 \leq y + \varepsilon.$$

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

So we finish the proof. □

Next, given a function  $g$  approximating  $f$  well on  $[0, 1]$  except for the trifling region, Lemma 3.10 below shows how to use the  $\text{mid}(x_1, x_2, x_3)$  function to construct a new function  $\phi$  uniformly approximating  $f$  well on  $[0, 1]$ , leveraging the useful property of  $\text{mid}(x_1, x_2, x_3)$  in Lemma 3.9.

**Lemma 3.10.** *Given any  $\varepsilon > 0$ ,  $K \in \mathbb{N}^+$ , and  $\delta \in (0, \frac{1}{3K}]$ , assume  $f$  is a continuous function in  $C([0, 1])$  and  $g : \mathbb{R} \rightarrow \mathbb{R}$  is a general function satisfying*

$$|g(x) - f(x)| \leq \varepsilon, \text{ i.e., } f(x) \in \mathcal{B}(g(x), \varepsilon), \quad (3.11)$$

for any  $x \in [0, 1] \setminus \Omega([0, 1], K, \delta)$ . Then

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

where

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

*Proof.* Divide  $[0, 1]$  into  $K$  small intervals denoted by  $Q_k = [\frac{k}{K}, \frac{k+1}{K}]$  for  $k = 0, 1, \dots, K-1$ . For each  $k$ , we further partition  $Q_k$  into four small closed intervals as shown in Figure 3.10. To be exact,

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

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

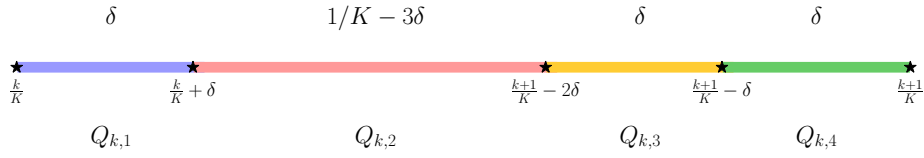


Figure 3.10: An illustration of  $Q_{k,i}$  for  $i = 1, 2, 3, 4$ .

Recall that  $\Omega([0, 1], K, \delta)$  is the trifling region defined in Equation (2.1). Clearly,  $Q_{K-1,4} \subseteq [0, 1] \setminus \Omega([0, 1], K, \delta)$  and  $Q_{k,i} \subseteq [0, 1] \setminus \Omega([0, 1], K, \delta)$  for  $k = 0, 1, \dots, K-1$  and  $i = 1, 2, 3$ .

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

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

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

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

It follows from Equation (3.11) that

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

and

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

By Lemma 3.9, we get

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

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

If  $x \in Q_{k,2}$ , then  $x - \delta, x, x + \delta \in [0, 1] \setminus \Omega([0, 1], K, \delta)$ . It follows from Equation (3.11) that

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

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

and

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

Then, by Lemma 3.9, we have

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

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

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

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

It follows from Equation (3.11) that

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

and

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

By Lemma 3.9, we get

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

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

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

- If  $k \in \{0, 1, \dots, K - 2\}$ , then  $x - \delta \in Q_{k,3} \in [0, 1] \setminus \Omega([0, 1], K, \delta)$  and  $x + \delta \in Q_{k+1,1} \subseteq [0, 1] \setminus \Omega([0, 1], K, \delta)$ . It follows from Equation (3.11) that

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

and

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

By Lemma 3.9, we get

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

- If  $k = K - 1$ , then  $x \in Q_{k,4} = Q_{K-1,4} \subseteq [0, 1] \setminus \Omega([0, 1], K, \delta)$  and  $x - \delta \in Q_{k,3} \subseteq [0, 1] \setminus \Omega([0, 1], K, \delta)$ . It follows from Equation (3.11) that

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

and

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

By Lemma 3.9, we get

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

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

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

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

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

So we finish the proof.  $\square$

The next lemma below is an analog of Lemma 3.10 for the multidimensional case.

**Lemma 3.11.** *Given any  $\varepsilon > 0$ ,  $K \in \mathbb{N}^+$ , and  $\delta \in (0, \frac{1}{3K}]$ , assume  $f$  is a continuous function in  $C([0, 1]^d)$  and  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  is a general function satisfying*

$$|g(\mathbf{x}) - f(\mathbf{x})| \leq \varepsilon, \text{ i.e., } f(\mathbf{x}) \in \mathcal{B}(g(\mathbf{x}), \varepsilon), \quad (3.12)$$

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

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

where  $\phi := \phi_d$  is defined by induction through

$$\phi_{i+1}(\mathbf{x}) := \text{mid}(\phi_i(\mathbf{x} - \delta \mathbf{e}_{i+1}), \phi_i(\mathbf{x}), \phi_i(\mathbf{x} + \delta \mathbf{e}_{i+1})), \quad \text{for } i = 0, 1, \dots, d-1, \quad (3.13)$$

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

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

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

Clearly,  $E_0 = [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)$  and  $E_d = [0, 1]^d$ . See Figure 3.11 for the illustrations of  $E_\ell$  for  $\ell = 0, 1, \dots, d$  when  $K = 4$  and  $d = 2$ .

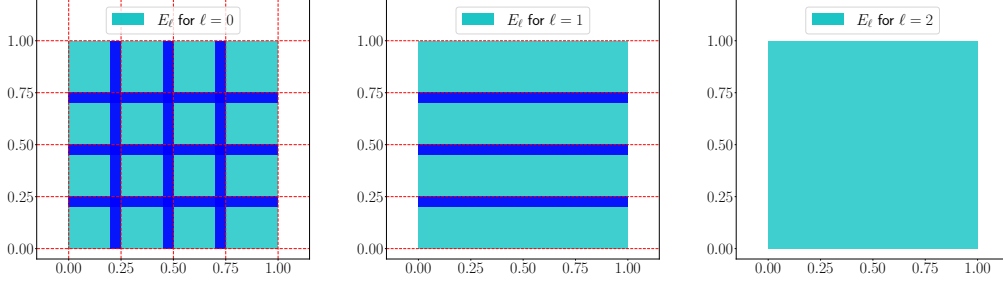


Figure 3.11: Illustrations of  $E_\ell$  for  $\ell = 0, 1, 2$  when  $K = 4$  and  $d = 2$ .

We would like to construct a sequence of functions  $\phi_0, \phi_1, \dots, \phi_d$  by induction, based on the iteration equation (3.13), such that, for each  $\ell \in \{0, 1, \dots, d\}$ ,

$$\phi_\ell(\mathbf{x}) \in \mathcal{B}(f(\mathbf{x}), \varepsilon + \ell \cdot \omega_f(\delta)), \quad \text{for any } \mathbf{x} \in E_\ell. \quad (3.14)$$

Let us first consider the case  $\ell = 0$ . Note that  $[0, 1]^d \setminus \Omega([0, 1]^d, K, \delta) = E_0$ . Then, by Equation (3.12), we have

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

That is, Equation (3.14) is true for  $\ell = 0$ .

Now assume Equation (3.14) is true for  $\ell = i \in \{0, 1, \dots, d-1\}$ . We will prove that it also holds for  $\ell = i+1$ . By the induction hypothesis, we have

$$\phi_i(x_1, \dots, x_i, t, x_{i+2}, \dots, x_d) \in \mathcal{B}\left(f(x_1, \dots, x_i, t, x_{i+2}, \dots, x_d), \varepsilon + i \cdot \omega_f(\delta)\right), \quad (3.15)$$

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

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

$$\mathbf{x}^{[i]} := (x_1, \dots, x_i, x_{i+2}, \dots, x_d) \in \mathbb{R}^{d-1}.$$



Then define

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

and

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

It follows from Equation (3.15) that

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

Then by Lemma 3.10 (set  $g = \psi_{\mathbf{x}^{[i]}}$  and  $f = f_{\mathbf{x}^{[i]}}$  therein), we get, for any  $t \in [0, 1]$ ,

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

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

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

Note that  $x_1, \dots, x_i \in [0, 1]$  and  $x_{i+2}, \dots, x_d \in [0, 1] \setminus \Omega([0, 1], K, \delta)$  are arbitrary.

Thus, for any  $\mathbf{x} \in E_{i+1}$ ,

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

implying

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

So Equation (3.14) holds for  $\ell = i + 1$ , which means we finish the process of mathematical induction.

By the principle of induction, we have

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

Therefore,

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

which means we finish the proof.  $\square$

Now we are ready to prove Theorem 3.7.

*Proof of Theorem 3.7.* Set  $\phi_0 = \tilde{\phi}$  and define  $\phi_i$  for  $i \in \{1, 2, \dots, d\}$  by induction through

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

where  $\{\mathbf{e}_i\}_{i=1}^d$  is the standard basis in  $\mathbb{R}^d$ . Then by Lemma 3.11 with  $\phi = \phi_d$ , we have

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

It remains to determine the network architecture implementing  $\phi = \phi_d$ .

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

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

Note that  $\phi_0 = \tilde{\phi} \in \mathcal{NN}(\text{width} \leq N; \text{depth} \leq L)$ . Hence,  $\phi_0(\cdot - \delta \mathbf{e}_1)$ ,  $\phi_0(\cdot)$ , and  $\phi_0(\cdot + \delta \mathbf{e}_1)$  can be generated by three networks with the same number of hidden layers, and these three networks are all with width  $N$  and depth  $L$ . Therefore, by putting these three networks in parallel (share the inputs), we have  $\Phi_0 \in \mathcal{NN}(\#input = d; \text{width} \leq 3N; \text{depth} \leq L; \#output = 3)$ .

Recall that  $\text{mid}(\cdot, \cdot, \cdot) \in \mathcal{NN}(\text{width} \leq 14; \text{depth} \leq 2)$  by Lemma 3.8. Therefore,

by Lemma 2.1 (ii),

$$\phi_1 = \min(\cdot, \cdot, \cdot) \circ \Phi_0 \in \mathcal{NN}(\text{width} \leq \max\{3N, 14\} \leq 3(N+4); \text{depth} \leq L+2)$$

Similarly, by the iterative formula

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

it is easy to verify

$$\phi = \phi_d \in \mathcal{NN}(\text{width} \leq 3^d(N+4); \text{depth} \leq L+2d).$$

So we finish the proof.  $\square$

### 3.4 Approximation of step functions

As mentioned earlier in Section 2.3, we need to construct a ReLU sub-network to project a cube to a point. We only need to approximate one-dimensional step functions, because in the multidimensional case we can simply set

$$\Phi(\mathbf{x}) = (\phi(x_1), \phi(x_2), \dots, \phi(x_d)), \quad \text{for any } \mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d,$$

where  $\phi$  is a one-dimensional step function. The theorem below, Theorem 3.12, shows that ReLU networks with width  $\mathcal{O}(N^{1/d})$  and depth  $\mathcal{O}(L)$  can implement one-dimensional step functions with  $\mathcal{O}(K) = \mathcal{O}(N^{2/d}L^{2/d})$  “steps” outside the trifling region for any  $d \in \mathbb{N}^+$ . See Figure 3.12 for an example.

**Theorem 3.12.** *For any  $N, L, d \in \mathbb{N}^+$  and  $\delta \in (0, \frac{1}{3K}]$  with  $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor$ , there exists a one-dimensional function  $\phi$  implemented by a ReLU network with*

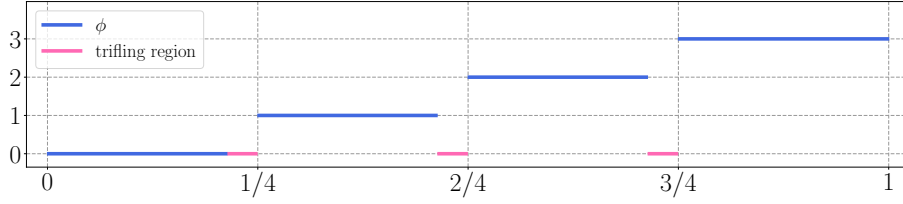


Figure 3.12: An example of a step function for the case  $K = 4$  and  $d = 1$ . We do not need to care about the values of  $\phi$  in the trifling region while constructing a ReLU network to implement  $\phi$ .

width  $4\lfloor N^{1/d} \rfloor + 3$  and depth  $4L + 5$  such that

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

The setting  $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor = \mathcal{O}(N^{2/d} L^{2/d})$  is not neat here, but it is very convenient for later use. Now, let us present the detailed proof of Theorem 3.12.

*Proof of Theorem 3.12.* We divide the proof into two cases:  $d = 1$  and  $d \geq 2$ .

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

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

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

Its size is  $2M + 1 = N \cdot ((2NL - 1) + 1) + 1$ . By Theorem 3.2 (set  $m = N$  and  $n = 2NL - 1$  therein), there exists

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

such that

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

2.

- $\phi_1$  is linear on  $[\frac{M-1}{M}, 1]$  and each interval  $[\frac{m}{M}, \frac{m+1}{M} - \delta]$  for  $m = 0, 1, \dots, M-2$ .

Then, for  $m = 0, 1, \dots, M-1$ , we have

$$\phi_1(x) = m, \quad \text{for any } x \in [\frac{m}{M}, \frac{m+1}{M} - \delta \cdot 1_{\{m \leq M-2\}}]. \quad (3.16)$$

Now consider the another sample set

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

Its size is  $2L+1 = 1 \cdot ((2L-1)+1) + 1$ . By Theorem 3.2 (set  $m = 1$  and  $n = 2L-1$  therein), there exists

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

such that

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

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

$$\phi_2(x - \frac{m}{M}) = \ell, \quad \text{for any } x \in [\frac{mL+\ell}{ML}, \frac{mL+\ell+1}{ML} - \delta \cdot 1_{\{\ell \leq L-2\}}]. \quad (3.17)$$

$K = ML$  implies any  $k \in \{0, 1, \dots, K-1\}$  can be unique represented by  $k = mL + \ell$  for  $m = 0, 1, \dots, M-1$  and  $\ell = 0, 1, \dots, L-1$ . Then the desired function  $\phi$  can be implemented by a ReLU network shown in Figure 3.13.

Clearly,

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

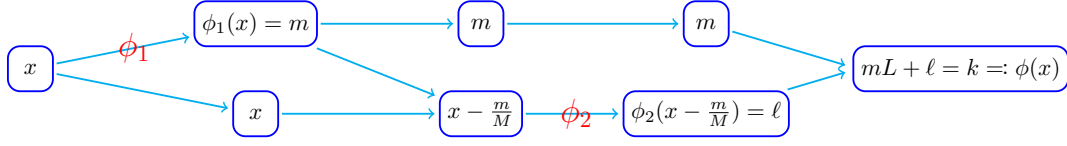


Figure 3.13: An illustration of the network architecture implementing  $\phi$  based on Equation (3.16) and (3.17) for  $x \in [\frac{k}{K}, \frac{k}{K} - \delta \cdot 1_{\{k \leq K-2\}}] = [\frac{mL+\ell}{ML}, \frac{mL+\ell+1}{ML} - \delta \cdot 1_{\{m \leq M-2 \text{ or } \ell \leq L-2\}}]$ , where  $k = mL + \ell$  for  $m = 0, 1, \dots, M-1$  and  $\ell = 0, 1, \dots, L-1$ .

By Theorem 3.1, we have

$$\phi_1 \in \mathcal{NN}(\text{widthvec} = [2N, 4NL - 1]) \subseteq \mathcal{NN}(\text{width} \leq 4N + 2; \text{depth} \leq 2L + 1)$$

and

$$\phi_2 \in \mathcal{NN}(\text{widthvec} = [2, 4L - 1]) \subseteq \mathcal{NN}(\text{width} \leq 6; \text{depth} \leq 2L + 1),$$

implying by Figure 3.13 that  $\phi$  can be implemented by a ReLU network with width

$$\max\{4N + 2 + 1, 6 + 1\} = 4N + 3$$

and depth

$$(2L + 1) + 2 + (2L + 1) + 1 = 4L + 5.$$

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

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

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

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

Its size is

$$2K + 1 = \lfloor N^{1/d} \rfloor ((2 \lfloor N^{1/d} \rfloor \lfloor L^{2/d} \rfloor - 1) + 1) + 1.$$

By Theorem 3.2 (set  $m = \lfloor N^{1/d} \rfloor$  and  $n = 2\lfloor N^{1/d} \rfloor \lfloor L^{2/d} \rfloor - 1$  therein), there exists

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

such that

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

Then, for  $k = 0, 1, \dots, K - 1$ , we have

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

By Theorem 3.1,

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

Thus, we finish the proof for the case  $d \geq 2$ . □

This page is intentionally left blank.



# Approximation by ReLU networks

## 4.1 Approximation of polynomials

In this section, we show how to construct a ReLU network to approximate a multidimensional polynomial  $P(\mathbf{x})$ .

### 4.1.1 Main theorem

For simplicity, we may assume a polynomial  $P(\mathbf{x})$  has only one term with coefficient one, namely,  $P(\mathbf{x}) = \mathbf{x}^\alpha = x_1^{\alpha_1} x_2^{\alpha_2} \cdots x_d^{\alpha_d}$  for some  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_d) \in \mathbb{N}^d$ . As shown in the following theorem, Theorem 4.1, ReLU networks can uniformly approximate polynomials on  $[0, 1]^d$  with exponential errors.

**Theorem 4.1.** *Given any  $k \in \mathbb{N}^+$  and  $\alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{N}^d$ , assume  $P(\mathbf{x}) = \mathbf{x}^\alpha = x_1^{\alpha_1} x_2^{\alpha_2} \cdots x_d^{\alpha_d}$  with  $\|\alpha\|_1 \leq k$ . For any  $N, L \in \mathbb{N}^+$ , there exists a function  $\phi$  implemented by a ReLU network with width  $9(N+1) + k - 1$  and depth  $7k^2L$  such that*

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

The choice of depth  $7k^2L$  is not neat here, but it is convenient for later use. Theorem 4.1 shows that ReLU networks with width  $\mathcal{O}(N)$  and depth  $\mathcal{O}(L)$  are able to approximate polynomials on  $[0, 1]^d$  within an error  $\mathcal{O}(N^{-L})$ . This reveals the

power of depth in ReLU networks for approximating polynomials, from function compositions. Theorem 4.1 can be generalized to the case of polynomials defined on an arbitrary hypercube  $[a, b]^d$  by scaling. To prove Theorem 4.1, we will construct ReLU networks to approximate polynomials following the four steps below.

- $f(x) = x^2$ . We approximate  $f(x) = x^2$  by the combinations and compositions of “sawtooth” functions as shown in Figure 4.1 and 4.2.
- $f(x, y) = xy$ . To approximate  $f(x, y) = xy$ , we use the result of the previous step and the fact  $xy = 2\left(\left(\frac{x+y}{2}\right)^2 - \left(\frac{x}{2}\right)^2 - \left(\frac{y}{2}\right)^2\right)$ .
- $f(x_1, x_2, \dots, x_k) = x_1 x_2 \cdots x_k$ . We approximate  $f(x_1, x_2, \dots, x_k) = x_1 x_2 \cdots x_k$  via mathematical induction based on the result of the previous step.
- A general polynomial  $P(\mathbf{x}) = \mathbf{x}^\alpha = x_1^{\alpha_1} x_2^{\alpha_2} \cdots x_d^{\alpha_d}$  with  $\|\alpha\|_1 \leq k \in \mathbb{N}^+$ .  $P(\mathbf{x})$  can be written as  $P(\mathbf{x}) = z_1 z_2 \cdots z_k$ , where  $\mathbf{z} = (z_1, z_2, \dots, z_k)$  is a vector with each of its entries equal to 1 or an entry of  $\mathbf{x}$ . Then use the result of the previous step.

### 4.1.2 Approximation of $x^2$

Let us show how to approximate  $f(x) = x^2$  by linear combinations of “sawtooth” functions, which can be easily implemented by ReLU networks. The idea of using “sawtooth” functions (see Figure 4.1) was first raised in [58] for approximating  $x^2$  using networks with depth  $\mathcal{O}(L)$  and a constant width, and achieving an error  $\mathcal{O}(2^{-L})$ . Our construction is different to and more general than that in [58], working for ReLU networks of width  $3N$  and depth  $L$  for any  $N, L \in \mathbb{N}^+$ , and achieving an error  $N^{-L}$  as shown in Lemma 4.2.

**Lemma 4.2.** *For any  $N, L \in \mathbb{N}^+$ , there exists a function  $\phi$  implemented by a ReLU network with width  $3N$  and depth  $L$  such that*

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

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

Let

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

and

$$T_i = T_{i-1} \circ T_1, \quad \text{for } i = 2, 3, 4, \dots$$

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

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

See Figure 4.1 for illustrations of  $T_i$  for  $i = 1, 2, 3, 4$ .

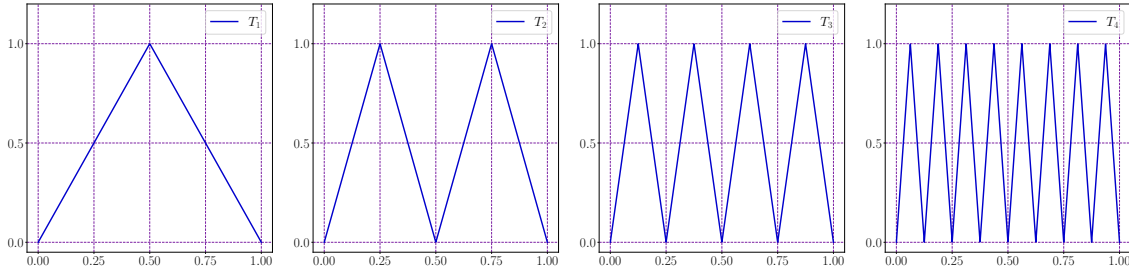


Figure 4.1: Examples of “sawtooth” functions  $T_1$ ,  $T_2$ ,  $T_3$ , and  $T_4$ .

Define piecewise linear functions  $f_s : [0, 1] \rightarrow [0, 1]$  for  $s \in \mathbb{N}^+$  satisfying the following two requirements (see Figure 4.2 for several examples of  $f_s$ ).

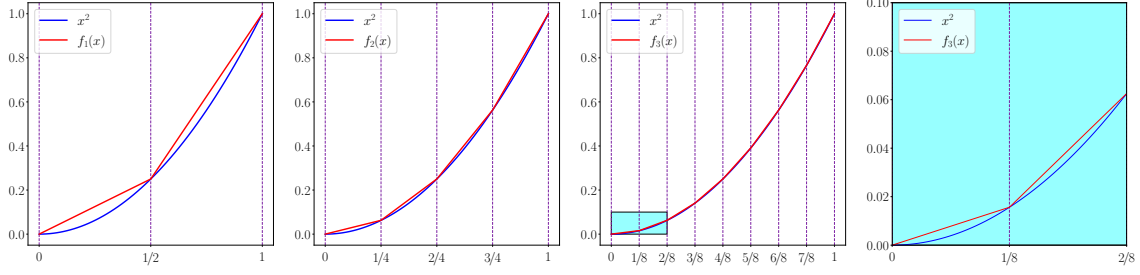
- $f_s(x) = x^2$  for any  $x \in \{\frac{j}{2^s} : j = 0, 1, 2, \dots, 2^s\}$ .
- $f_s(x)$  is linear between any two adjacent points of  $\{\frac{j}{2^s} : j = 0, 1, 2, \dots, 2^s\}$ .

Recall the fact

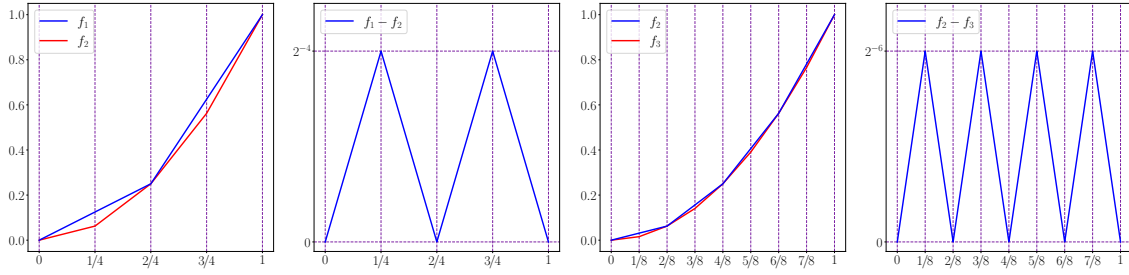
$$0 \leq tx_1^2 + (1-t)x_2^2 - \left(tx_1 + (1-t)x_2\right)^2 \leq \frac{(x_2 - x_1)^2}{4}, \quad \text{for any } t, x_1, x_2 \in [0, 1].$$

Thus, we have

$$0 \leq f_s(x) - x^2 \leq \frac{(2^{-s})^2}{4} = 2^{-2(s+1)}, \quad \text{for any } x \in [0, 1] \text{ and } s \in \mathbb{N}^+. \quad (4.1)$$

Figure 4.2: Illustrations of  $f_1$ ,  $f_2$ , and  $f_3$  for approximating  $x^2$ .

Note that  $f_{i-1}(x) = f_i(x) = x^2$  for each  $x \in \{\frac{j}{2^{i-1}} : j = 0, 1, 2, \dots, 2^{i-1}\}$  and the graph of  $f_{i-1} - f_i$  is a symmetric “sawtooth” between any two adjacent points of  $\{\frac{j}{2^{i-1}} : j = 0, 1, 2, \dots, 2^{i-1}\}$ . See Figure 4.3 for illustrations.

Figure 4.3: Illustrations of  $f_1 - f_2$  and  $f_2 - f_3$ .

Thus, it is easy to verify

$$f_{i-1}(x) - f_i(x) = \frac{T_i(x)}{2^{2i}}, \quad \text{for any } x \in [0, 1] \text{ and } i = 2, 3, 4, \dots.$$

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

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

Given any  $N \in \mathbb{N}^+$ , there exists a unique  $k \in \mathbb{N}^+$  such that

$$(k-1)2^{k-1} + 1 \leq N \leq k2^k.$$

For this  $k$ , we can construct a ReLU network as shown in Figure 4.4 to implement

a function  $\phi = f_{Lk}$  approximating  $x^2$  well. Note that  $T_i$  has  $2^{i-1}$  “sawteeth” and hence has  $2^i + 1$  breakpoints including the endpoints for any  $i \in \mathbb{N}^+$ . Then, by Lemma 3.3,  $T_i$  can be implemented by a one-hidden-layer ReLU network with width  $2^i$ . Therefore, the network in Figure 4.4 has width  $k2^k + 1 \leq 3N$  <sup>①</sup> and depth  $2L$ .

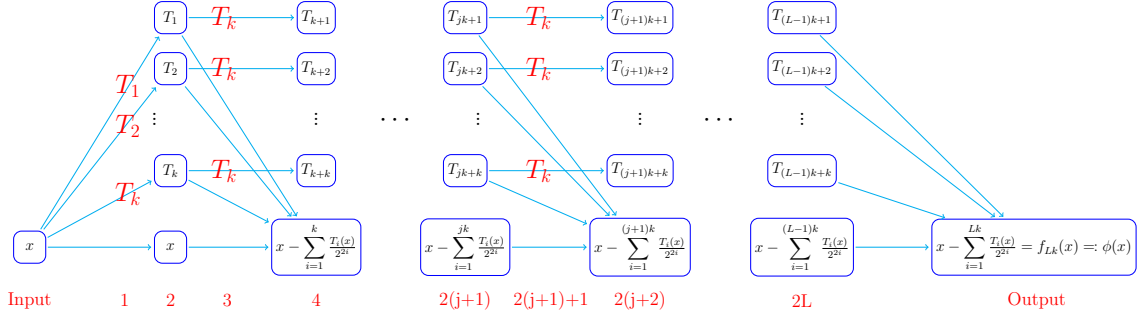


Figure 4.4: An illustration of the target network architecture for approximating  $x^2$  on  $x \in [0, 1]$ .  $T_i$  can be implemented by a one-hidden-layer ReLU network with width  $2^i$  for  $i = 1, 2, \dots, k$ . The red numbers below the architecture indicate the order of hidden layers.

As shown in Figure 4.4, all neurons in  $(2\ell)$ -th hidden layer of the network have the identity function as their activation functions for  $\ell = 1, 2, \dots, L$ . Thus, the network in Figure 4.4 can be interpreted as a ReLU network with width  $3N$  and depth  $L$ . In fact, if all activation functions in a certain hidden layer are identity maps, the depth can be reduced by one via combining adjacent two affine linear transforms into one, the idea of which is similar to that of Lemma 2.1. For example, suppose  $\mathbf{W}_1 \in \mathbb{R}^{N_2 \times N_1}$ ,  $\mathbf{b}_1 \in \mathbb{R}^{N_2}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{N_3 \times N_2}$ ,  $\mathbf{b}_2 \in \mathbb{R}^{N_3}$ ,  $\varrho$  is an identity map that can be applied to vectors or matrices elementwisely, and  $\mathcal{L}_i$  is an affine linear map given by  $\mathcal{L}_i(\mathbf{x}) = \mathbf{W}_i \cdot \mathbf{x} + \mathbf{b}_i$  for  $i = 1, 2$ . Then, we have

$$\mathcal{L}_2 \circ \varrho \circ \mathcal{L}_1(\mathbf{x}) = \mathbf{W}_2 \varrho(\mathbf{W}_1 \cdot \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2 = \mathbf{W}_3 \cdot \mathbf{x} + \mathbf{b}_3, \quad \text{for any } \mathbf{x} \in \mathbb{R}^{N_1},$$

where  $\mathbf{W}_3 = \mathbf{W}_2 \cdot \mathbf{W}_1 \in \mathbb{R}^{N_3 \times N_1}$ ,  $\mathbf{b}_3 = \mathbf{W}_2 \cdot \mathbf{b}_1 + \mathbf{b}_2 \in \mathbb{R}^{N_3}$ .

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

It remains to estimate the approximation error of  $\phi(x) \approx x^2$ . By Equation (4.1), for any  $x \in [0, 1]$ , we have

$$0 \leq \phi(x) - x^2 = f_{Lk}(x) - x^2 \leq 2^{-2(Lk+1)} \leq (2^{2k})^{-L} \leq N^{-L},$$

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

### 4.1.3 Approximation of $x_1 x_2 \cdots x_k$

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

**Lemma 4.3.** *For any  $N, L \in \mathbb{N}^+$ , there exists a function  $\phi$  implemented by a ReLU network with width  $9N$  and depth  $L$  such that*

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

*Proof.* By Lemma 4.2, there exists a function  $\psi$  implemented by a ReLU network with width  $3N$  and depth  $L$  such that

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

Recall the fact

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

The target function  $\phi$  is defined as

$$\phi(x, y) := 2\left(\psi\left(\frac{x+y}{2}\right) - \psi\left(\frac{x}{2}\right) - \psi\left(\frac{y}{2}\right)\right), \quad \text{for any } x, y \in \mathbb{R}. \quad (4.2)$$

Then  $\phi$  can be implemented by the network architecture in Figure 4.5.

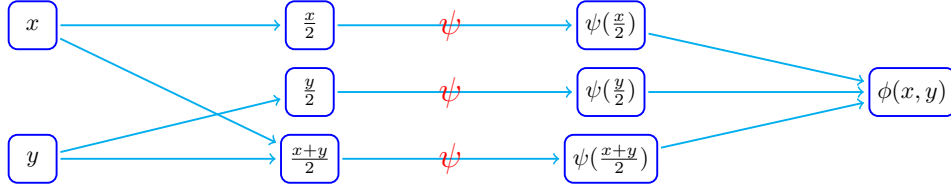


Figure 4.5: An illustration of the network architecture implementing  $\phi$  for approximating  $xy$  on  $[0, 1]^2$ .

It follows from  $\psi \in \mathcal{NN}(\text{width} \leq 3N; \text{depth} \leq L)$  that the network in Figure 4.5 is with width  $9N$  and depth  $L + 2$ . Similar to the discussion in the proof of Lemma 4.2, the network in Figure 4.5 can be interpreted as a ReLU network with width  $9N$  and depth  $L$ , since two of hidden layers have the identify map as their activation functions. Moreover, for any  $x, y \in [0, 1]$ ,

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

Therefore, we have finished the proof.  $\square$

We would like to remark that we can also use Lemma 2.1 to verify the function  $\phi$  defined in Equation (4.2) can be implemented by a ReLU network with width  $9N$  and depth  $L$ , since  $\psi \in \mathcal{NN}(\text{width} \leq 3N; \text{depth} \leq L)$ .

Now let us show how to construct a ReLU network to approximate  $f(x, y) = xy$  on  $[a, b]^2$  with arbitrary  $a < b$ , i.e., a rescaled version of Lemma 4.3.

**Lemma 4.4.** *For any  $N, L \in \mathbb{N}^+$  and  $a, b \in \mathbb{R}$  with  $a < b$ , there exists a function  $\phi$  implemented by a ReLU network with width  $9N + 1$  and depth  $L$  such that*

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

*Proof.* By Lemma 4.3, there exists a function  $\psi$  implemented by a ReLU network with width  $9N$  and depth  $L$  such that

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

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

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

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

$$\left| (b-a)^2 \psi\left(\frac{x-a}{b-a}, \frac{y-a}{b-a}\right) + a(x+y) - a^2 - xy \right| \leq 6(b-a)^2 N^{-L}. \quad (4.3)$$

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

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

Then  $\phi$  can be implemented by the network architecture in Figure 4.6.

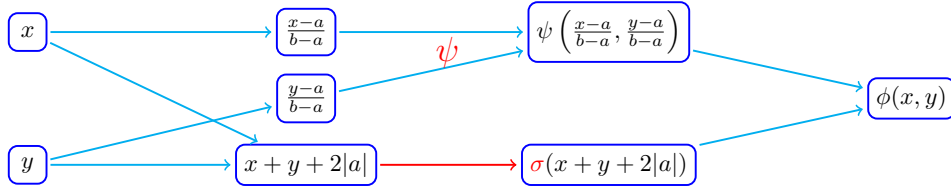


Figure 4.6: An illustration of the network architecture implementing  $\phi$  for approximating  $xy$  on  $[a, b]^2$ . Two of hidden layers has the identify function as their activation functions, since the red “ $\sigma$ ” comes from the red arrow “ $\rightarrow$ ”, where the red arrow “ $\rightarrow$ ” represents a ReLU network with width 1 and depth  $L \geq 1$ .

It follows from  $\psi \in \mathcal{NN}(\text{width} \leq 9N; \text{depth} \leq L)$  that the network in Figure 4.6 is with width  $9N + 1$  and depth  $L + 2$ . Similar to the discussion in the proof of Lemma 4.2, the network in Figure 4.6 can be interpreted as a ReLU network with width  $9N + 1$  and depth  $L$ , since two of hidden layers have the identify function as their activation functions.

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

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



Hence, by Equation (4.3), we have

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

So we finish the proof.  $\square$

The next lemma constructs a ReLU network to approximate a multivariable function  $f(x_1, x_2, \dots, x_k) = x_1 x_2 \cdots x_k$  on  $[0, 1]^k$ .

**Lemma 4.5.** *For any  $N, L, k \in \mathbb{N}^+$  with  $k \geq 2$ , there exists a function  $\phi$  implemented by a ReLU network with width  $9(N + 1) + k - 1$  and depth  $7kL(k - 1)$  such that*

$$|\phi(\mathbf{x}) - x_1 x_2 \cdots x_k| \leq 9(k - 1)(N + 1)^{-7kL},$$

for any  $\mathbf{x} = (x_1, x_2, \dots, x_k) \in [0, 1]^k$ .

*Proof.* By Lemma 4.4, there exists a function  $\phi_1$  implemented by a ReLU network with width  $9(N + 1) + 1$  and depth  $7kL$  such that

$$|\phi_1(x, y) - xy| \leq 6(1.2)^2(N + 1)^{-7kL} \leq 9(N + 1)^{-7kL}, \quad \text{for any } x, y \in [-0.1, 1.1]. \quad (4.4)$$

This equation means the case  $k = 2$  is clear. We may assume  $k \geq 3$  below. We would like to construct a sequence of functions  $\phi_i : [0, 1]^{i+1} \rightarrow [0, 1]$  for any  $i \in \{1, 2, \dots, k - 1\}$  by induction such that

(i)  $\phi_i \in \mathcal{NN}(\text{width} \leq 9(N + 1) + i; \text{depth} \leq 7kLi)$  for any  $i \in \{1, 2, \dots, k - 1\}$ .

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

$$|\phi_i(x_1, x_2, \dots, x_{i+1}) - x_1 x_2 \cdots x_{i+1}| \leq 9i(N + 1)^{-7kL}. \quad (4.5)$$

First, let us consider the case  $i = 1$ . It is obvious that the two required conditions are true: 1)  $9(N + 1) + i = 9(N + 1) + 1$  and  $7kLi = 7kL$  in the case  $i = 1$ ; 2) Equation (4.4) implies Equation (4.5) for  $i = 1$ .

Now assume  $\phi_i$  has been defined for some  $i \in \{1, 2, \dots, k-2\}$ , we define

$$\phi_{i+1}(x_1, x_2, \dots, x_{i+2}) := \phi_1(\phi_i(x_1, x_2, \dots, x_{i+1}), \sigma(x_{i+2})), \quad (4.6)$$

for any  $x_1, x_2, \dots, x_{i+2} \in \mathbb{R}$ . By the induction hypothesis, we have

$$\phi_i \in \mathcal{NN}(\text{width} \leq 9(N+1) + i; \text{depth} \leq 7kLi)$$

and

$$|\phi_i(x_1, x_2, \dots, x_{i+1}) - x_1 x_2 \cdots x_{i+1}| \leq 9i(N+1)^{-7kL}, \quad (4.7)$$

for any  $x_1, x_2, \dots, x_{i+1} \in [0, 1]$ . Clearly,  $\phi_1 \in \mathcal{NN}(\text{width} \leq 9(N+1) + 1; \text{depth} \leq 7kL)$ . Then  $\phi_{i+1}$ , defined in Equation (4.6), can be implemented via a ReLU network with width

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

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

Note that  $9i(N+1)^{-7kL} \leq 9k2^{-7k} \leq 0.1$  for any  $N, L, k \in \mathbb{N}^+$  and  $i \in \{1, 2, \dots, k\}$ .

It follows from Equation (4.7) that

$$\phi_i(x_1, x_2, \dots, x_{i+1}) \in [-0.1, 1.1], \quad \text{for any } x_1, x_2, \dots, x_{i+1} \in [0, 1].$$

Therefore, by Equation (4.4) and (4.7), we have

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

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

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

$$\phi = \phi_{k-1} \in \mathcal{NN}(\text{width} \leq 9(N+1) + k - 1; \text{depth} \leq 7kL(k-1))$$

and

$$\begin{aligned} |\phi(x_1, x_2, \dots, x_k) - x_1 x_2 \cdots x_k| &= |\phi_{k-1}(x_1, x_2, \dots, x_k) - x_1 x_2 \cdots x_k| \\ &\leq 9(k-1)(N+1)^{-7kL}, \end{aligned}$$

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

#### 4.1.4 Proof of main theorem

With Lemma 4.5 in hand, we are ready to prove Theorem 4.1 for approximating general multivariable polynomials by ReLU networks.

*Proof of Theorem 4.1.* The case  $k = 1$  is trivial, so we assume  $k \geq 2$  below. Set  $\tilde{k} = \|\boldsymbol{\alpha}\|_1 \leq k$ , denote  $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_d)$ , and let  $(z_1, z_2, \dots, z_{\tilde{k}}) \in \mathbb{R}^{\tilde{k}}$  be the vector satisfying

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

That is,

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

Then we have  $P(\mathbf{x}) = \mathbf{x}^{\boldsymbol{\alpha}} = x_1^{\alpha_1} x_2^{\alpha_2} \cdots x_d^{\alpha_d} = z_1 z_2 \cdots z_{\tilde{k}}$ .

We construct the target ReLU network in two steps. First, there exists an affine linear map  $\mathcal{L} : \mathbb{R}^d \rightarrow \mathbb{R}^k$  that duplicates  $\mathbf{x}$  to form a new vector

$$(z_1, z_2, \dots, z_{\tilde{k}}, 1, \dots, 1) \in \mathbb{R}^k,$$

i.e.,  $\mathcal{L}(\mathbf{x}) = (z_1, z_2, \dots, z_{\tilde{k}}, 1, \dots, 1) \in \mathbb{R}^k$  for any  $\mathbf{x} \in \mathbb{R}^d$ . Second, by Lemma 4.5, there exists a function  $\psi : \mathbb{R}^k \rightarrow \mathbb{R}$  implemented by a ReLU network with width  $9(N+1)+k-1$  and depth  $7kL(k-1)$  such that  $\psi$  maps  $(z_1, z_2, \dots, z_{\tilde{k}}, 1, \dots, 1) \in \mathbb{R}^k$  to  $z_1 z_2 \dots z_{\tilde{k}}$  within an error  $9(k-1)(N+1)^{-7kL}$  for any  $z_1, z_2, \dots, z_{\tilde{k}} \in [0, 1]$ .

Hence, we can construct our final target function via  $\phi := \psi \circ \mathcal{L}$ . Then by Lemma 2.1 (i),  $\phi$  can be implemented by a ReLU network with width  $9(N+1)+k-1$  and depth  $7kL(k-1) \leq 7k^2L$ , and

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

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

## 4.2 Approximation of continuous functions

In this section, let us focus on constructing ReLU networks to approximate continuous functions on  $[0, 1]^d$ .

### 4.2.1 Main theorem and its proof

Theorem 4.6 below shows that ReLU networks with width  $\mathcal{O}(N)$  and depth  $\mathcal{O}(L)$  can approximate  $f \in C([0, 1]^d)$  with an approximation error  $19\sqrt{d}\omega_f(N^{-2/d}L^{-2/d})$ .

**Theorem 4.6.** *Given a continuous function  $f \in C([0, 1]^d)$ , for any  $N, L \in \mathbb{N}^+$  and  $p \in [1, \infty]$ , there exists a function  $\phi$  implemented by a ReLU network with width  $C_1 \max\{d\lfloor N^{1/d} \rfloor, N+1\}$  and depth  $12L + C_2$  such that*

$$\|\phi - f\|_{L^p([0, 1]^d)} \leq 19\sqrt{d}\omega_f(N^{-2/d}L^{-2/d}),$$

where  $C_1 = 12$  and  $C_2 = 14$  if  $p \in [1, \infty)$ ;  $C_1 = 3^{d+3}$  and  $C_2 = 14 + 2d$  if  $p = \infty$ .

The approximation error becomes  $19\sqrt{d}\lambda N^{-2\alpha/d}L^{-2\alpha/d}$  when Theorem 4.6 is applied to a function  $f \in \text{Hölder}([0, 1]^d, \alpha, \lambda)$  as shown in the corollary below, where  $\text{Hölder}([0, 1]^d, \alpha, \lambda)$  is the space of Hölder continuous functions of order  $\alpha \in (0, 1]$  with a Hölder constant  $\lambda > 0$ .

**Corollary 4.7.** *Given a function  $f \in \text{Hölder}([0, 1]^d, \alpha, \lambda)$ , for any  $N, L \in \mathbb{N}^+$  and  $p \in [1, \infty]$ , there exists a function  $\phi$  implemented by a ReLU network with width  $C_1 \max\{d\lfloor N^{1/d} \rfloor, N + 1\}$  and depth  $12L + C_2$  such that*

$$\|\phi - f\|_{L^p([0, 1]^d)} \leq 19\sqrt{d}\lambda N^{-2\alpha/d}L^{-2\alpha/d},$$

where  $C_1 = 12$  and  $C_2 = 14$  if  $p \in [1, \infty)$ ;  $C_1 = 3^{d+3}$  and  $C_2 = 14 + 2d$  if  $p = \infty$ .

The next question is: How much we can improve the approximation error in Theorem 4.6 and Corollary 4.7? In fact, for the Hölder continuous function space, the approximation error in Corollary 4.7 is nearly optimal based on VC-dimension as we shall see later in Section 4.4.

To prove Theorem 4.6, we introduce a theorem below, Theorem 4.8, to construct ReLU networks to uniformly approximate continuous functions outside the trifling region, which can deduce Theorem 4.6 easily by applying Theorem 3.7.

**Theorem 4.8.** *Given a continuous function  $f \in C([0, 1]^d)$ , for any  $N, L \in \mathbb{N}^+$ , there exists a function  $\phi$  implemented by a ReLU network with width  $\max\{4d\lfloor N^{1/d} \rfloor + 3d, 12N + 8\}$  and depth  $12L + 14$  such that  $\|\phi\|_{L^\infty(\mathbb{R}^d)} \leq |f(\mathbf{0})| + \omega_f(\sqrt{d})$  and*

$$|\phi(\mathbf{x}) - f(\mathbf{x})| \leq 18\sqrt{d}\omega_f(N^{-2/d}L^{-2/d}), \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta),$$

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

Now we are ready to prove Theorem 4.6 by assuming Theorem 4.8 is true, which will be proved later in this section.

*Proof of Theorem 4.6.* Let us first consider the case  $p \in [1, \infty)$ . We may assume  $f$  is not a constant function since it is a trivial case. Then  $\omega_f(r) > 0$  for any  $r > 0$ .

Set  $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor$  and choose a small  $\delta \in (0, \frac{1}{3K}]$  such that

$$\begin{aligned} Kd\delta(2|f(\mathbf{0})| + 2\omega_f(\sqrt{d}))^p &= \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor d\delta(2|f(\mathbf{0})| + 2\omega_f(\sqrt{d}))^p \\ &\leq (\omega_f(N^{-2/d}L^{-2/d}))^p. \end{aligned} \quad (4.8)$$

By Theorem 4.8, there exists a function  $\phi$  implemented by a ReLU network with width

$$\max \{4d\lfloor N^{1/d} \rfloor + 3d, 12N + 8\} \leq 12 \max \{d\lfloor N^{1/d} \rfloor, N + 1\}$$

and depth  $12L + 14$  such that  $\|\phi\|_{L^\infty(\mathbb{R}^d)} \leq |f(\mathbf{0})| + \omega_f(\sqrt{d})$  and

$$|f(\mathbf{x}) - \phi(\mathbf{x})| \leq 18\sqrt{d}\omega_f(N^{-2/d}L^{-2/d}), \quad \text{for any } \mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta).$$

Note that  $\mu(\Omega([0, 1]^d, K, \delta)) \leq Kd\delta$  and  $\|f\|_{L^\infty([0, 1]^d)} \leq |f(\mathbf{0})| + \omega_f(\sqrt{d})$ . Then, by Equation (4.8), we have

$$\begin{aligned} \|f - \phi\|_{L^p([0, 1]^d)}^p &= \int_{\Omega([0, 1]^d, K, \delta)} |f(\mathbf{x}) - \phi(\mathbf{x})|^p d\mathbf{x} + \int_{[0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)} |f(\mathbf{x}) - \phi(\mathbf{x})|^p d\mathbf{x} \\ &\leq Kd\delta(2|f(\mathbf{0})| + 2\omega_f(\sqrt{d}))^p + (18\sqrt{d}\omega_f(N^{-2/d}L^{-2/d}))^p \\ &\leq (\omega_f(N^{-2/d}L^{-2/d}))^p + (18\sqrt{d}\omega_f(N^{-2/d}L^{-2/d}))^p \\ &\leq (19\sqrt{d}\omega_f(N^{-2/d}L^{-2/d}))^p. \end{aligned}$$

Hence,  $\|f - \phi\|_{L^p([0, 1]^d)} \leq 19\sqrt{d}\omega_f(N^{-2/d}L^{-2/d})$ .

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

$$d \cdot \omega_f(\delta) \leq \omega_f(N^{-2/d}L^{-2/d}).$$

By Theorem 4.8, there exists a function implemented  $\tilde{\phi}$  by a ReLU network with width  $\max \{4d\lfloor N^{1/d} \rfloor + 3d, 12N + 8\}$  and depth  $12L + 14$  such that

$$|f(\mathbf{x}) - \tilde{\phi}(\mathbf{x})| \leq 18\sqrt{d}\omega_f(N^{-2/d}L^{-2/d}) =: \varepsilon, \quad \text{for } \mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta).$$

By Theorem 3.7, there exists a function  $\phi$  implemented by a ReLU network with width

$$3^d \left( \max \{4d \lfloor N^{1/d} \rfloor + 3d, 12N + 8\} + 4 \right) \leq 3^{d+3} \max \{d \lfloor N^{1/d} \rfloor, N + 1\}$$

and depth  $12L + 14 + 2d$  such that

$$|f(\mathbf{x}) - \phi(\mathbf{x})| \leq \varepsilon + d \cdot \omega_f(\delta) \leq 19\sqrt{d} \omega_f(N^{-2/d} L^{-2/d}), \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

So we finish the proof.  $\square$

It remains to prove Theorem 4.8. To this end, we establish a proposition below to warm up the proof of Theorem 4.8.

**Proposition 4.9.** *For any  $\varepsilon > 0$  and arbitrary  $N, L, J \in \mathbb{N}^+$  with  $J \leq N^2 L^2$ , given  $J$  samples  $y_j \geq 0$  for  $j = 0, 1, \dots, J-1$  with*

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

*Then there exists  $\phi \in \mathcal{NN}(\#input = 1; \text{width} \leq 12N+8; \text{depth} \leq 4L+9; \#output = 1)$  such that*

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

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

### 4.2.2 Proof of auxiliary theorem

We essentially construct an almost piecewise constant function implemented by a ReLU network with width  $\mathcal{O}(N)$  and depth  $\mathcal{O}(L)$  to approximate the target continuous function  $f$  on  $[0, 1]^d$ . We assume  $f$  is not a constant since it is a trivial case. Then  $\omega_f(r) > 0$  for any  $r > 0$ . It is 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  $0 \leq \tilde{f}(\mathbf{x}) \leq 2\omega_f(\sqrt{d})$  for any

$\mathbf{x} \in [0, 1]^d$ . Let  $M = N^2L$ ,  $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor$ , and  $\delta$  be an arbitrary number in  $(0, \frac{1}{3K}]$ . The proof can be divided into four steps as follows.

1. 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  $\mathbf{x}_\beta$  as the vertex of  $Q_\beta$  with minimum  $\|\cdot\|_1$  norm for each  $\beta \in \{0, 1, \dots, K-1\}^d$ . See Figure 4.7 for the illustrations of  $\Omega([0, 1]^d, K, \delta)$ ,  $Q_\beta$ , and  $\mathbf{x}_\beta$  for any  $\beta \in \{0, 1, \dots, K-1\}^d$ .
2. Construct a sub-network to implement a vector-valued function  $\Phi_1$  projecting the whole cube  $Q_\beta$  to the  $d$ -dimensional index  $\beta$  for each  $\beta$ , *i.e.*,  $\Phi(\mathbf{x}) = \beta$  for all  $\mathbf{x} \in Q_\beta$  and each  $\beta \in \{0, 1, \dots, K-1\}^d$ .
3. Construct a sub-network to implement a function  $\phi_2$  mapping the index  $\beta$  approximately to  $\tilde{f}(\mathbf{x}_\beta)$  for each  $\beta$ . This step can be further divided into three sub-steps.
  - 3.1. Construct an affine linear map  $\psi_1 : \mathbb{R}^d \rightarrow \mathbb{R}$  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 4.8 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 two conditions.
    - Assign the value of  $g$  at  $\psi(\beta) \in \mathcal{A}_1$  as  $\tilde{f}(\mathbf{x}_\beta)$ , *i.e.*,  $g \circ \psi_1(\beta) = \tilde{f}(\mathbf{x}_\beta)$  for each  $\beta \in \{0, 1, \dots, K-1\}^d$ .
    - The values of  $g$  at breakpoints in  $\mathcal{A}_2 \cup \{1\}$  are properly assigned for applying Proposition 4.9.
  - 3.3. Apply Proposition 4.9 to construct a sub-network to implement a function  $\psi_2$  approximating  $g$  well on  $\mathcal{A}_1 \cup \mathcal{A}_2 \cup \{1\}$ . Then  $\phi_2 = \psi_2 \circ \psi_1$  satisfies  $\phi_2(\beta) = \psi_2 \circ \psi_1(\beta) \approx g \circ \psi_1(\beta) = \tilde{f}(\mathbf{x}_\beta)$  for each  $\beta \in \{0, 1, \dots, K-1\}^d$ .
4. Construct the final target network to implement the desired function  $\phi$  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  $\mathbf{x} \in Q_\beta$  and each  $\beta \in \{0, 1, \dots, K-1\}^d$ .



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

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

for each  $\beta = (\beta_1, \beta_2, \dots, \beta_d) \in \{0, 1, \dots, K-1\}^d$ . Recall that  $\Omega([0, 1]^d, K, \delta)$  is the trifling region defined in Equation (2.1). Apparently,  $\mathbf{x}_\beta$  is the vertex of  $Q_\beta$  with minimum  $\|\cdot\|_1$  norm and

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

See Figure 4.7 for illustrations of  $\Omega([0, 1]^d, K, \delta)$ ,  $Q_\beta$ , and  $\mathbf{x}_\beta$  for any  $\beta \in \{0, 1, \dots, K-1\}^d$ .

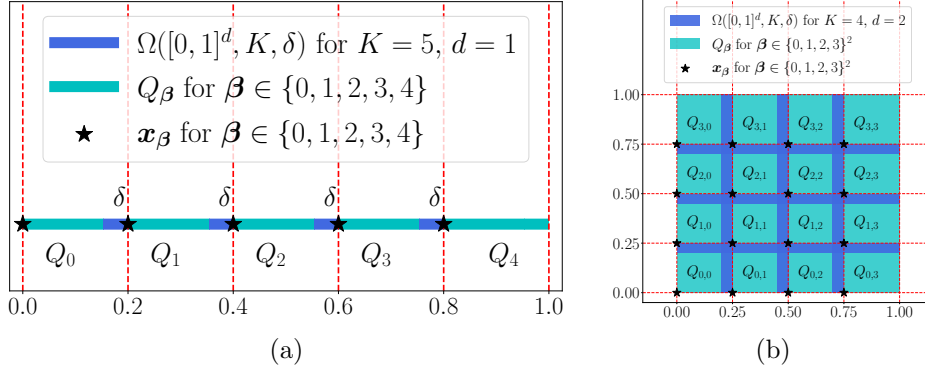


Figure 4.7: Illustrations of  $\Omega([0, 1]^d, K, \delta)$ ,  $Q_\beta$ , and  $\mathbf{x}_\beta$  for any  $\beta \in \{0, 1, \dots, K-1\}^d$ . (a)  $K = 5$  and  $d = 1$ . (b)  $K = 4$  and  $d = 2$ .

**Step 2:** Construct  $\Phi_1$  mapping  $\mathbf{x} \in Q_\beta$  to  $\beta$ .

By Theorem 3.12, there exists  $\phi_1 \in \mathcal{NN}(\text{width} \leq 4\lfloor N^{1/d} \rfloor + 3; \text{depth} \leq 4L + 5)$  and

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

By defining

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

we have  $\Phi_1(\mathbf{x}) = (\phi_1(x_1), \phi_1(x_2), \dots, \phi_1(x_d)) = \beta$  for all  $\mathbf{x} \in Q_\beta$  and each  $\beta \in \{0, 1, \dots, K-1\}^d$ .

**Step 3:** Construct  $\phi_2$  mapping  $\beta$  approximately to  $\tilde{f}(\mathbf{x}_\beta)$ .

The construction of the sub-network implementing  $\phi_2$  is essentially based on Proposition 4.9. To meet the requirements of applying Proposition 4.9, we first define two auxiliary sets  $\mathcal{A}_1$  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,  $\mathcal{A}_1 \cup \mathcal{A}_2 \cup \{1\} = \{\frac{j}{2K^d} : j = 0, 1, \dots, 2K^d\}$  and  $\mathcal{A}_1 \cap \mathcal{A}_2 = \emptyset$ . See Figure 4.8 for an illustration of  $\mathcal{A}_1$  and  $\mathcal{A}_2$ . Next, we divide this step into three sub-steps.

**Step 3.1:** Construct  $\psi_1$  bijectively mapping  $\{0, 1, \dots, K-1\}^d$  to  $\mathcal{A}_1$ .

Inspired by the binary representation, we define

$$\psi_1(\mathbf{x}) := \frac{x_d}{2K^d} + \sum_{i=1}^{d-1} \frac{x_i}{K^i}, \quad \text{for any } \mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d. \quad (4.9)$$

Then  $\psi_1$  is an affine linear function bijectively mapping the index set  $\{0, 1, \dots, K-1\}^d$

$1\}^d$  to

$$\begin{aligned} & \left\{ \frac{\beta_d}{2K^d} + \sum_{i=1}^{d-1} \frac{\beta_i}{K^i} : \boldsymbol{\beta} = (\beta_1, \dots, \beta_d) \in \{0, 1, \dots, K-1\}^d \right\} \\ &= \left\{ \frac{i}{K^{d-1}} + \frac{k}{2K^d} : i = 0, 1, \dots, K^{d-1} - 1 \quad \text{and} \quad k = 0, 1, \dots, K-1 \right\} = \mathcal{A}_1. \end{aligned}$$

**Step 3.2:** Construct  $g$  to satisfy  $g \circ \psi_1(\boldsymbol{\beta}) = \tilde{f}(\mathbf{x}_{\boldsymbol{\beta}})$  and to meet the requirements of applying Proposition 4.9.

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

$$g(\psi_1(\boldsymbol{\beta})) = \tilde{f}(\mathbf{x}_{\boldsymbol{\beta}}), \quad \text{for any } \boldsymbol{\beta} \in \{0, 1, \dots, K-1\}^d. \quad (4.10)$$

- At the breakpoint 1, let  $g(1) = \tilde{f}(\mathbf{1})$ , where  $\mathbf{1} = (1, 1, \dots, 1) \in \mathbb{R}^d$ .
- 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 4.9. 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  $g$  be 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 4.8 for an example) and satisfies

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

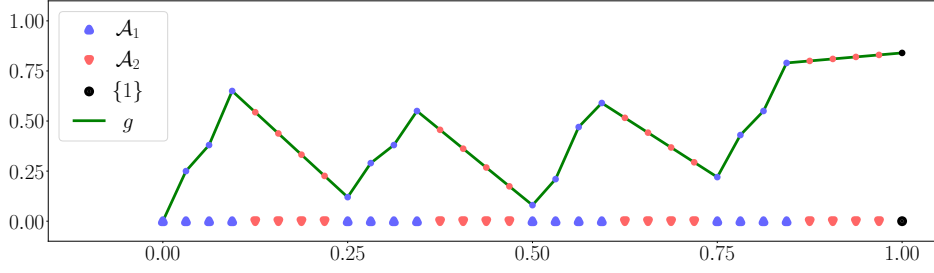


Figure 4.8: An illustration of  $\mathcal{A}_1$ ,  $\mathcal{A}_2$ ,  $\{1\}$ , and  $g$  for the case  $d = 2$  and  $K = 4$ .

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  $\psi_2$  approximating  $g$  well on  $\mathcal{A}_1 \cup \mathcal{A}_2 \cup \{1\}$ .

Since  $2K^d = 2(\lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor)^d \leq 2(N^2 L^2) \leq N^2 \tilde{L}^2$ , where  $\tilde{L} = 2L$ , by Proposition 4.9 (set  $y_j = g(\frac{j}{2K^d})$  and  $\varepsilon = \omega_f(\frac{\sqrt{d}}{K}) > 0$  therein), there exists  $\tilde{\psi}_2 \in \mathcal{NN}(\text{width} \leq 12N + 8; \text{depth} \leq 4\tilde{L} + 9) = \mathcal{NN}(\text{width} \leq 12N + 8; \text{depth} \leq 8L + 9)$  such that

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

and

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

Define  $\psi_2(x) := \tilde{\psi}_2(2K^d x)$  for any  $x \in \mathbb{R}$ . Then, we have  $\psi_2 \in \mathcal{NN}(\text{width} \leq 12N + 8; \text{depth} \leq 8L + 9)$  by Lemma 2.1 (i),

$$0 \leq \psi_2(x) = \tilde{\psi}_2(2K^d x) \leq 2\omega_f(\sqrt{d}), \quad \text{for any } x \in \mathbb{R}, \quad (4.11)$$

and

$$|\psi_2(\frac{j}{2K^d}) - g(\frac{j}{2K^d})| = |\tilde{\psi}_2(j) - g(\frac{j}{2K^d})| \leq \omega_f(\frac{\sqrt{d}}{K}), \quad (4.12)$$

for  $j = 0, 1, \dots, 2K^d - 1$ .

The desired function  $\phi_2$  can be defined as  $\phi_2 := \psi_2 \circ \psi_1$ . Note that  $\psi_1 : \mathbb{R}^d \rightarrow \mathbb{R}$  is

an affine linear map and  $\psi_2 \in \mathcal{NN}(\#input = 1; \text{width} \leq 12N + 8; \text{depth} \leq 8L + 9)$ . Thus, by Lemma 2.1 (i),  $\phi_2 = \psi_2 \circ \psi_1 \in \mathcal{NN}(\#input = d; \text{width} \leq 12N + 8; \text{depth} \leq 8L + 9)$ . By Equation (4.10) and (4.12), we have

$$|\phi_2(\boldsymbol{\beta}) - \tilde{f}(\mathbf{x}_{\boldsymbol{\beta}})| = |\psi_2(\psi_1(\boldsymbol{\beta})) - g(\psi_1(\boldsymbol{\beta}))| \leq \omega_f\left(\frac{\sqrt{d}}{K}\right), \quad (4.13)$$

for any  $\boldsymbol{\beta} \in \{0, 1, \dots, K-1\}^d$ . Equation (4.11) and  $\phi_2 = \psi_2 \circ \psi_1$  implies

$$0 \leq \phi_2(\mathbf{x}) \leq 2\omega_f(\sqrt{d}), \quad \text{for any } \mathbf{x} \in \mathbb{R}^d. \quad (4.14)$$

**Step 4:** Construct the final network to implement the desired function  $\phi$ .

Define  $\phi := \phi_2 \circ \Phi_1 + f(0) - \omega_f(\sqrt{d})$ .

$$\phi_1 \in \mathcal{NN}(\#input = 1; \text{width} \leq 4\lfloor N^{1/d} \rfloor + 3; \text{depth} \leq 4L + 5; \#output = 1),$$

implies

$$\Phi_1 \in \mathcal{NN}(\#input = d; \text{width} \leq 4d\lfloor N^{1/d} \rfloor + 3d; \text{depth} \leq 4L + 5; \#output = d).$$

Note that  $\phi_2 \in \mathcal{NN}(\#input = d; \text{width} \leq 12N + 8; \text{depth} \leq 8L + 9)$ . Thus, by Lemma 2.1,  $\phi = \phi_2 \circ \Phi_1 + f(\mathbf{0}) - \omega_f(\sqrt{d})$  is in

$$\mathcal{NN}(\text{width} \leq \max\{4d\lfloor N^{1/d} \rfloor + 3d, 12N + 8\}; \text{depth} \leq (4L + 5) + (8L + 9) = 12L + 14).$$

Finally, let us estimate the approximation error. Recall that  $f = \tilde{f} + f(\mathbf{0}) - \omega_f(\sqrt{d})$ . By Equation (4.13), for any  $\mathbf{x} \in Q_{\boldsymbol{\beta}}$  and each  $\boldsymbol{\beta} \in \{0, 1, \dots, K-1\}^d$ , we

have

$$\begin{aligned}
|f(\mathbf{x}) - \phi(\mathbf{x})| &= |\tilde{f}(\mathbf{x}) - \phi_2 \circ \Phi_1(\mathbf{x})| = |\tilde{f}(\mathbf{x}) - \phi_2(\boldsymbol{\beta})| \\
&\leq |\tilde{f}(\mathbf{x}) - \tilde{f}(\mathbf{x}_\beta)| + |\tilde{f}(\mathbf{x}_\beta) - \phi_2(\boldsymbol{\beta})| \\
&\leq \omega_f\left(\frac{\sqrt{d}}{K}\right) + \omega_f\left(\frac{\sqrt{d}}{K}\right) \leq 2\omega_f(8\sqrt{d}N^{-2/d}L^{-2/d}),
\end{aligned}$$

where the last inequality comes from the fact  $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor \geq \frac{N^{2/d}L^{2/d}}{8}$  for any  $N, L \in \mathbb{N}^+$ . Recall the fact  $\omega_f(nr) \leq n\omega_f(r)$  for any  $n \in \mathbb{N}^+$  and  $r \in [0, \infty)$ . Therefore, for any  $\mathbf{x} \in \cup_{\beta \in \{0,1,\dots,K-1\}^d} Q_\beta = [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)$ , we have

$$\begin{aligned}
|f(\mathbf{x}) - \phi(\mathbf{x})| &\leq 2\omega_f(8\sqrt{d}N^{-2/d}L^{-2/d}) \leq 2\lceil 8\sqrt{d} \rceil \omega_f(N^{-2/d}L^{-2/d}) \\
&\leq 18\sqrt{d}\omega_f(N^{-2/d}L^{-2/d})
\end{aligned}$$

It remains to show the upper bound of  $\phi$ . By Equation (4.14) and  $\phi = \phi_2 \circ \Phi_1 + f(\mathbf{0}) - \omega_f(\sqrt{d})$ , we have  $\|\phi\|_{L^\infty(\mathbb{R}^d)} \leq |f(\mathbf{0})| + \omega_f(\sqrt{d})$ . Thus, we finish the proof.

### 4.2.3 Proof of key proposition for auxiliary theorem

Let us prove Proposition 4.9 to end Section 4.2. We apply Theorem 3.5 to prove Lemma 4.10 below, which is a key intermediate conclusion to prove Proposition 4.9.

**Lemma 4.10.** *For any  $\varepsilon > 0$  and  $N, L \in \mathbb{N}^+$ , denote  $M = N^2L$  and assume  $y_{m,\ell} \geq 0$  for  $m = 0, 1, \dots, M-1$  and  $\ell = 0, 1, \dots, L-1$  are samples with*

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

*Then there exists  $\phi \in \mathcal{NN}(\#\text{input} = 2; \text{width} \leq 12N+8; \text{depth} \leq 3L+6; \#\text{output} = 1)$  such that*

$$(i) \quad |\phi(m, \ell) - y_{m,\ell}| \leq \varepsilon \text{ for } m = 0, 1, \dots, M-1 \quad \text{and} \quad \ell = 0, 1, \dots, L-1.$$

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

*Proof.* Define

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

We will construct a function implemented by a ReLU network to map the index  $(m, \ell)$  to  $a_{m,\ell}\varepsilon = \lfloor y_{m,\ell}/\varepsilon \rfloor \varepsilon \approx y_{m,\ell}$  for  $m = 0, 1, \dots, M-1$  and  $\ell = 0, 1, \dots, L-1$ .

Define  $b_{m,0} := 0$  and  $b_{m,\ell} := a_{m,\ell} - a_{m,\ell-1}$  for  $m = 0, 1, \dots, M-1$  and  $\ell = 1, \dots, L-1$ . Since  $|y_{m,\ell} - y_{m,\ell-1}| \leq \varepsilon$  for all  $m$  and  $\ell$ , we have

$$b_{m,\ell} = a_{m,\ell} - a_{m,\ell-1} = \lfloor y_{m,\ell}/\varepsilon \rfloor - \lfloor y_{m,\ell-1}/\varepsilon \rfloor \in \{-1, 0, 1\}.$$

Hence, there exist  $c_{m,\ell}$  and  $d_{m,\ell} \in \{0, 1\}$  such that  $b_{m,\ell} = c_{m,\ell} - d_{m,\ell}$ , which implies

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

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

Consider the sample set

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

Its size is  $M+1 = N \cdot ((NL-1)+1) + 1$ . By Theorem 3.2 (set  $m = N$  and  $n = NL-1$  therein), there exists

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

such that

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

By Theorem 3.5, there exist  $\psi_2, \psi_3 \in \mathcal{NN}(\text{width} \leq 4N+3; \text{depth} \leq 3L+3)$

such that

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

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

$$a_{m,\ell} = a_{m,0} + \sum_{j=0}^{\ell} c_{m,j} - \sum_{j=0}^{\ell} d_{m,j} = \psi_1(m) + \psi_2(m, \ell) - \psi_3(m, \ell), \quad (4.15)$$

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

Define

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

Recall that, for any  $x_1, x_2 \in \mathbb{R}$ , we have

$$\min\{x_1, x_2\} = \frac{x_1 + x_2 - |x_1 - x_2|}{2} = \frac{\sigma(x_1 + x_2) - \sigma(-x_1 - x_2) - \sigma(x_1 - x_2) - \sigma(-x_1 + x_2)}{2}. \quad (4.16)$$

Then the desired function can be implemented by the composition of two sub-networks shown in Figure 4.9.

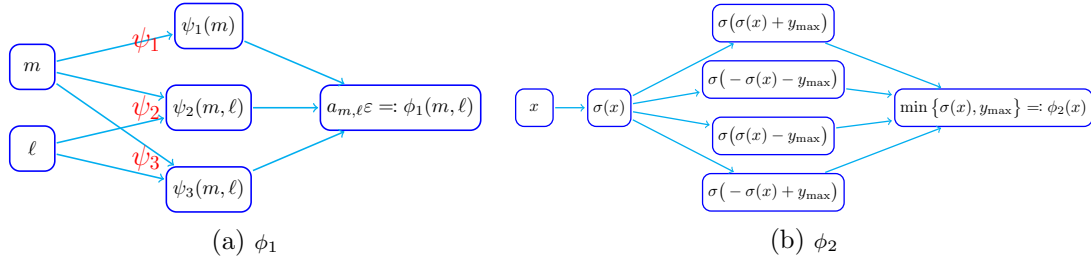


Figure 4.9: Illustrations of two sub-network architectures for implementing the desired function  $\phi = \phi_2 \circ \phi_1$  based on Equation (4.15) and (4.16) for  $m = 0, 1, \dots, M-1$  and  $\ell = 0, 1, \dots, L-1$ .

By Theorem 3.1,  $\psi_1 \in \mathcal{NN}(\text{widthvec} = [2N, 2NL - 1]) \subseteq \mathcal{NN}(\text{width} \leq 4N + 2; \text{depth} \leq L + 1)$ . Note that  $\psi_2, \psi_3 \in \mathcal{NN}(\text{width} \leq 4N + 3; \text{depth} \leq 3L + 3)$ . Thus,  $\phi_1 \in \mathcal{NN}(\text{width} \leq (4N + 2) + 2(4N + 3) = 12N + 8; \text{depth} \leq (3L + 3) + 1 = 3L + 4)$  as shown in Figure 4.9. It is clear that  $\phi_2 \in \mathcal{NN}(\text{width} \leq 4; \text{depth} \leq 2)$ , implying



$\phi = \phi_2 \circ \phi_1 \in \mathcal{NN}(\text{width} \leq 12N + 8; \text{depth} \leq (3L + 4) + 2 = 3L + 6)$  by Lemma 2.1 (ii).

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

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

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

With Lemma 4.10 in hand, we are ready to prove Proposition 4.9.

*Proof of Proposition 4.9.* Let  $M = N^2L$ , then we may assume  $J = ML$  since we can set  $y_{J-1} = y_J = y_{J+1} = \dots = y_{ML-1}$  if  $J < ML$ .

Consider the sample set

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

whose size is  $2M + 1 = N \cdot ((2NL - 1) + 1) + 1$ . By Theorem 3.2 (set  $m = N$  and  $n = NL - 1$  therein), there exist

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

such that

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

It follows that

$$\phi_1(j) = m, \quad \text{and} \quad j - L\phi_1(j) = \ell, \quad \text{where } j = mL + \ell, \quad (4.17)$$

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

Note that any number  $j$  in  $\{0, 1, \dots, J-1\}$  can be uniquely indexed as  $j = mL + \ell$  for  $m = 0, 1, \dots, M-1$  and  $\ell = 0, 1, \dots, L-1$ . So we can denote  $y_j = y_{mL+\ell}$  by  $y_{m,\ell}$ . Then by Lemma 4.10, there exists  $\phi_2 \in \mathcal{NN}(\text{width} \leq 12N + 8; \text{depth} \leq 3L + 6)$  such that

$$|\phi_2(m, \ell) - y_{m,\ell}| \leq \varepsilon, \quad \text{for } m = 0, 1, \dots, M-1 \text{ and } \ell = 0, 1, \dots, L-1, \quad (4.18)$$

and

$$0 \leq \phi_2(x_1, x_2) \leq y_{\max}, \quad \text{for any } x_1, x_2 \in \mathbb{R}, \quad (4.19)$$

where  $y_{\max} := \max\{y_{m,\ell} : m = 0, 1, \dots, M-1 \text{ and } \ell = 0, 1, \dots, L-1\} = \max\{y_j : j = 0, 1, \dots, ML-1\}$ . Then the desired function  $\phi$  can be implemented by the network in Figure 4.10.

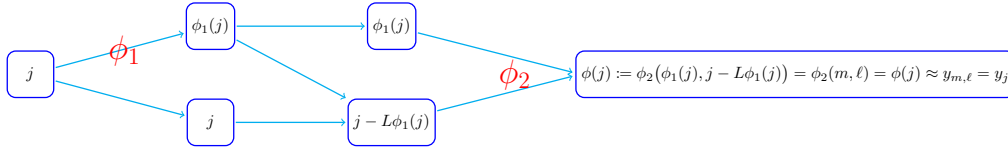


Figure 4.10: An illustration of the network architecture implementing the desired function  $\phi$  based Equation (4.17). The index  $j \in \{0, 1, \dots, ML-1\}$  is unique represented by  $j = mL + \ell$  for  $m = 0, 1, \dots, M-1$  and  $\ell = 0, 1, \dots, L-1$ .

Note that  $\phi_1 \in \mathcal{NN}(\text{widthvec} = [2N, 4NL-1]) \subseteq \mathcal{NN}(\text{width} \leq 8N+2; \text{depth} \leq L+1)$  by Theorem 3.1 and  $\phi_2 \in \mathcal{NN}(\text{width} \leq 12N+8; \text{depth} \leq 3L+6)$ . So  $\phi \in \mathcal{NN}(\text{width} \leq \max\{8N+2+1, 12N+8\} = 12N+8; \text{depth} \leq (L+1)+2+(3L+6) = 4L+9)$  as shown in Figure 4.10.

By Equation (4.19) and Figure 4.10, we have

$$0 \leq \phi(x) \leq y_{\max} = \max\{y_j : j = 0, 1, \dots, ML-1\}, \quad \text{for any } x \in \mathbb{R}.$$

Represent  $j \in \{0, 1, \dots, ML-1\}$  via  $j = mL + \ell$  for  $m = 0, 1, \dots, M-1$  and

$\ell = 0, 1, \dots, L - 1$ , then we have, by Equation (4.18),

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

So we finish the proof.  $\square$

We would like to remark that the key idea in the proof of Proposition 4.9 is the “bit extraction” technique, which allows us to store  $L$  bits in a binary number  $\text{bin}0.\theta_1\theta_2\cdots\theta_L$  and extract each bit  $\theta_i$ . The extraction operator can be efficiently carried out via a deep ReLU neural network, demonstrating the power of depth.

## 4.3 Approximation of smooth functions

In Section 4.2, we show that the approximation of a function  $f \in C([0, 1]^d)$ , by ReLU networks with width  $\mathcal{O}(N)$  and depth  $\mathcal{O}(L)$ , admits an approximation error  $19\sqrt{d}\omega_f(N^{-2/d}L^{-2/d})$  in the  $L^p$ -norm for  $p \in [1, \infty]$ . The next question is whether the smoothness of functions can improve the approximation error. In this section, we investigate the approximation of smooth functions by ReLU networks.

### 4.3.1 Main theorem and its proof

Theorem 4.11 below shows that ReLU networks with width  $\mathcal{O}(N \ln N)$  and depth  $\mathcal{O}(L \ln L)$  can approximate a function  $f \in C^s([0, 1]^d)$  with a nearly optimal approximation error  $\mathcal{O}(\|f\|_{C^s([0, 1]^d)} N^{-2s/d} L^{-2s/d})$ . See Section 4.4.2 for the optimality discussion.

**Theorem 4.11.** *Given a smooth function  $f \in C^s([0, 1]^d)$  with  $s \in \mathbb{N}^+$ , for any  $N, L \in \mathbb{N}^+$ , there exists a function  $\phi$  implemented by a ReLU network with width  $C_1(N + 2) \log_2(8N)$  and depth  $C_2(L + 2) \log_2(4L) + 2d$  such that*

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

where  $C_1 = 17s^{d+1}3^d$ ,  $C_2 = 18s^2$ , and  $C_3 = 85(s+1)^d8^s$ .

As we can see from Theorem 4.11, the smoothness improves the approximation error in  $N$  and  $L$ . However, we would like to remark that the improved approximation error is at the price of much larger constants.

In Theorem 4.11, the logarithmic terms in width and depth can be further reduced if the approximation error is weakened. Note that for any  $\tilde{N}, \tilde{L} \in \mathbb{N}^+$  with

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

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

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

and

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

It follows that

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

and

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

Thus, we have an immediate corollary.

**Corollary 4.12.** *Given a function  $f \in C^s([0,1]^d)$  with  $s \in \mathbb{N}^+$ , for any  $\tilde{N}, \tilde{L} \in \mathbb{N}^+$ , there exist a function  $\phi$  implemented by a ReLU network with width  $\tilde{N}$  and depth  $\tilde{L}$  such that*

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

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

To prove Theorem 4.11, we first introduce Theorem 4.13, a simplified version of Theorem 4.11 ignoring the approximation error in the trifling region  $\Omega([0, 1]^d, K, \delta)$ . Then Theorem 4.11 can be easily proved by combining Theorem 3.7 and 4.13 together. Recall that  $C_u^s([0, 1]^d)$  is the closed unit ball of  $C^s([0, 1]^d)$ .

**Theorem 4.13.** *Given a smooth function  $f \in C_u^s([0, 1]^d)$ , for any  $N, L \in \mathbb{N}^+$ , there exists a function  $\phi$  implemented by ReLU network with width  $16s^{d+1}d(N + 2)\log_2(8N)$  and depth  $18s^2(L + 2)\log_2(4L)$  such that*

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

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

Theorem 4.13 will be proved in Section 4.3.3. By assuming Theorem 4.13 is true, we can prove Theorem 4.11 based on Theorem 3.7.

*Proof of Theorem 4.11.* We may assume  $\|f\|_{C^s([0, 1]^d)} > 0$  since  $\|f\|_{C^s([0, 1]^d)} = 0$  is a trivial case. Define  $\tilde{f} := \frac{f}{\|f\|_{C^s([0, 1]^d)}} \in C_u^s([0, 1]^d)$ , set  $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor$ , and choose a small  $\delta \in (0, \frac{1}{3K}]$  such that

$$d \cdot \omega_f(\delta) \leq N^{-2s/d} L^{-2s/d}.$$

By Theorem 4.13, there exists a function  $\hat{\phi}$  implemented by a ReLU network with width  $16s^{d+1}d(N + 2)\log_2(8N)$  and depth  $18s^2(L + 2)\log_2(4L)$  such that

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

By Theorem 3.7, there exists a new function  $\tilde{\phi}$  implemented by a ReLU network

with width

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

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

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

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

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

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

It remains to prove Theorem 4.13, a weaker version of Theorem 4.11 targeting a ReLU network constructed to approximate a smooth function outside the trifling region. We discuss the ideas of the proof in Section 4.3.2 and give the detailed proof in Section 4.3.3.

### 4.3.2 Ideas of proving auxiliary theorem

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

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

Clearly,  $[0,1]^d = \Omega([0,1]^d, K, \delta) \cup \left( \bigcup_{\beta \in \{0,1,\dots,K-1\}^d} Q_\beta \right)$  and  $\mathbf{x}_\beta$  is the vertex of  $Q_\beta$  with minimum  $\|\cdot\|_1$  norm. See Figure 4.11 for the illustrations of  $\Omega([0,1]^d, K, \delta)$ ,

$Q_\beta$ , and  $\mathbf{x}_\beta$  for any  $\beta \in \{0, 1, \dots, K-1\}^d$ .

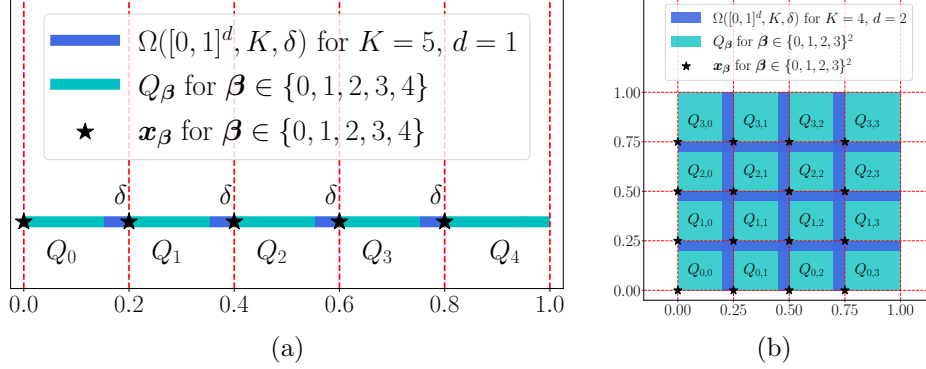


Figure 4.11: Illustrations of  $\Omega([0, 1]^d, K, \delta)$ ,  $Q_\beta$ , and  $\mathbf{x}_\beta$  for any  $\beta \in \{0, 1, \dots, K-1\}^d$ . (a)  $K = 5$  and  $d = 1$ . (b)  $K = 4$  and  $d = 2$ .

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

$$f(\mathbf{x}) = \underbrace{\sum_{\|\alpha\|_1 \leq s-1} \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \mathbf{h}^\alpha}_{\mathcal{T}_1} + \underbrace{\sum_{\|\alpha\|_1 = s} \frac{\partial^\alpha f(\mathbf{x}_\beta + \xi_{\mathbf{x}} \mathbf{h})}{\alpha!} \mathbf{h}^\alpha}_{\mathcal{T}_2} =: \mathcal{T}_1 + \mathcal{T}_2, \quad \textcircled{2}$$

where  $\mathbf{h} = \mathbf{x} - \mathbf{x}_\beta = \mathbf{x} - \beta/K$ . It is clear that the magnitude of  $\mathcal{T}_2$  is bounded by  $\mathcal{O}(K^{-s}) = \mathcal{O}(N^{-2s/d} L^{-2s/d})$ . So we only need to construct a function in  $\mathcal{NN}$  (width  $\leq \mathcal{O}(N \ln N)$ ; depth  $\leq \mathcal{O}(L \ln L)$ ) to approximate

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

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

- Construct a ReLU network to implement a vector-valued function  $\Psi : \mathbb{R}^d \rightarrow \mathbb{R}^d$  projecting the whole cube  $Q_\beta$  to the point  $\mathbf{x}_\beta = \beta/K$ , i.e.,  $\Psi(\mathbf{x}) = \mathbf{x}_\beta$  for any  $\mathbf{x} \in Q_\beta$  and each  $\beta \in \{0, 1, \dots, K-1\}^d$ .
- Construct a ReLU network to implement a function  $P_\alpha : \mathbb{R}^d \rightarrow \mathbb{R}$  approximating the polynomial  $\mathbf{h}^\alpha$  for each  $\alpha \in \mathbb{N}^d$  with  $\|\alpha\|_1 \leq s-1$ .

---

<sup>②</sup> $\sum_{\|\alpha\|_1 = s}$  is short for  $\sum_{\|\alpha\|_1 = s, \alpha \in \mathbb{N}^d}$ . The same notation is used throughout this dissertation.

- Construct a ReLU network to implement a function  $\phi_{\alpha} : \mathbb{R}^d \rightarrow \mathbb{R}$  approximating  $\partial^{\alpha} f$  via solving a point fitting problem, *i.e.*,  $\phi_{\alpha}$  should fit  $\partial^{\alpha} f$  well at all points in  $\{\mathbf{x}_{\beta} : \beta \in \{0, 1, \dots, K-1\}^d\}$  for each  $\alpha \in \mathbb{N}^d$  with  $\|\alpha\|_1 \leq s-1$ . That is, for each  $\alpha \in \mathbb{N}^d$  with  $\|\alpha\|_1 \leq s-1$ , we need to design  $\phi_{\alpha}$  to make the following equation true.

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

Note that the first and second steps are done by Theorem 3.12 and 4.1, respectively. We will establish a proposition for the last step, which will be applied to support the construction of the desired ReLU networks. Its proof will be available later in Section 4.3.4. In fact, we can construct ReLU networks with width  $\mathcal{O}(sN \ln N)$  and depth  $\mathcal{O}(L \ln L)$  to fit  $\mathcal{O}(N^2 L^2)$  points with an error  $\mathcal{O}(N^{-2s} L^{-2s}) \leq \mathcal{O}(N^{-2s/d} L^{-2s/d})$  as shown in Proposition 4.14 below.

**Proposition 4.14.** *Given any  $N, L, s \in \mathbb{N}^+$  and  $\xi_i \in [0, 1]$  for  $i = 0, 1, \dots, N^2 L^2 - 1$ , there exists  $\phi \in \mathcal{NN}(\#input = 1; \text{width} \leq 16s(N+1)\log_2(8N); \text{depth} \leq 5(L+2)\log_2(4L); \#output = 1)$  such that*

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

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

The proof of Proposition 4.14 can be found in Section 4.3.4. Finally, let us summarize the main ideas of proving Theorem 4.13 in Table 4.1. See the detailed proof in Section 4.3.3.

### 4.3.3 Proof of auxiliary theorem

According to the key ideas of proving Theorem 4.13 we summarized in Section 4.3.2, we are ready to present the detailed proof.

*Proof of Theorem 4.13.* The detailed proof can be divided into three steps as follows.



Table 4.1: Key ideas of approximating a smooth function. Note that  $\mathbf{h} = \mathbf{x} - \Psi(\mathbf{x}) = \mathbf{x} - \mathbf{x}_\beta$  for any  $\mathbf{x} \in Q_\beta$  and each  $\beta \in \{0, 1, \dots, K-1\}^d$ .

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

**Step 1:** Set up.

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

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

Clearly,  $[0, 1]^d = \Omega([0, 1]^d, K, \delta) \cup \left( \cup_{\beta \in \{0, 1, \dots, K-1\}^d} Q_\beta \right)$  and  $\mathbf{x}_\beta$  is the vertex of  $Q_\beta$  with minimum  $\|\cdot\|_1$  norm. See Figure 4.11 for the illustrations of  $\Omega([0, 1]^d, K, \delta)$ ,  $Q_\beta$ , and  $\mathbf{x}_\beta$  for any  $\beta \in \{0, 1, \dots, K-1\}^d$ .

By Theorem 3.12, there exists  $\psi \in \mathcal{NN}(\text{width} \leq 4N+3; \text{depth} \leq 4L+5)$  such that

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

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

Define

$$\Psi(\mathbf{x}) := (\psi(x_1), \psi(x_2), \dots, \psi(x_d))/K, \quad \text{for any } \mathbf{x} \in \mathbb{R}^d,$$

then

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

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

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

**Step 2:** Construct the desired function  $\phi$ .

By Lemma 4.4, there exists  $\varphi \in \mathcal{NN}(\text{width} \leq 9(N+1)+1; \text{depth} \leq 2s(L+1))$  such that

$$|\varphi(x_1, x_2) - x_1 x_2| \leq 216(N+1)^{-2s(L+1)} =: \mathcal{E}_1, \quad \text{for any } x_1, x_2 \in [-3, 3]. \quad (4.20)$$

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

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

such that

$$|P_\alpha(\mathbf{x}) - \mathbf{x}^\alpha| \leq 9s(N+1)^{-7sL} =: \mathcal{E}_2, \quad \text{for any } \mathbf{x} \in [0, 1]^d. \quad (4.21)$$

For each  $i = 0, 1, \dots, K^d - 1$ , define

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

such that  $\sum_{j=1}^d \eta_j K^{j-1} = i$ . Such a map  $\boldsymbol{\eta}$  is a bijection from  $\{0, 1, \dots, K^d - 1\}$  to

$\{0, 1, \dots, K-1\}^d$ . For each  $\alpha \in \mathbb{N}^d$  with  $\|\alpha\|_1 \leq s-1$ , define

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

Note that  $K^d = (\lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor)^d \leq N^2 L^2$  and  $\xi_{\alpha,i} = (\partial^\alpha f(\frac{\eta(i)}{K}) + 1)/2 \in [0, 1]$  for  $i = 0, 1, \dots, K^d - 1$ . By Proposition 4.14, there exists

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

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

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

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

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

which implies by Lemma 2.1 that

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

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

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

Thus, for each  $\beta \in \{0, 1, \dots, K-1\}^d$ , we have

$$\left| \phi_\alpha(\mathbf{x}_\beta) - \partial^\alpha f(\mathbf{x}_\beta) \right| = \left| \phi_\alpha\left(\frac{\beta}{K}\right) - \partial^\alpha f\left(\frac{\beta}{K}\right) \right| \leq 2N^{-2s} L^{-2s} =: \mathcal{E}_3. \quad (4.22)$$

Now we can construct the target function  $\phi$  as

$$\phi(\mathbf{x}) := \sum_{\|\alpha\|_1 \leq s-1} \varphi\left(\frac{\phi_\alpha(\Psi(\mathbf{x}))}{\alpha!}, P_\alpha(\mathbf{x} - \Psi(\mathbf{x}))\right), \quad \text{for any } \mathbf{x} \in \mathbb{R}^d. \quad (4.23)$$

**Step 3:** Estimate approximation error.

**Fix**  $\beta \in \{0, 1, \dots, K-1\}^d$ , let us estimate the approximation error for a **fixed**  $\mathbf{x} \in Q_\beta$ . Recall that  $\Psi(\mathbf{x}) = \mathbf{x}_\beta$  and  $\mathbf{h} = \mathbf{x} - \Psi(\mathbf{x}) = \mathbf{x} - \mathbf{x}_\beta$ . Then, it is easy to verify that  $|f(\mathbf{x}) - \phi(\mathbf{x})|$  is bounded by

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

Recall that

$$\sum_{\|\alpha\|_1 = s} 1 = |\{\alpha \in \mathbb{N}^d : \|\alpha\|_1 = s\}| \leq (s+1)^{d-1} \quad \textcircled{3}$$

and

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

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

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

---

<sup>③</sup>In fact, we have  $|\{\alpha \in \mathbb{N}^d : \|\alpha\|_1 = s\}| = \binom{s+d-1}{d-1}$ , implying  $(s/d+1)^{d-1} \leq \sum_{\|\alpha\|_1 = s} 1 \leq (s+1)^{d-1}$ . Thus, the lower bound of the estimate is still exponentially large in  $d$ . To the best of our knowledge, we cannot avoid a constant prefactor that is exponentially large in  $d$  when Taylor expansion is used in the analysis.

Now let us estimate the second part  $\mathcal{J}_2$  as follows.

$$\begin{aligned}
\mathcal{J}_2 &= \sum_{\|\alpha\|_1 \leq s-1} \left| \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \mathbf{h}^\alpha - \varphi\left(\frac{\phi_\alpha(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) \right| \\
&\leq \underbrace{\sum_{\|\alpha\|_1 \leq s-1} \left| \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \mathbf{h}^\alpha - \varphi\left(\frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) \right|}_{\mathcal{J}_{2,1}} + \underbrace{\sum_{\|\alpha\|_1 \leq s-1} \left| \varphi\left(\frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) - \varphi\left(\frac{\phi_\alpha(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) \right|}_{\mathcal{J}_{2,2}} \\
&=: \mathcal{J}_{2,1} + \mathcal{J}_{2,2}.
\end{aligned}$$

Note that  $\mathcal{E}_2 = 9s(N+1)^{-7sL} \leq 9s(2)^{-7s} \leq 2$ . By Equation (4.21), it easy to verify that  $P_\alpha(\mathbf{h}) \in [-2, 3] \subseteq [-3, 3]$  for each  $\alpha \in \mathbb{N}^d$  with  $\|\alpha\|_1 \leq s-1$ . Clearly,  $\mathbf{h} \in [0, 1]^d$  and  $\frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \in [-1, 1] \subseteq [-3, 3]$  for each  $\alpha$ . Then, by Equation (4.20) and (4.21), we have

$$\begin{aligned}
\mathcal{J}_{2,1} &= \sum_{\|\alpha\|_1 \leq s-1} \left| \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \mathbf{h}^\alpha - \varphi\left(\frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) \right| \\
&\leq \sum_{\|\alpha\|_1 \leq s-1} \left( \left| \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \mathbf{h}^\alpha - \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} P_\alpha(\mathbf{h}) \right| + \underbrace{\left| \frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} P_\alpha(\mathbf{h}) - \varphi\left(\frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) \right|}_{\leq \mathcal{E}_1 \text{ by Eq. (4.20)}} \right) \\
&\leq \sum_{\|\alpha\|_1 \leq s-1} \left( \underbrace{\left| \frac{1}{\alpha!} \mathbf{h}^\alpha - P_\alpha(\mathbf{h}) \right|}_{\leq \mathcal{E}_2 \text{ by Eq. (4.21)}} + \mathcal{E}_1 \right) \leq \sum_{\|\alpha\|_1 \leq s-1} (\mathcal{E}_2 + \mathcal{E}_1) \leq s^d (\mathcal{E}_1 + \mathcal{E}_2).
\end{aligned}$$

To estimate  $\mathcal{J}_{2,2}$ , we need the following fact derived from Equation (4.20):

$$\begin{aligned}
|\varphi(x_1, x_2) - \varphi(\tilde{x}_1, x_2)| &\leq \underbrace{|\varphi(x_1, x_2) - x_1 x_2|}_{\leq \mathcal{E}_1 \text{ by Eq. (4.20)}} + \underbrace{|\varphi(\tilde{x}_1, x_2) - \tilde{x}_1 x_2|}_{\leq \mathcal{E}_1 \text{ by Eq. (4.20)}} + |x_1 x_2 - \tilde{x}_1 x_2| \\
&\leq 2\mathcal{E}_1 + 3|x_1 - \tilde{x}_1|,
\end{aligned} \tag{4.24}$$

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

Since  $\mathcal{E}_3 = 2N^{-2s}L^{-2s} \leq 2$  and  $\frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!} \in [-1, 1]$  for each  $\alpha \in \mathbb{N}^d$  with  $\|\alpha\|_1 \leq s-1$ , we have  $\frac{\phi_\alpha(\mathbf{x}_\beta)}{\alpha!} \in [-3, 3]$  by Equation (4.22). Recall that  $P_\alpha(\mathbf{h}) \in [-3, 3]$ .

Then, by Equation (4.22) and (4.24), we have

$$\begin{aligned}
\mathcal{J}_{2,2} &= \sum_{\|\alpha\|_1 \leq s-1} \left| \underbrace{\varphi\left(\frac{\partial^\alpha f(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right) - \varphi\left(\frac{\phi_\alpha(\mathbf{x}_\beta)}{\alpha!}, P_\alpha(\mathbf{h})\right)}_{\leq 2\mathcal{E}_1 + \frac{3}{\alpha!}|\partial^\alpha f(\mathbf{x}_\beta) - \phi_\alpha(\mathbf{x}_\beta)| \text{ by Eq. (4.24)}} \right| \\
&\leq \sum_{\|\alpha\|_1 \leq s-1} \left( 2\mathcal{E}_1 + \frac{3}{\alpha!} \underbrace{|\partial^\alpha f(\mathbf{x}_\beta) - \phi_\alpha(\mathbf{x}_\beta)|}_{\leq \mathcal{E}_3 \text{ by Eq. (4.22)}} \right) \leq \sum_{\|\alpha\|_1 \leq s-1} (2\mathcal{E}_1 + 3\mathcal{E}_3) \leq s^d(2\mathcal{E}_1 + 3\mathcal{E}_3).
\end{aligned}$$

Therefore, we have

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

Recall the fact  $[0, 1]^d \setminus \Omega([0, 1]^d, K, \delta) = \cup_{\beta \in \{0, 1, \dots, K-1\}^d} Q_\beta$ . Since  $\beta \in \{0, 1, \dots, K-1\}^d$  and  $\mathbf{x} \in Q_\beta$  are arbitrary, we have,

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

for any  $\mathbf{x} \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)$ . Note that  $K = \lfloor N^{1/d} \rfloor^2 \lfloor L^{2/d} \rfloor \geq \frac{N^{2/d} L^{2/d}}{8}$  and

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

Then we have

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

It remains to estimate the width and depth of the network implementing  $\phi$ .

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

$$\begin{cases} \Psi \in \mathcal{NN}(\text{width} \leq d(4N+3); \text{depth} \leq 4L+5), \\ \phi_\alpha \in \mathcal{NN}(\text{width} \leq 16s(N+1)\log_2(8N); \text{depth} \leq 5(L+2)\log_2(4L)), \\ P_\alpha \in \mathcal{NN}(\text{width} \leq 9(N+1)+s-1; \text{depth} \leq 7s^2L), \\ \varphi \in \mathcal{NN}(\text{width} \leq 9N+10; \text{depth} \leq 2s(L+1)). \end{cases}$$

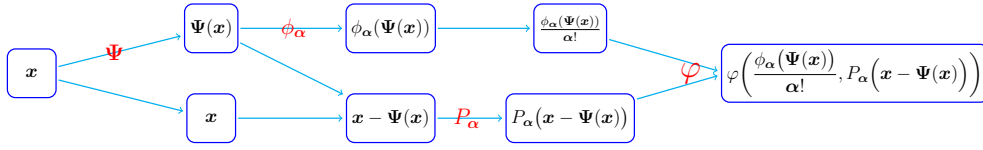


Figure 4.12: An illustration of the sub-network architecture implementing  $\varphi\left(\frac{\phi_\alpha(\Psi(\mathbf{x}))}{\alpha!}, P_\alpha(\mathbf{x} - \Psi(\mathbf{x}))\right)$  for each  $\alpha \in \mathbb{N}^d$  with  $\|\alpha\| \leq s-1$  when  $\mathbf{x} \in Q_\beta$  for each  $\beta \in \{0, 1, \dots, K-1\}^d$ .

By Equation (4.23) and Figure 4.12, it easy to verify  $\phi$  can be implemented by a ReLU network with width

$$\begin{aligned} \sum_{\|\alpha\|_1 \leq s-1} 16sd(N+2)\log_2(8N) &\leq s^d \cdot 16sd(N+2)\log_2(8N) \\ &= 16s^{d+1}d(N+2)\log_2(8N) \end{aligned}$$

and depth

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

So we finish the proof.  $\square$

#### 4.3.4 Proof of key proposition for auxiliary theorem

Let us discuss the construction of ReLU networks to fit a collection of points in  $\mathbb{R}^2$ . It is trivial to fit  $n$  points via one-hidden-layer ReLU networks with  $\mathcal{O}(n)$  parameters. However, to prove Proposition 4.14, we need to fit  $n$  points with much

fewer parameters, which is the main difficulty of our proof. Our proof below is mainly based on the “bit extraction” technique and the idea of function compositions.

*Proof of Proposition 4.14.* Set  $J = \lceil 2s \log_2(NL + 1) \rceil \in \mathbb{N}^+$ . For each  $\xi_i \in [0, 1]$ , there exist  $\xi_{i,1}, \xi_{i,2}, \dots, \xi_{i,J} \in \{0, 1\}$  such that

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

By Theorem 3.4, there exist

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

such that

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

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

$$\begin{aligned} \left| \sum_{j=1}^J 2^{-j} \phi_j(i) - \xi_i \right| &= \left| \sum_{j=1}^J 2^{-j} \xi_{i,j} - \xi_i \right| \\ &= |\text{bin}0.\xi_{i,1}\xi_{i,2}\dots\xi_{i,J} - \xi_i| \leq 2^{-J} \leq N^{-2s} L^{-2s}, \end{aligned} \tag{4.25}$$

where the last inequality comes from

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

Recall that

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

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



network architecture in Figure 4.13 to implement a function  $\tilde{\phi}$  such that

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

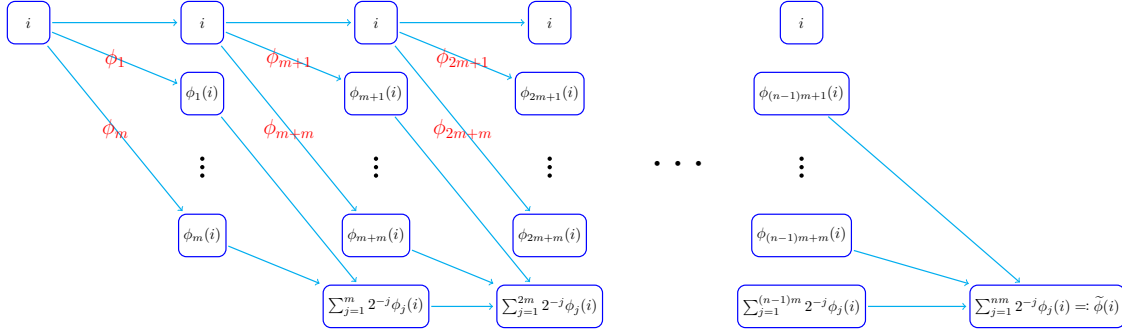


Figure 4.13: An illustration of the network architecture implementing  $\tilde{\phi}(i) = \sum_{j=1}^J 2^{-j} \phi_j(i)$  for any  $i \in \{0, 1, \dots, N^2 L^2 - 1\}$ . We assume  $J = mn$ , where  $m = 2s \lceil \log_2(4N) \rceil$  and  $n = \lceil \log_2(2L) \rceil$ , since we can set  $\phi_{J+1} = \dots = \phi_{nm} = 0$  if  $J < nm$ .

Clearly, the network architecture in Figure 4.13 is with width

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

and depth

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

Finally, we define

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

See Figure 4.14 for the network architecture implementing  $\phi$ . Then  $0 \leq \phi(x) \leq 1$  for any  $x \in \mathbb{R}$  and  $\phi$  can be implemented by a ReLU network with width  $16s(N + 1) \log_2(8N)$  and depth  $(5L + 8) \log_2(4L) + 3 \leq 5(L + 2) \log_2(4L)$ .

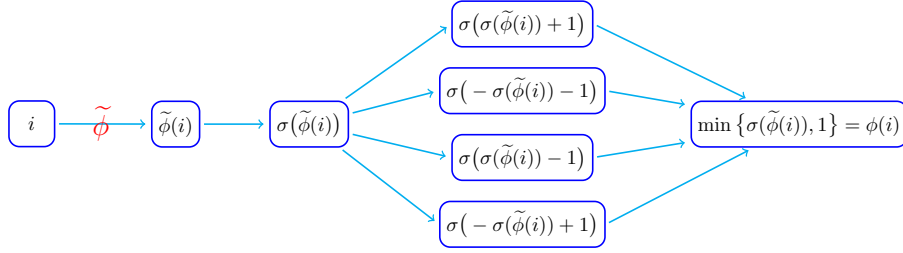


Figure 4.14: An illustration of the network architecture implementing the desired function  $\phi$  for  $i = 0, 1, \dots, N^2 L^2 - 1$ , based on the fact  $\min\{x_1, x_2\} = \frac{x_1 + x_2 - |x_1 - x_2|}{2} = \frac{\sigma(x_1 + x_2) - \sigma(-x_1 - x_2) - \sigma(x_1 - x_2) - \sigma(-x_1 + x_2)}{2}$ .

Note that

$$\tilde{\phi}(i) = \sum_{j=1}^J 2^{-j} \phi_j(i) = \sum_{j=1}^J 2^{-j} \xi_{i,j} = \text{bin} 0.\xi_{i,1}\xi_{i,2}\cdots\xi_{i,J} \in [0, 1],$$

for  $i = 0, 1, \dots, N^2 L^2 - 1$ , implying

$$|\phi(i) - \xi_i| = \left| \min\{\sigma(\tilde{\phi}(i)), 1\} - \xi_i \right| = |\tilde{\phi}(i) - \xi_i| = \left| \sum_{j=1}^J 2^{-j} \xi_{i,j} - \xi_i \right| \leq N^{-2s} L^{-2s},$$

where the last inequality comes from Equation (4.25). So the proof is complete.  $\square$

## 4.4 Optimality of approximation by networks

In this section, we will study the best possible approximation errors for several function spaces approximated by ReLU networks. To this end, we adopt the method in [38, 52, 53, 58, 59, 60] via studying the connection between the approximation error and VC-dimension. Thus, let us first present the definitions of VC-dimension and related concepts.

**Definition 4.15** (Growth function, VC-dimension, Shattering). Let  $H$  be a class of functions mapping from a general domain  $\mathcal{X}$  to  $\{0, 1\}$ . For any  $m \in \mathbb{N}^+$ , we define

the growth function of  $H$  as

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

where  $|S|$  denotes the size of a set  $S$ .

We say  $H$  shatters the set  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\} \subseteq \mathcal{X}$  if

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

The Vapnik-Chervonenkis (VC) dimension of  $H$ , denoted by  $\text{VCDim}(H)$ , is the size of the largest shattered set, namely, the largest  $m$  such that  $\Pi_H(m) = 2^m$ . By convention,  $\text{VCDim}(H) = \infty$  if  $\Pi_H(m) = 2^m$  for all  $m \in \mathbb{N}^+$ .

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).

**Definition 4.16.** Let  $Q(\mathbf{x}_0, \eta) \subseteq \mathbb{R}^d$  denote the closed cube with center  $\mathbf{x}_0$  and sidelength  $\eta$ . For any cube  $Q = Q(\mathbf{x}_0, \eta)$ ,  $rQ$  denote the closed cube satisfying two conditions: 1)  $rQ$  has the same center as  $Q$ ; 2) the sidelength of  $rQ$  is equal to the multiplication of  $r$  and that of  $Q$ .

#### 4.4.1 Hölder continuous functions

Let us first consider the Hölder continuous function space  $\text{Hölder}([0, 1]^d, \alpha, \lambda)$ . Without loss of generality, we assume  $\lambda = 1$ . Theorem 4.17 below shows that the best possible approximation error of functions in  $\text{Hölder}([0, 1]^d, \alpha, 1)$  approximated by functions in  $\mathcal{F}$  is bounded by a formula characterized by  $\text{VCDim}(\mathcal{F})$ .

**Theorem 4.17.** *Given any  $\varepsilon \in (0, 2/9)$  and a function set  $\mathcal{F}$  with all elements defined on  $[0, 1]^d$ , if*

$$\inf_{\phi \in \mathcal{F}} \|\phi - f\|_{L^\infty([0,1]^d)} \leq \varepsilon, \quad \text{for any } f \in \text{Hölder}([0, 1]^d, \alpha, 1), \quad (4.26)$$

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

This theorem investigates the connection between VC-dimension of  $\mathcal{F}$  and the approximation errors of functions in  $\text{Hölder}([0, 1]^d, \alpha, 1)$  approximated by elements of  $\mathcal{F}$ . In other words, 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  $\text{Hölder}([0, 1]^d, \alpha, 1)$ . It is shown in [4] that VC-dimension of ReLU networks with a fixed architecture with  $W$  parameters and  $L$  layers has an upper bound  $\mathcal{O}(WL \ln W)$ . It follows that VC-dimension of ReLU networks with width  $N$  and depth  $L$  is bounded by  $\mathcal{O}(N^2 L \cdot L \cdot \ln(N^2 L)) \leq \mathcal{O}(N^2 L^2 \ln(NL))$ . That is,  $\text{VCDim}(\mathcal{F}) \leq \mathcal{O}(N^2 L^2 \ln(NL))$ , where

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

We denote the best approximation error of functions in  $\text{Hölder}([0, 1]^d, \alpha, 1)$  approximated by ReLU networks with width  $N$  and depth  $L$  as

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

for any  $N, L \in \mathbb{N}^+$ . Then, by Theorem 4.17 and Corollary 4.7, we have

$$\underbrace{C_1(\alpha, d) \cdot \left( N^2 L^2 \ln(NL) \right)^{-\alpha/d}}_{\text{implied by Theorem 4.17}} \leq \mathcal{E}_{\alpha,d}(N, L) \leq \underbrace{C_2(\alpha, d) \cdot \left( N^2 L^2 \right)^{-\alpha/d}}_{\text{shown in Corollary 4.7}},$$

for any  $N, L \in \mathbb{N}^+$ ,<sup>④</sup> where  $C_1(\alpha, d)$  and  $C_2(\alpha, d)$  are two positive constants determined by  $\alpha$  and  $d$ , and  $C_2(\alpha, d)$  can be **explicitly** represented. Therefore, the approximation error in Corollary 4.7 is nearly optimal.

Now let us present the detailed proof of Theorem 4.17.

*Proof of Theorem 4.17.* 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  $f_\chi$  and Equation (4.26) such that  $\{\phi_\chi : \chi \in \mathcal{B}\} \subseteq \mathcal{F}$  also shatters  $\mathcal{O}(\varepsilon^{-d/\alpha})$  points.

The details of these two steps can be found below.

**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.

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

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

for any index vector  $\beta = (\beta_1, \beta_2, \dots, \beta_d) \in \{0, 1, \dots, K-1\}^d$ .

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

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

---

<sup>④</sup>To make this equation hold for any  $N, L \in \mathbb{N}^+$ , one needs to choose  $C_1(\alpha, d)$  and  $C_2(\alpha, d)$  carefully based on Theorem 4.17 and Corollary 4.7.

Define

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

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

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

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

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

By Equation (4.26), for each  $\chi \in \mathcal{B}$ , there exists  $\phi_\chi \in \mathcal{F}$  such that

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

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

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

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

$$|\phi_\chi(\mathbf{x}) - f_\chi(\mathbf{x})| \leq \frac{82}{81}\varepsilon, \quad \text{for any } \chi \in \mathcal{B} \text{ and } \mathbf{x} \in [0, 1]^d \setminus \mathcal{H}. \quad (4.27)$$

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

$$|f_\chi(\mathbf{x})| = \zeta_{Q_{\boldsymbol{\beta}}}(\mathbf{x}) \geq \frac{9}{10} \zeta_{Q_{\boldsymbol{\beta}}}(\mathbf{x}_{Q_{\boldsymbol{\beta}}}) = \frac{9}{10} \left( \frac{1}{2 \lfloor (9\varepsilon/2)^{-1/\alpha} \rfloor} \right)^\alpha / 2 \geq \frac{81}{80}\varepsilon, \quad (4.28)$$

where  $\mathbf{x}_{Q_{\boldsymbol{\beta}}}$  is the center of  $Q_{\boldsymbol{\beta}}$ .

It follows from  $\mu((\frac{1}{10}Q_\beta) \setminus \mathcal{H}) > 0$  that  $(\frac{1}{10}Q_\beta) \setminus \mathcal{H}$  is not empty for each  $\beta \in \{0, 1, \dots, K-1\}^d$ . Thus, by Equation (4.27) and (4.28), for each  $\beta \in \{0, 1, \dots, K-1\}^d$  and each  $\chi \in \mathcal{B}$ , there exists  $\mathbf{x}_\beta \in (\frac{1}{10}Q_\beta) \setminus \mathcal{H}$  such that

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

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

$$\text{VCDim}(\mathcal{F}) \geq \text{VCDim}(\{\phi_\chi : \chi \in \mathcal{B}\}) \geq K^d = \lfloor (9\varepsilon/2)^{-1/\alpha} \rfloor^d \geq (9\varepsilon)^{-d/\alpha}, \quad (4.29)$$

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.  $\square$

#### 4.4.2 Smooth functions

Next, let us consider another function space  $C_u^s([0, 1]^d)$ , which is the closed unit ball of the smooth function space  $C^s([0, 1]^d)$ . Theorem 4.18 below shows that the best possible approximation error of functions in  $C_u^s([0, 1]^d)$  approximated by functions in  $\mathcal{F}$  is bounded by a formula characterized by  $\text{VCDim}(\mathcal{F})$ .

**Theorem 4.18.** *Given any  $s, d \in \mathbb{N}^+$ , there exists a small positive constant  $C_{s,d}$  determined by  $s$  and  $d$  such that: For any  $\varepsilon \in (0, (2^d C_{s,d})^{s/d}]$  and a function set  $\mathcal{F}$  with all elements defined on  $[0, 1]^d$ , if*

$$\inf_{\phi \in \mathcal{F}} \|\phi - f\|_{L^\infty([0, 1]^d)} \leq \varepsilon, \quad \text{for any } f \in C_u^s([0, 1]^d), \quad (4.30)$$

then  $\text{VCDim}(\mathcal{F}) \geq C_{s,d} \varepsilon^{-d/s}$ .  $\textcircled{5}$

This theorem demonstrates the connection between VC-dimension of  $\mathcal{F}$  and the

---

$\textcircled{5}$ In fact,  $C_{s,d}$  can be expressed by  $s$  and  $d$  with a **explicitly** formula as we remark in the proof of this theorem. However, the formula may be very complicated.

approximation error using elements of  $\mathcal{F}$  to approximate functions in  $C_u^s([0, 1]^d)$ . To be precise, the best possible approximation error is controlled by  $\mathcal{O}(\text{VCDim}(\mathcal{F})^{-s/d})$ . A typical application of this theorem is to prove the optimality of approximation errors when using ReLU networks to approximate functions in  $C_u^s([0, 1]^d)$ . It is shown in [4] that VC-dimension of ReLU networks with a fixed architecture with  $W$  parameters and  $L$  layers has an upper bound  $\mathcal{O}(WL \ln W)$ . It follows that VC-dimension of ReLU networks with width  $N$  and depth  $L$  is bounded by  $\mathcal{O}(N^2 L \cdot L \cdot \ln(N^2 L)) \leq \mathcal{O}(N^2 L^2 \ln(NL))$ . That is,  $\text{VCDim}(\mathcal{F}) \leq \mathcal{O}(N^2 L^2 \ln(NL))$ , where

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

We denote the best approximation error of functions in  $C_u^s([0, 1]^d)$  approximated by ReLU networks with width  $N$  and depth  $L$  as

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

for any  $N, L \in \mathbb{N}^+$ , where  $C_u^s([0, 1]^d)$  is the closed unit ball of  $C^s([0, 1]^d)$ . Then, by Theorem 4.18 and Corollary 4.12, we have

$$\underbrace{C_1(s, d) \cdot \left( N^2 L^2 \ln(NL) \right)^{-s/d}}_{\text{implied by Theorem 4.18}} \leq \mathcal{E}_{s,d}(N, L) \leq \underbrace{C_2(s, d) \cdot \left( \frac{N^2 L^2}{(\ln N \ln L)^2} \right)^{-s/d}}_{\text{shown in Corollary 4.12}},$$

for any  $N, L \in \mathbb{N}^+$  with  $N \geq 2$  and  $L \geq 2$ ,<sup>⑥</sup> where  $C_1(s, d)$  and  $C_2(s, d)$  are two positive constants in  $s$  and  $d$ , and  $C_2(s, d)$  can be **explicitly** expressed. Therefore, the approximation errors in Theorem 4.11 and Corollary 4.12 are nearly optimal.

Now let us present the detailed proof of Theorem 4.18.

*Proof of Theorem 4.18.* To find a subset of  $\mathcal{F}$  shattering  $\mathcal{O}(\varepsilon^{-d/s})$  points in  $[0, 1]^d$ ,

---

<sup>⑥</sup>To make this equation hold for any  $N, L \in \mathbb{N}^+$  with  $N \geq 2$  and  $L \geq 2$ , one needs to choose  $C_1(s, d)$  and  $C_2(s, d)$  carefully based on Theorem 4.18 and Corollary 4.12.



we divide the proof into two steps.

- Construct  $\{f_\chi : \chi \in \mathcal{B}\} \subseteq C_u^s([0, 1]^d)$  that scatters  $\mathcal{O}(\varepsilon^{-d/s})$  points, where  $\mathcal{B}$  is a set defined later.
- Design  $\phi_\chi \in \mathcal{F}$ , for each  $\chi \in \mathcal{B}$ , based on  $f_\chi$  and Equation (4.30) such that  $\{\phi_\chi : \chi \in \mathcal{B}\} \subseteq \mathcal{F}$  also shatters  $\mathcal{O}(\varepsilon^{-d/s})$  points.

The details of these two steps can be found below.

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

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

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

for any index vector  $\beta = (\beta_1, \beta_2, \dots, \beta_d) \in \{0, 1, \dots, K-1\}^d$ .

There exists a “bump function”  $\tilde{g} \in C^\infty(\mathbb{R}^d)$  such that  $\tilde{g}(\mathbf{0}) = 1$  and  $\tilde{g}(\mathbf{x}) = 0$  for  $\|\mathbf{x}\|_2 \geq 1/3$ . For example, we can define  $\tilde{g}$  as

$$\tilde{g}(\mathbf{x}) := \begin{cases} \exp\left(\frac{1}{\|\mathbf{x}\|_2^2 - 1} + 1\right), & \text{if } \|\mathbf{x}\|_2 < 1/3, \\ 0, & \text{otherwise,} \end{cases}$$

where  $\exp(x) = e^x$  for any  $x \in \mathbb{R}$  and  $e \approx 2.7$  is the natural logarithmic base. Then, we have  $g := \tilde{g}/\tilde{C}_{s,d} \in C_u^s([0, 1]^d)$  by setting  $\tilde{C}_{s,d} := \|\tilde{g}\|_{C^s([0, 1]^d)}$ .

Define

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

and

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

where  $\mathbf{x}_{Q_\beta}$  is the center of  $Q_\beta$ . Then, we have

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

Next, for each  $\chi \in \mathcal{B}$ , we can define  $f_\chi$  via

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

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

- By the definition of  $g_\beta$  and  $\chi$ , we have

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

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

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

$$\text{implying } |\partial^\alpha f_\chi(\mathbf{x})| = |K^{-(s-\|\boldsymbol{\alpha}\|_1)} \partial^\alpha g(K(\mathbf{x} - \mathbf{x}_\beta))| \leq 1.$$

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

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

By Equation (4.30), for each  $\chi \in \mathcal{B}$ , there exists  $\phi_\chi \in \mathcal{F}$  such that

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

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

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

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

$$|\phi_\chi(\mathbf{x}) - f_\chi(\mathbf{x})| \leq \frac{3}{2}\varepsilon, \quad \text{for any } \chi \in \mathcal{B} \text{ and } \mathbf{x} \in [0, 1]^d \setminus \mathcal{H}. \quad (4.31)$$

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

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

where  $\mathbf{x}_{Q_\beta}$  is the center of  $Q_\beta$ .

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

$$|f_\chi(\mathbf{x}_\beta)| = g_\beta(\mathbf{x}_\beta) \geq \frac{1}{2}g_\beta(\mathbf{x}_{Q_\beta}) = \frac{1}{2}K^{-s}g(\mathbf{0}) = \frac{1}{2}K^{-s}/\tilde{C}_{s,d} \geq 2\varepsilon, \quad (4.32)$$

where the last inequality is attained by setting  $K = \lfloor (4\varepsilon\tilde{C}_{s,d})^{-1/s} \rfloor$ . Since our proof is invalid when  $K = 0$ , it is necessary to guarantee  $K = \lfloor (4\varepsilon\tilde{C}_{s,d})^{-1/s} \rfloor \geq 1$ , which will be verified later.

By Equation (4.31) and (4.32), we have, for each  $\beta \in \{0, 1, \dots, K-1\}^d$  and each  $\chi \in \mathcal{B}$ ,

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

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

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

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

Finally, set

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

This means  $C_{s,d}$  can be computed by an explicit mathematical formula based on the function  $\tilde{g}$  defined previously. Moreover, we have

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

and

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

where the last inequality comes from  $\varepsilon \in (0, (2^d C_{s,d})^{s/d}]$ . So we finish the proof.  $\square$

## Approximation by Floor-ReLU networks

As shown in Section 4.1, an exponential approximation error  $O(N^{-L})$  can be achieved when using ReLU networks with width  $\mathcal{O}(N)$  and depth  $\mathcal{O}(L)$  to approximate polynomials on  $[0, 1]^d$ . But such an exponential error is not true for general function spaces as discussed in Section 4.4. The limitation of ReLU networks motivates us to explore other types of network architectures to admit (nearly) exponential approximation errors.

In particular, we introduce new networks built with either Floor ( $\lfloor x \rfloor$ ) or ReLU ( $\max\{0, x\}$ ) as the activation function<sup>①</sup> in each neuron. We call such networks Floor-ReLU networks. See Figure 5.1 for an example. We will prove in this chapter that Floor-ReLU networks with fixed architectures can attain nearly exponential approximation errors for approximating (Hölder) continuous functions on  $[0, 1]^d$ .

### 5.1 Main theorem and its proof

In Theorem 5.1 below, we show by construction that Floor-ReLU networks, with fixed architectures, with width  $\max\{d, 5N + 13\}$  and depth  $64dL + 3$  can uniformly approximate an arbitrary continuous function  $f$  on  $[0, 1]^d$  with a nearly exponential

---

<sup>①</sup>Our results can be easily generalized to Ceiling-ReLU networks, namely, feed-forward fully connected neural networks with either Ceiling ( $\lceil x \rceil$ ) or ReLU ( $\max\{0, x\}$ ) as the activation function in each neuron.

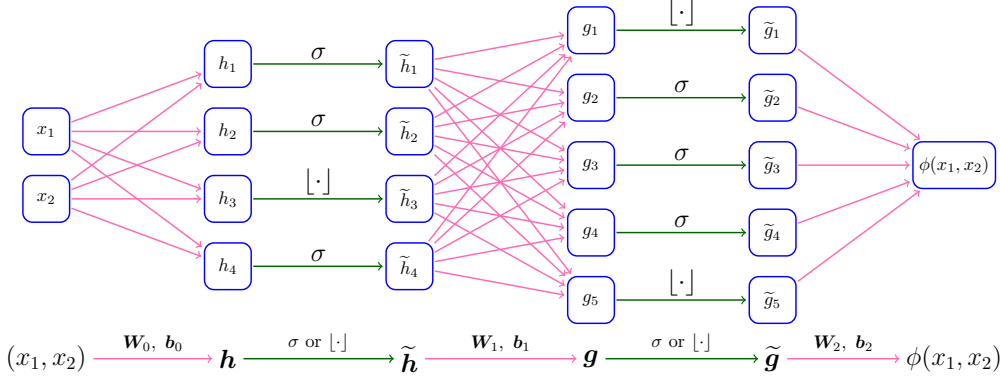


Figure 5.1: An example of a Floor-ReLU network with width 5 and depth 2.

approximation error  $\omega_f(\sqrt{d} N^{-\sqrt{L}}) + 2\omega_f(\sqrt{d})N^{-\sqrt{L}}$ .

**Theorem 5.1.** *Given any  $N, L, d \in \mathbb{N}^+$ , there exists a fixed Floor-ReLU network architecture with width  $\max\{d, 5N + 13\}$  and depth  $64dL + 3$  such that: For any continuous function  $f \in C([0, 1]^d)$ , there exists a function  $\phi$ , implemented by this Floor-ReLU network architecture with proper parameters, satisfying*

$$|\phi(\mathbf{x}) - f(\mathbf{x})| \leq \omega_f(\sqrt{d} N^{-\sqrt{L}}) + 2\omega_f(\sqrt{d})N^{-\sqrt{L}}, \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

With Theorem 5.1, we have an immediate corollary.

**Corollary 5.2.** *Given any  $\bar{N}, \bar{L}, d \in \mathbb{N}^+$ , there exists a fixed Floor-ReLU network architecture with width  $\bar{N}$  and depth  $\bar{L}$  such that: For any continuous function  $f \in C([0, 1]^d)$ , there exists a function  $\phi$ , implemented by this Floor-ReLU network architecture with proper parameters, satisfying*

$$|\phi(\mathbf{x}) - f(\mathbf{x})| \leq \omega_f\left(\sqrt{d} \left\lfloor \frac{\bar{N}-13}{5} \right\rfloor^{-\sqrt{\left\lfloor \frac{\bar{L}-3}{64d} \right\rfloor}}\right) + 2\omega_f(\sqrt{d}) \left\lfloor \frac{\bar{N}-13}{5} \right\rfloor^{-\sqrt{\left\lfloor \frac{\bar{L}-3}{64d} \right\rfloor}},$$

for any  $\mathbf{x} \in [0, 1]^d$  and  $\bar{N}, \bar{L} \in \mathbb{N}^+$  with  $\bar{N} \geq \max\{d, 18\}$  and  $\bar{L} \geq 64d + 3$ .

We would like to remark that the Floor-ReLU network architectures in Theorem 5.1 and Corollary 5.2 are independent of the target function  $f$ . That is, only the values of the parameters rely on the target function  $f$ . In particular, the choice

of activation functions (Floor or ReLU) in each neuron is independent of the target function  $f$ .

In Theorem 5.1, the error in  $\omega_f(\sqrt{d}N^{-\sqrt{L}})$  implicitly depends on  $N$  and  $L$  through the modulus of continuity of  $f$ , while the error in  $2\omega_f(\sqrt{d})N^{-\sqrt{L}}$  is explicit in  $N$  and  $L$ . Simplifying the implicit approximation error to make it explicitly depending on  $N$  and  $L$  is challenging in general. However, if  $f$  is a Hölder continuous function on  $[0, 1]^d$  of order  $\alpha \in (0, 1]$  with a constant  $\lambda$ , i.e.,  $f \in \text{Hölder}([0, 1]^d, \alpha, \lambda)$ , then we have

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq \lambda \|\mathbf{x} - \mathbf{y}\|_2^\alpha, \quad \text{for any } \mathbf{x}, \mathbf{y} \in [0, 1]^d, \quad (5.1)$$

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

**Corollary 5.3.** *Given any  $N, L, d \in \mathbb{N}^+$ , there exist a fixed Floor-ReLU network architecture with width  $\max\{d, 5N + 13\}$  and depth  $64dL + 3$  such that: For any function  $f \in \text{Hölder}([0, 1]^d, \alpha, \lambda)$ , there exists a function  $\phi$ , implemented by this Floor-ReLU network architecture with proper parameters, satisfying*

$$|\phi(\mathbf{x}) - f(\mathbf{x})| \leq 3\lambda d^{\alpha/2}N^{-\alpha\sqrt{L}}, \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

First, Theorem 5.1 and Corollary 5.3 show that the approximation capacity of deep networks for continuous functions can be nearly exponentially improved by increasing the network depth, and the approximation error can be explicitly characterized in terms of the width  $\mathcal{O}(N)$  and depth  $\mathcal{O}(L)$ . 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  $\omega_f(\sqrt{d}N^{-\sqrt{L}})$ . Finally, applying piecewise constant and integer-valued functions as activation functions and integer numbers as parameters has been explored in the

study of quantized neural networks [6, 26, 61] with efficient training algorithms for low computational complexity [56]. The floor function ( $\lfloor x \rfloor$ ) is a piecewise constant function and can be easily implemented numerically at very little cost. Hence, the evaluation of the proposed network could be efficiently implemented in practical computation. Though there might not be an existing optimization algorithm to identify an approximant with an approximation error, Theorem 5.1 and Corollary 5.3 can provide an expected accuracy before a learning task and how much the current optimization algorithms could be improved. Designing an efficient optimization algorithm for Floor-ReLU networks will be left as future work with several possible directions discussed later.

In particular, we let  $N = 2$  and  $L = W$  in Theorem 5.1, then the width is  $\max\{d, 23\}$ , the depth is  $64dW + 3$ , and the total number of parameters is bounded by  $\mathcal{O}(\max\{d^2, 23^2\}(64dW + 3)) = \mathcal{O}(W)$ . Therefore, we can prove Corollary 5.4 below stating that Floor-ReLU networks can provide a nearly exponential approximation error in terms of the number of parameters.

**Corollary 5.4.** *Given any  $W, d \in \mathbb{N}^+$ , there exists a fixed Floor-ReLU network architecture with  $\mathcal{O}(W)$  parameters, width  $\max\{d, 23\}$ , and depth  $64dW + 3$ , such that: For any continuous function  $f \in C([0, 1]^d)$ , there exists a function  $\phi$ , implemented by this Floor-ReLU network architecture with proper parameters, satisfying*

$$|\phi(\mathbf{x}) - f(\mathbf{x})| \leq \omega_f(\sqrt{d}2^{-\sqrt{W}}) + 2\omega_f(\sqrt{d})2^{-\sqrt{W}}, \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

We would like to point out that the derivative of Floor is zero almost everywhere, which may lead to the failure of the backpropagation algorithm. To overcome this, we propose three possible methods as follows.

- First, we can consider gradient-free optimization methods, *e.g.*, particle swarm optimization [30], consensus-based optimization [11, 50].
- The second method is to apply optimization algorithms for quantized networks that also have piecewise constant activation functions [6, 9, 26, 36, 56, 61]. For



example, an empirical method is to use a straight through estimator (STE) by setting the incoming gradients to the activation function equal to its outgoing gradients, disregarding the derivative of the activation function itself.

- The final method is to use the linear combination of ReLU and Floor, *i.e.*,  $p\sigma(x) + (1-p)\lfloor x \rfloor$  for  $p \in (0, 1)$ , to replace Floor ( $\lfloor x \rfloor$ ) to avoid zero derivative. Similar to Theorem 5.1 and Corollary 5.3, the nearly exponential errors can be attained with the new activation functions ( $\sigma$  and  $p\sigma(x) + (1-p)\lfloor x \rfloor$ ).

The proof of Theorem 5.1 is an immediate result of Theorem 5.5 below.

**Theorem 5.5.** *Given a continuous function  $f \in C([0, 1]^d)$ , for any  $N, L \in \mathbb{N}^+$ , there exists a function  $\phi$  implemented by a Floor-ReLU network, with a fixed architecture independent of  $f$ , with width  $\max\{d, 2N^2 + 5N\}$  and depth  $7dL^2 + 3$  such that*

$$|\phi(\mathbf{x}) - f(\mathbf{x})| \leq \omega_f(\sqrt{d} N^{-L}) + 2\omega_f(\sqrt{d})2^{-NL}, \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

Theorem 5.5 will be proved later in this section. Now let us prove Theorem 5.1 based on Theorem 5.5.

*Proof of Theorem 5.1.* Given any  $N, L \in \mathbb{N}^+$ , there exist  $\tilde{N}, \tilde{L} \in \mathbb{N}^+$  with  $\tilde{N} \geq 2$  and  $\tilde{L} \geq 3$  such that

$$(\tilde{N} - 1)^2 \leq N < \tilde{N}^2 \quad \text{and} \quad (\tilde{L} - 1)^2 \leq 4L < \tilde{L}^2.$$

By Theorem 5.5, there exists a function  $\phi$  implemented by a Floor-ReLU network, with a fixed architecture independent of  $f$ , with width  $\max\{d, 2\tilde{N}^2 + 5\tilde{N}\}$  and depth  $7d\tilde{L}^2 + 3$  such that

$$|\phi(\mathbf{x}) - f(\mathbf{x})| \leq \omega_f(\sqrt{d} \tilde{N}^{-\tilde{L}}) + 2\omega_f(\sqrt{d})2^{-\tilde{N}\tilde{L}}, \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

Note that

$$2^{-\tilde{N}\tilde{L}} \leq \tilde{N}^{-\tilde{L}} = (\tilde{N}^2)^{-\frac{1}{2}\sqrt{\tilde{L}^2}} \leq N^{-\frac{1}{2}\sqrt{4L}} \leq N^{-\sqrt{L}}.$$

Then we have

$$|\phi(\mathbf{x}) - f(\mathbf{x})| \leq \omega_f(\sqrt{d}N^{-\sqrt{L}}) + 2\omega_f(\sqrt{d})N^{-\sqrt{L}}, \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

For  $\tilde{N}, \tilde{L} \in \mathbb{N}^+$  with  $\tilde{N} \geq 2$  and  $\tilde{L} \geq 3$ , we have

$$2\tilde{N}^2 + 5\tilde{N} \leq 5(\tilde{N} - 1)^2 + 13 \leq 5N + 13 \quad \text{and} \quad 7\tilde{L}^2 \leq 16(\tilde{L} - 1)^2 \leq 64L.$$

Therefore,  $\phi$  can be implemented by a Floor-ReLU network, with a fixed architecture independent of  $f$ , with width  $\max\{d, 2\tilde{N}^2 + 5\tilde{N}\} \leq \max\{d, 5N + 13\}$  and depth  $7d\tilde{L}^2 + 3 \leq 64dL + 3$ , as desired. So we finish the proof.  $\square$

## 5.2 Proof of auxiliary theorem

To prove Theorem 5.5, we first present the general ideas of the proof. Shortly speaking, we construct piecewise constant functions implemented by Floor-ReLU networks to approximate continuous functions on  $[0, 1]^d$ . There are four 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  $\{Q_\beta\}_{\beta \in \{0, 1, \dots, K-1\}^d}$ , and denote  $\mathbf{x}_\beta$  as the vertex of  $Q_\beta$  with minimum  $\|\cdot\|_1$  norm, where  $K$  is an integer determined later. See Figure 5.2 for the illustrations of  $Q_\beta$  and  $\mathbf{x}_\beta$  for any  $\beta \in \{0, 1, \dots, K-1\}^d$ .
2. Construct a Floor-ReLU sub-network to implement a vector-valued function  $\Phi_1 : \mathbb{R}^d \rightarrow \mathbb{R}^d$  projecting the whole cube  $Q_\beta$  to the index  $\beta$  for each  $\beta$ , *i.e.*,  $\Phi_1(\mathbf{x}) = \beta$  for all  $\mathbf{x} \in Q_\beta$  and each  $\beta \in \{0, 1, \dots, K-1\}^d$ .
3. Construct a Floor-ReLU sub-network to implement a function  $\phi_2 : \mathbb{R}^d \rightarrow \mathbb{R}$  mapping  $\beta \in \{0, 1, \dots, K-1\}^d$  approximately to  $\tilde{f}(\mathbf{x}_\beta)$  for each  $\beta$ , *i.e.*,  $\phi_2(\beta) \approx \tilde{f}(\mathbf{x}_\beta)$ . Then  $\phi_2 \circ \Phi_1(\mathbf{x}) = \phi_2(\beta) \approx \tilde{f}(\mathbf{x}_\beta)$  for any  $\mathbf{x} \in Q_\beta$  and each

$\beta \in \{0, 1, \dots, K-1\}^d$ , implying  $\tilde{\phi} := \phi_2 \circ \Phi_1$  approximates  $\tilde{f}$  within an error  $\mathcal{O}(\omega_f(1/K))$  on  $[0, 1]^d$ .

4. Re-scale and shift  $\tilde{\phi}$  to obtain the desired function  $\phi$  approximating  $f$  well and determine the final Floor-ReLU network to implement  $\phi$ .

It is not difficult to construct Floor-ReLU networks with the desired width and depth to implement  $\Phi_1$ . The most technical part is the construction of a Floor-ReLU network with the desired width and depth implementing  $\phi_2$ , which needs the following proposition based on the “bit extraction” technique introduced in [5].

**Proposition 5.6.** *Given any  $N, L \in \mathbb{N}^+$  and arbitrary  $\theta_m \in \{0, 1\}$  for  $m = 1, 2, \dots, N^L$ , there exists a function  $\phi$  implemented by a Floor-ReLU network, with a fixed architecture independent of  $\theta_m \in \{0, 1\}$  for  $m = 1, 2, \dots, N^L$ , with width  $2N + 2$  and depth  $7L - 2$  such that*

$$\phi(m) = \theta_m, \quad \text{for } m = 1, 2, \dots, N^L.$$

The proof of this proposition is presented in Section 5.3. It is easy to prove that the VC-dimension of Floor-ReLU networks with a constant width and depth  $\mathcal{O}(L)$  has a lower bound  $2^L$ . Such a lower bound is much larger than  $\mathcal{O}(L^2)$ , which is a VC-dimension upper bound of ReLU networks with the same width and depth due to Theorem 8 of [4]. This means Floor-ReLU networks are much more powerful than ReLU networks from the perspective of VC-dimension.

Now let us give the detailed proof of Theorem 5.5 as follows.

*Proof of Theorem 5.5.* The proof consists of four 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

$$\tilde{f} := \frac{f - f(\mathbf{0}) + \omega_f(\sqrt{d})}{2\omega_f(\sqrt{d})}. \quad (5.2)$$

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

Set  $K = N^L$ ,  $E_{K-1} = [\frac{K-1}{K}, 1]$ , and  $E_k = [\frac{k}{K}, \frac{k+1}{K})$  for  $k = 0, 1, \dots, K-2$ . Define

$$Q_\beta := \left\{ \mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d : x_j \in E_{\beta_j} \text{ for } j = 1, 2, \dots, d \right\},$$

for any  $\beta = (\beta_1, \beta_2, \dots, \beta_d) \in \{0, 1, \dots, K-1\}^d$ . See Figure 5.2 for the examples of  $Q_\beta$  and  $\mathbf{x}_\beta$  for any  $\beta \in \{0, 1, \dots, K-1\}^d$  with for  $K = 4$  and  $d = 1, 2$ .

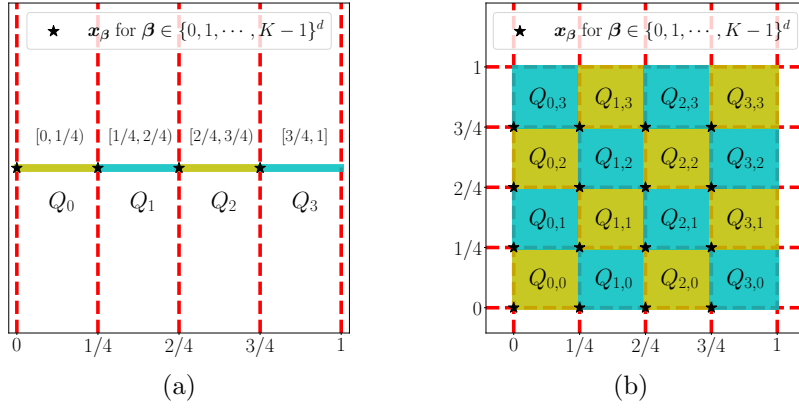


Figure 5.2: Illustrations of  $Q_\beta$  and  $\mathbf{x}_\beta$  for any  $\beta \in \{0, 1, \dots, K-1\}^d$ . (a)  $K = 4, d = 1$ . (b)  $K = 4, d = 2$ .

**Step 2:** Construct  $\Phi_1$  mapping  $\mathbf{x} \in Q_\beta$  to  $\beta$ .

Define a step function  $\phi_1$  as

$$\phi_1(x) := \lfloor -\sigma(-Kx + K - 1) + K - 1 \rfloor, \quad \text{for any } x \in \mathbb{R}. \textcircled{2}$$

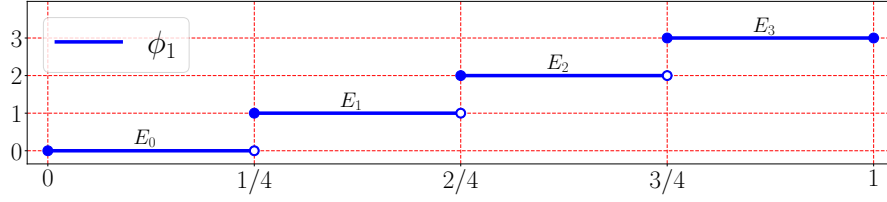
See Figure 5.3 for an illustration of  $\phi_1$  when  $K = 4$ . It follows from the definition of  $\phi_1$  that

$$\phi_1(x) = k, \quad \text{if } x \in E_k, \quad \text{for } k = 0, 1, \dots, K-1.$$

Define

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

<sup>②</sup>If we just define  $\phi_1(x) = \lfloor Kx \rfloor$ , then  $\phi_1(1) = K \neq K-1$  even though  $1 \in E_{K-1}$ .

Figure 5.3: An illustration of  $\phi_1$  on  $[0, 1]$  for the case  $K = 4$ .

Clearly, we have, for all  $\mathbf{x} \in Q_\beta$  and each  $\beta \in \{0, 1, \dots, K-1\}^d$ ,

$$\Phi_1(\mathbf{x}) = (\phi_1(x_1), \phi_1(x_2), \dots, \phi_1(x_d)) = (\beta_1, \beta_2, \dots, \beta_d) = \beta.$$

**Step 3:** Construct  $\phi_2$  mapping  $\beta \in \{0, 1, \dots, K-1\}^d$  approximately to  $\tilde{f}(\mathbf{x}_\beta)$ .

Using the idea of  $K$ -ary representation, we define

$$\psi_1(\mathbf{x}) := 1 + \sum_{j=1}^d x_j K^{j-1}, \quad \text{for any } \mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d.$$

It follows that  $\psi_1$  is a bijection from  $\{0, 1, \dots, K-1\}^d$  to  $\{1, 2, \dots, K^d\}$ .

Given any  $i \in \{1, 2, \dots, K^d\}$ , there exists a unique  $\beta \in \{0, 1, \dots, K-1\}^d$  such that  $i = \psi_1(\beta)$ . Then define

$$\xi_i := \tilde{f}(\mathbf{x}_\beta) \in [0, 1], \quad \text{for } i = \psi_1(\beta) \text{ and } \beta \in \{0, 1, \dots, K-1\}^d,$$

where  $\tilde{f}$  is the normalization of  $f$  defined in Equation (5.2). It follows that there exists  $\xi_{i,j} \in \{0, 1\}$  for  $j = 1, 2, \dots, NL$  such that

$$|\xi_i - \text{bin}_{0.\xi_{i,1}\xi_{i,2}\dots\xi_{i,NL}}| \leq 2^{-NL}, \quad \text{for } i = 1, 2, \dots, K^d.$$

By  $K^d = (NL)^d = N^{dL}$  and Proposition 5.6, there exists a function  $\psi_{2,j}$  implemented by a Floor-ReLU network, with a fixed architecture independent of  $\xi_{i,j}$  for

all  $i$ , with width  $2N + 2$  and depth  $7dL - 2$ , for each  $j = 1, 2, \dots, NL$ , such that

$$\psi_{2,j}(i) = \xi_{i,j}, \quad \text{for } i = 1, 2, \dots, K^d.$$

Define

$$\psi_2 := \sum_{j=1}^{NL} 2^{-j} \psi_{2,j} \quad \text{and} \quad \phi_2 := \psi_2 \circ \psi_1.$$

Then we have, for  $i = \psi_1(\boldsymbol{\beta})$  and  $\boldsymbol{\beta} \in \{0, 1, \dots, K - 1\}^d$ ,

$$\begin{aligned} |\tilde{f}(\mathbf{x}_{\boldsymbol{\beta}}) - \phi_2(\boldsymbol{\beta})| &= |\tilde{f}(\mathbf{x}_{\boldsymbol{\beta}}) - \psi_2(\psi_1(\boldsymbol{\beta}))| = |\xi_i - \psi_2(i)| = \left| \xi_i - \sum_{j=1}^{NL} 2^{-j} \psi_{2,j}(i) \right| \\ &= |\xi_i - \text{bin} 0.\xi_{i,1}\xi_{i,2} \cdots \xi_{i,NL}| \leq 2^{-NL}. \end{aligned} \quad (5.3)$$

**Step 4:** Determine the final network to implement the desired function  $\phi$ .

Define  $\tilde{\phi} := \phi_2 \circ \Phi_1$ , i.e., for any  $\mathbf{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$ ,

$$\tilde{\phi}(\mathbf{x}) = \phi_2 \circ \Phi_1(\mathbf{x}) = \phi_2(\phi_1(x_1), \phi_1(x_2), \dots, \phi_1(x_d)).$$

Note that  $|\mathbf{x} - \mathbf{x}_{\boldsymbol{\beta}}| \leq \frac{\sqrt{d}}{K}$  for any  $\mathbf{x} \in Q_{\boldsymbol{\beta}}$  and  $\boldsymbol{\beta} \in \{0, 1, \dots, K - 1\}^d$ . Then we have, for any  $\mathbf{x} \in Q_{\boldsymbol{\beta}}$  and  $\boldsymbol{\beta} \in \{0, 1, \dots, K - 1\}^d$ ,

$$\begin{aligned} |\tilde{f}(\mathbf{x}) - \tilde{\phi}(\mathbf{x})| &\leq |\tilde{f}(\mathbf{x}) - \tilde{f}(\mathbf{x}_{\boldsymbol{\beta}})| + |\tilde{f}(\mathbf{x}_{\boldsymbol{\beta}}) - \tilde{\phi}(\mathbf{x})| \\ &\leq \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{K}\right) + |\tilde{f}(\mathbf{x}_{\boldsymbol{\beta}}) - \phi_2(\Phi_1(\mathbf{x}))| \\ &\leq \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{K}\right) + |\tilde{f}(\mathbf{x}_{\boldsymbol{\beta}}) - \phi_2(\boldsymbol{\beta})| \leq \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{K}\right) + 2^{-NL}, \end{aligned}$$

where the last inequality comes from Equation (5.3).

Note that  $[0, 1]^d = \cup_{\boldsymbol{\beta} \in \{0, 1, \dots, K-1\}^d} Q_{\boldsymbol{\beta}}$ . Since  $\mathbf{x} \in Q_{\boldsymbol{\beta}}$  and  $\boldsymbol{\beta} \in \{0, 1, \dots, K - 1\}^d$  are arbitrary, we have

$$|\tilde{f}(\mathbf{x}) - \tilde{\phi}(\mathbf{x})| \leq \omega_{\tilde{f}}\left(\frac{\sqrt{d}}{K}\right) + 2^{-NL}, \quad \text{for any } \mathbf{x} \in [0, 1]^d.$$

Define

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

By  $K = N^L$  and  $\omega_f(r) = 2\omega_f(\sqrt{d}) \cdot \omega_{\tilde{f}}(r)$  for any  $r \geq 0$ , we have, for any  $\mathbf{x} \in [0, 1]^d$ ,

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

It remains to determine the width and depth of the Floor-ReLU network implementing  $\phi$ . Clearly,  $\phi_2$  can be implemented by the architecture in Figure 5.4.

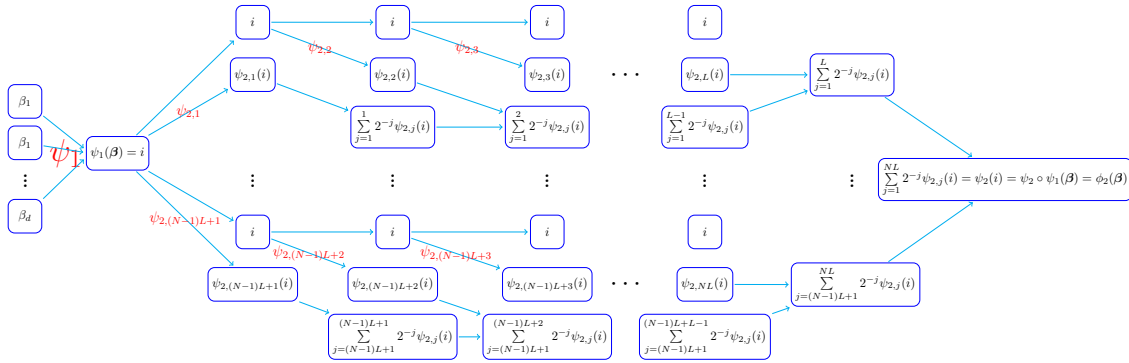


Figure 5.4: An illustration of the desired network architecture implementing  $\phi_2 = \psi_2 \circ \psi_1$  for any input  $\beta \in \{0, 1, \dots, K-1\}^d$ , where  $i = \psi_1(\beta)$ .

As we can see from Figure 5.4,  $\phi_2$  can be implemented by a Floor-ReLU network with width

$$N((2N+2)+3) = 2N^2 + 5N$$

and depth

$$1 + L((7dL-2)+1) + 1 = L(7dL-1) + 2.$$

Note that, for each  $j$ ,  $\psi_{2,j}$  is implemented by a Floor-ReLU network, with a fixed architecture independent of  $\xi_{i,j}$  that is essentially determined by the target function  $f$  for all  $i$ . Thus, the Floor-ReLU network implementing  $\phi_2$  has a fixed architecture

independent of  $f$ , as shown in Figure 5.4.

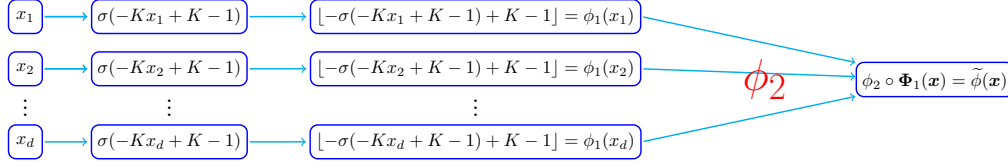


Figure 5.5: An illustration of the network architecture implementing  $\tilde{\phi} = \phi_2 \circ \Phi_1$  for any  $\mathbf{x} = (x_1, x_2, \dots, x_d) \in [0, 1]^d$ .

Note that  $\phi$  is defined via re-scaling and shifting  $\tilde{\phi}$ . As shown in Figure 5.5,  $\phi$  and  $\tilde{\phi}$  can be implemented by a Floor-ReLU network, with a fixed architecture independent of  $f$ , with width  $\max\{d, 2N^2 + 5N\}$  and depth  $2 + L(7dL - 1) + 2 \leq 7dL^2 + 3$ . So we finish the proof.  $\square$

### 5.3 Proof of key proposition for auxiliary theorem

The proof of Proposition 5.6 mainly relies on the “bit extraction” technique. As we shall see later, our key idea is to apply the Floor activation function to make “bit extraction” more powerful to reduce network sizes. In particular, Floor-ReLU networks can extract much more bits than ReLU networks with the same network size.

Let us first establish a basic lemma to extract  $1/N$  of total bits stored in a new binary number from an input binary number.

**Lemma 5.7.** *Given any  $J, N \in \mathbb{N}^+$ , there exists a function  $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}$  implemented by a Floor-ReLU network with width  $2N$  and depth 4 such that, for any  $\theta_j \in \{0, 1\}$ ,  $j = 1, 2, \dots, NJ$ , we have*

$$\phi(\text{bin}0.\theta_1\theta_2\cdots\theta_{NJ}, n) = \text{bin}0.\theta_{(n-1)J+1}\theta_{(n-1)J+2}\cdots\theta_{nJ}, \quad \text{for } n = 1, 2, \dots, N.$$

*Proof.* Given any  $\theta_j \in \{0, 1\}$  for  $j = 1, 2, \dots, NJ$ , denote

$$s = \text{bin}0.\theta_1\theta_2\cdots\theta_{NJ} \quad \text{and} \quad s_n = \text{bin}0.\theta_{(n-1)J+1}\theta_{(n-1)J+2}\cdots\theta_{nJ},$$



for  $n = 1, 2, \dots, N$ .

Then our goal is to construct a function  $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}$  implemented by a Floor-ReLU network with the desired width and depth that satisfies

$$\phi(s, n) = s_n, \quad \text{for } n = 1, 2, \dots, N.$$

Based on the properties of the binary representation, it is easy to check that

$$s_n = \lfloor 2^{nJ} s \rfloor / 2^J - \lfloor 2^{(n-1)J} s \rfloor, \quad \text{for } n = 1, 2, \dots, N. \quad (5.4)$$

With formulas to return  $s_1, s_2, \dots, s_N$ , it is still technical to construct a network outputting  $s_n$  for a given index  $n \in \{1, 2, \dots, N\}$ .

Set  $\delta = 2^{-J}$  and define  $g$  (see Figure 5.6) as

$$g(x) := \sigma\left(\sigma(x) - \sigma\left(\frac{x+\delta-1}{\delta}\right)\right), \quad \text{for any } x \in \mathbb{R}.$$

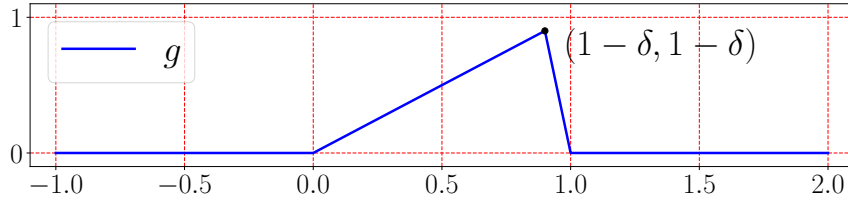


Figure 5.6: An illustration of  $g(x) = \sigma\left(\sigma(x) - \sigma\left(\frac{x+\delta-1}{\delta}\right)\right)$ .

Since  $s_n \in [0, 1 - \delta]$  for  $n = 1, 2, \dots, N$ , we have

$$s_n = \sum_{k=1}^N g(s_k + k - n), \quad \text{for } n = 1, 2, \dots, N. \quad (5.5)$$

As shown in Figure 5.7, the desired function  $\phi$  can be implemented by a Floor-ReLU network with width  $2N$  and depth 4. Moreover, it holds that

$$\phi(s, n) = s_n, \quad \text{for } n = 1, 2, \dots, N.$$

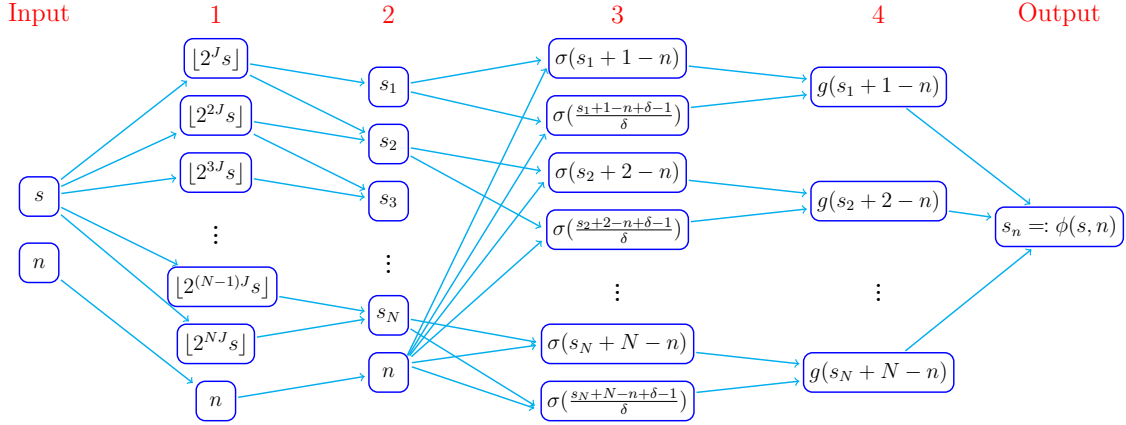


Figure 5.7: An illustration of the Floor-ReLU network implementing the desired function  $\phi$  based on Equation (5.4) and (5.5). All parameters in this network are essentially determined by Equation (5.4) and (5.5), which are valid no matter what  $\theta_1, \dots, \theta_{NJ} \in \{0, 1\}$  are. Thus, the desired function  $\phi$  implemented by this network is independent of  $\theta_1, \dots, \theta_{NJ} \in \{0, 1\}$ .

So we finish the proof.  $\square$

The next lemma constructs a Floor-ReLU network that can extract any bit from a binary number according to a specific index.

**Lemma 5.8.** *Given any  $N, L \in \mathbb{N}^+$ , there exists a function  $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}$  implemented by a Floor-ReLU network with width  $2N + 2$  and depth  $7L - 3$  such that, for any  $\theta_m \in \{0, 1\}$ ,  $m = 1, 2, \dots, N^L$ , we have*

$$\phi(\text{bin}0.\theta_1\theta_2\cdots\theta_{N^L}, m) = \theta_m, \quad \text{for } m = 1, 2, \dots, N^L.$$

*Proof.* The proof is based on repeated applications of Lemma 5.7. To be exact, we construct a sequence of functions  $\phi_1, \phi_2, \dots, \phi_L$  implemented by Floor-ReLU networks by induction to satisfy the following two conditions for each  $\ell \in \{1, 2, \dots, L\}$ .

- (i)  $\phi_\ell : \mathbb{R}^2 \rightarrow \mathbb{R}$  can be implemented by a Floor-ReLU network with width  $2N + 2$  and depth  $7\ell - 3$ .

(ii) For any  $\theta_m \in \{0, 1\}$ ,  $m = 1, 2, \dots, N^\ell$ , we have

$$\phi_\ell(\text{bin}0.\theta_1\theta_2\cdots\theta_{N^\ell}, m) = \text{bin}0.\theta_m, \quad \text{for } m = 1, 2, \dots, N^\ell.$$

First, consider the case  $\ell = 1$ . By Lemma 5.7 (set  $J = 1$  therein), there exists a function  $\phi_1$  implemented by a Floor-ReLU network with width  $2N \leq 2N + 2$  and depth  $4 = 7 - 3$  such that, for any  $\theta_m \in \{0, 1\}$ ,  $m = 1, 2, \dots, N$ , we have

$$\phi_1(\text{bin}0.\theta_1\theta_2\cdots\theta_N, m) = \text{bin}0.\theta_m, \quad \text{for } m = 1, 2, \dots, N.$$

It follows that Condition (i) and (ii) hold for  $\ell = 1$ .

Next, assume Condition (i) and (ii) hold for  $\ell = k$ . We would like to construct  $\phi_{k+1}$  to make Condition (i) and (ii) true for  $\ell = k + 1$ . By Lemma 5.7 (set  $J = N^k$  therein), there exists a function  $\psi$  implemented by a Floor-ReLU network with width  $2N$  and depth 4 such that, for any  $\theta_m \in \{0, 1\}$ ,  $m = 1, 2, \dots, N^{k+1}$ , we have

$$\psi(\text{bin}0.\theta_1\theta_2\cdots\theta_{N^{k+1}}, n) = \text{bin}0.\theta_{(n-1)N^k+1}\theta_{(n-1)N^k+2}\cdots\theta_{(n-1)N^k+N^k}, \quad (5.6)$$

for  $n = 1, 2, \dots, N$ . By the induction hypothesis, we have

- $\phi_k : \mathbb{R}^2 \rightarrow \mathbb{R}$  can be implemented by a Floor-ReLU network with width  $2N + 2$  and depth  $7k - 3$ .
- For any  $\theta_j \in \{0, 1\}$ ,  $j = 1, 2, \dots, N^k$ , we have

$$\phi_k(\text{bin}0.\theta_1\theta_2\cdots\theta_{N^k}, j) = \text{bin}0.\theta_j, \quad \text{for } j = 1, 2, \dots, N^k. \quad (5.7)$$

Given any  $m \in \{1, 2, \dots, N^{k+1}\}$ , there exist  $n \in \{1, 2, \dots, N\}$  and  $j \in \{1, 2, \dots, N^k\}$  such that  $m = (n - 1)N^k + j$ , and such  $n, j$  can be obtained by

$$n = \lfloor (m - 1)/N^k \rfloor + 1 \quad \text{and} \quad j = m - (n - 1)N^k. \quad (5.8)$$

Then the desired architecture of the Floor-ReLU network implementing  $\phi_{k+1}$  is shown in Figure 5.8.

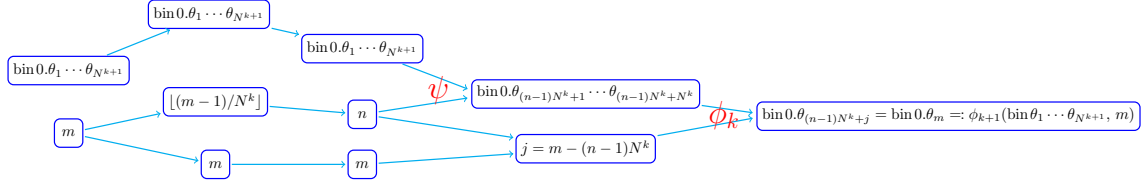


Figure 5.8: An illustration of the Floor-ReLU network architecture implementing  $\phi_{k+1}$ , based on Equation (5.6), (5.7), and (5.8) for any  $\theta_m \in \{0, 1\}$  and  $m \in \{1, 2, \dots, N^{k+1}\}$ .

Note that  $\psi$  can be implemented by a Floor-ReLU network with width  $2N$  and depth 4. Then the desired network implementing  $\phi_{k+1}$  is shown in Figure 5.8. Moreover, we have

- $\phi_{k+1} : \mathbb{R}^2 \rightarrow \mathbb{R}$  can be implemented by a Floor-ReLU network with width  $2N + 2$  and depth  $2 + 4 + 1 + (7k - 3) = 7(k + 1) - 3$ , which implies Condition (i) for  $\ell = k + 1$ .
- For any  $\theta_m \in \{0, 1\}$ ,  $m = 1, 2, \dots, N^{k+1}$ , we have

$$\phi_{k+1}(\text{bin} 0.\theta_1 \theta_2 \cdots \theta_{N^{k+1}}, m) = \text{bin} 0.\theta_m, \quad \text{for } m = 1, 2, \dots, N^{k+1}.$$

That is, Condition (ii) holds for  $\ell = k + 1$ .

So we finish the process of induction.

By the principle of induction, there exists a function  $\phi_L : \mathbb{R}^2 \rightarrow \mathbb{R}$  such that

- $\phi_L$  can be implemented by a Floor-ReLU network with width  $2N + 2$  and depth  $7L - 3$ .
- For any  $\theta_m \in \{0, 1\}$ ,  $m = 1, 2, \dots, N^L$ , we have

$$\phi_L(\text{bin} 0.\theta_1 \theta_2 \cdots \theta_{N^L}, m) = \text{bin} 0.\theta_m, \quad \text{for } m = 1, 2, \dots, N^L.$$

Finally, define  $\phi := 2\phi_L$ . Then  $\phi$  can also be implemented by a Floor-ReLU network with width  $2N+2$  and depth  $7L-3$ . Moreover, for any  $\theta_m \in \{0, 1\}$ ,  $m = 1, 2, \dots, N^L$ , we have

$$\phi(\text{bin}0.\theta_1\theta_2 \cdots \theta_{N^L}, m) = 2 \cdot \phi_L(\text{bin}0.\theta_1\theta_2 \cdots \theta_{N^L}, m) = 2 \cdot \text{bin}0.\theta_m = \theta_m,$$

for  $m = 1, 2, \dots, N^L$ . So we finish the proof.  $\square$

With Lemma 5.8 in hand, we are ready to prove Proposition 5.6.

*Proof of Proposition 5.6.* By Lemma 5.8, there exists a function  $\tilde{\phi} : \mathbb{R}^2 \rightarrow \mathbb{R}$  implemented by a Floor-ReLU network with width  $2N + 2$  and depth  $7L - 3$  such that, for any  $z_m \in \{0, 1\}$ ,  $m = 1, 2, \dots, N^L$ , we have

$$\tilde{\phi}(\text{bin}0.z_1z_2 \cdots z_{N^L}, m) = z_m, \quad \text{for } m = 1, 2, \dots, N^L.$$

Based on  $\theta_m \in \{0, 1\}$  for  $m = 1, 2, \dots, N^L$  given in Proposition 5.6, we define the final function  $\phi$  as

$$\phi(x) := \tilde{\phi}(\sigma(x \cdot 0 + \text{bin}0.\theta_1\theta_2 \cdots \theta_{N^L}), \sigma(x)), \quad \text{where } \sigma(x) = \max\{0, x\}.$$

Clearly,  $\phi$  can be implemented by a Floor-ReLU network, with a fixed architecture independent of  $\theta_m \in \{0, 1\}$  for  $m = 1, 2, \dots, N^L$ , with width  $2N + 2$  and depth  $1 + (7L - 3) = 7L - 2$ . In fact, only one parameter ( $\text{bin}0.\theta_1\theta_2 \cdots \theta_{N^L}$ ) of the network implementing  $\phi$  relies on  $\theta_m \in \{0, 1\}$  for  $m = 1, 2, \dots, N^L$ . Moreover, we have, for any  $m \in \{1, 2, \dots, N^L\}$ ,

$$\phi(m) = \tilde{\phi}(\sigma(m \cdot 0 + \text{bin}0.\theta_1\theta_2 \cdots \theta_{N^L}), \sigma(m)) = \tilde{\phi}(\text{bin}0.\theta_1\theta_2 \cdots \theta_{N^L}, m) = \theta_m.$$

So we finish the proof.  $\square$

We shall point out that only the properties of Floor on  $[0, \infty)$  are used in our

proof. Thus, the Floor can be replaced by the truncation function that can be easily implemented by truncating the decimal part.

Finally, we would like to remark that the key reason Floor-ReLU networks can attain much better approximation errors than ReLU networks is that Floor has infinite (constant) pieces, while ReLU has only two (linear) pieces. Thus, roughly speaking, one Floor activation function can do what many ReLU activation functions do in our construction. For this reason, compared to ReLU networks, Floor-ReLU networks attain significantly better approximation errors. In fact, one may replace Floor by other activation functions with “many pieces”. For example, it is shown in [60] that ReLU/Sin-activated networks can also attain nearly exponential approximation errors.

## Conclusion

This dissertation aims to study the approximation power of ReLU networks and Floor-ReLU networks. Based on the idea of function compositions, we construct ReLU networks to uniformly approximate polynomials, Hölder continuous functions, general continuous functions, and smooth functions on a  $d$ -dimensional hypercube  $[0, 1]^d$  with (nearly optimal) approximation errors. All the approximation error estimates are characterized by the width and depth simultaneously and have explicit formulas for the prefactors. Meanwhile, the optimality of the approximation by ReLU networks is discussed via studying the connection between the approximation error and VC-dimension.

To overcome the limitation of ReLU networks that (nearly) exponential approximation errors can only be attained for polynomials among all function spaces considered, we introduce Floor-ReLU networks. It is proved by construction that nearly exponential approximation errors can be attained when using Floor-ReLU networks with fixed architectures to uniformly approximate (Hölder) continuous functions on  $[0, 1]^d$ . In other words, approximation errors are improved from polynomial ones to nearly exponential ones via adding a simple activation function (Floor) to ReLU networks. The optimality of the approximation by Floor-ReLU networks is not discussed due to the nearly exponential approximation errors. All these results stated above completely solve the three problems (Problem 1, 2, and 3) listed in Chapter 1

for certain function spaces.

Finally, we would like to remark that our analysis is only for the feed-forward fully connected neural networks with two types of activation functions: ReLU and Floor. It would be an interesting direction to generalize our results to neural networks with other architectures (*e.g.*, convolutional neural networks and residual networks) and activation functions (*e.g.*, tanh and sigmoid functions). These will be left as future work.



---

## Bibliography

---

- [1] O. Abdel-Hamid, A. Mohamed, H. Jiang, L. Deng, G. Penn, and D. Yu. Convolutional neural networks for speech recognition. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 22(10):1533–1545, Oct 2014.
- [2] M. Anthony and P. L. Bartlett. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, New York, NY, USA, 1st edition, 2009.
- [3] A. R. Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information Theory*, 39(3):930–945, May 1993.
- [4] P. L. Bartlett, N. Harvey, C. Liaw, and A. Mehrabian. Nearly-tight VC-dimension and pseudodimension bounds for piecewise linear neural networks. *Journal of Machine Learning Research*, 20(63):1–17, 2019.
- [5] P. L. Bartlett, V. Maiorov, and R. Meir. Almost linear VC-dimension bounds for piecewise polynomial networks. *Neural Computation*, 10(8):2159–2173, 1998.
- [6] Y. Bengio, N. Léonard, and A. C. Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *CoRR*, abs/1308.3432, 2013.

- [7] 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(8):1553–1565, Aug 2014.
- [8] E. K. Blum and L. K. Li. Approximation theory and feedforward networks. *Neural Networks*, 4(4):511–515, 1991.
- [9] Y. Boo, S. Shin, and W. Sung. Quantized neural networks: Characterization and holistic optimization. In *2020 IEEE Workshop on Signal Processing Systems (SiPS)*, pages 1–6, 2020.
- [10] D. S. Broomhead and D. Lowe. Multivariable functional interpolation and adaptive networks. *Complex Systems* 2, pages 321–355, 1988.
- [11] J. A. Carrillo, S. Jin, L. Li, and Y. Zhu. A consensus-based global optimization method for high dimensional machine learning problems. *arXiv e-prints*, page arXiv:1909.09249, Sep 2019.
- [12] S. Chen and D. Donoho. Basis pursuit. In *Proceedings of 1994 28th Asilomar Conference on Signals, Systems and Computers*, volume 1, pages 41–44, Oct 1994.
- [13] D. Costarelli and A. R. Sambucini. Saturation classes for max-product neural network operators activated by sigmoidal functions. *Results in Mathematics*, 72(3):1555–1569, 2017.
- [14] D. Costarelli and A. R. Sambucini. Approximation results in Orlicz spaces for sequences of Kantorovich max-product neural network operators. *Results in Mathematics*, 73(1):1–15, 2018.
- [15] D. Costarelli and G. Vinti. Convergence for a family of neural network operators in Orlicz spaces. *Mathematische Nachrichten*, 290(2-3):226–235, 2017.
- [16] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems*, 2:303–314, 1989.

- [17] I. Daubechies, R. DeVore, S. Foucart, B. Hanin, and G. Petrova. Nonlinear approximation and (deep) ReLU networks. *arXiv e-prints*, page arXiv:1905.02199, May 2019.
- [18] R. DEVORE and A. RON. Approximation using scattered shifts of a multivariate function. *Transactions of the American Mathematical Society*, 362(12):6205–6229, 2010.
- [19] R. A. DeVore. Nonlinear approximation. *Acta Numerica*, 7:51–150, 1998.
- [20] W. E, J. Han, and A. Jentzen. Deep learning-based numerical methods for high-dimensional parabolic partial differential equations and backward stochastic differential equations. *Communications in Mathematics and Statistics*, 5(4):349–380, Dec 2017.
- [21] W. E and Q. Wang. Exponential convergence of the deep neural network approximation for analytic functions. *Science China Mathematics*, 61:1733–1740, 2018.
- [22] J. Han, A. Jentzen, and W. E. Solving high-dimensional partial differential equations using deep learning. *Proceedings of the National Academy of Sciences*, 115(34):8505–8510, 2018.
- [23] T. Hangelbroek and A. Ron. Nonlinear approximation using gaussian kernels. *Journal of Functional Analysis*, 259(1):203–219, 2010.
- [24] K. Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2):251–257, 1991.
- [25] K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359–366, 1989.
- [26] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio. Quantized neural networks: Training neural networks with low precision weights and activations. *Journal of Machine Learning Research*, 18(1):6869–6898, Jan 2017.

- 
- [27] K. Kawaguchi. Deep learning without poor local minima. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 586–594. Curran Associates, Inc., 2016.
  - [28] K. Kawaguchi and Y. Bengio. Depth with nonlinearity creates no bad local minima in resnets. *Neural Networks*, 118:167–174, 2019.
  - [29] M. J. Kearns and R. E. Schapire. Efficient distribution-free learning of probabilistic concepts. *Journal of Computer and System Sciences*, 48(3):464–497, Jun 1994.
  - [30] J. Kennedy and R. Eberhart. Particle swarm optimization. In *Proceedings of ICNN’95 - International Conference on Neural Networks*, volume 4, pages 1942–1948, 1995.
  - [31] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.
  - [32] V. Kůrková. Kolmogorov’s theorem and multilayer neural networks. *Neural Networks*, 5(3):501–506, 1992.
  - [33] G. Lewicki and G. Marino. Approximation of functions of finite variation by superpositions of a sigmoidal function. *Applied Mathematics Letters*, 17(10):1147–1152, 2004.
  - [34] S. Liang and R. Srikant. Why deep neural networks for function approximation? *arXiv e-prints*, page arXiv:1610.04161, Oct 2016.
  - [35] S. Lin, X. Liu, Y. Rong, and Z. Xu. Almost optimal estimates for approximation and learning by radial basis function networks. *Machine Learning*, 95(2):147–164, May 2014.

- [36] Y. Lin, M. Lei, and L. Niu. Optimization strategies in quantized neural networks: A review. In *2019 International Conference on Data Mining Workshops (ICDMW)*, pages 385–390, 2019.
- [37] B. Llanas and F. Sainz. Constructive approximate interpolation by neural networks. *Journal of Computational and Applied Mathematics*, 188(2):283–308, 2006.
- [38] J. Lu, Z. Shen, H. Yang, and S. Zhang. Deep network approximation for smooth functions. *arXiv e-prints*, page arXiv:2001.03040, Jan 2020.
- [39] Z. Lu, H. Pu, F. Wang, Z. Hu, and L. Wang. The expressive power of neural networks: A view from the width. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30, pages 6231–6239. Curran Associates, Inc., 2017.
- [40] V. Maiorov and A. Pinkus. Lower bounds for approximation by MLP neural networks. *Neurocomputing*, 25(1):81–91, 1999.
- [41] S. G. Mallat and Z. Zhang. Matching pursuits with time-frequency dictionaries. *IEEE Transactions on Signal Processing*, 41(12):3397–3415, Dec 1993.
- [42] H. Montanelli and Q. Du. New error bounds for deep ReLU networks using sparse grids. *SIAM Journal on Mathematics of Data Science*, 1(1):78–92, 2019.
- [43] H. Montanelli, H. Yang, and Q. Du. Deep ReLU networks overcome the curse of dimensionality for bandlimited functions. *arXiv e-prints*, page arXiv:1903.00735, Mar 2019.
- [44] G. F. Montufar, R. Pascanu, K. Cho, and Y. Bengio. On the number of linear regions of deep neural networks. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2924–2932. Curran Associates, Inc., 2014.

- [45] Q. Nguyen and M. Hein. The loss surface of deep and wide neural networks. In D. Precup and Y. W. Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 2603–2612, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR.
- [46] J. A. A. Opschoor, C. Schwab, and J. Zech. Exponential ReLU DNN expression of holomorphic maps in high dimension. Technical Report 2019-35, Seminar for Applied Mathematics, ETH Zürich, Switzerland., 2019.
- [47] J. Park and I. W. Sandberg. Universal approximation using radial-basis-function networks. *Neural Computation*, 3(2):246–257, Jun 1991.
- [48] P. Petersen and F. Voigtlaender. Optimal approximation of piecewise smooth functions using deep ReLU neural networks. *Neural Networks*, 108:296 – 330, 2018.
- [49] P. Petrushev. Multivariate n-term rational and piecewise polynomial approximation. *Journal of Approximation Theory*, 121(1):158–197, 2003.
- [50] R. Pinnau, C. Totzeck, O. Tse, and S. Martin. A consensus-based model for global optimization and its mean-field limit. *Mathematical Models and Methods in Applied Sciences*, 27(01):183–204, 2017.
- [51] A. Sakurai. Tight bounds for the VC-dimension of piecewise polynomial networks. In *Advances in Neural Information Processing Systems*, pages 323–329. Neural information processing systems foundation, 1999.
- [52] Z. Shen, H. Yang, and S. Zhang. Nonlinear approximation via compositions. *Neural Networks*, 119:74–84, 2019.
- [53] Z. Shen, H. Yang, and S. Zhang. Deep network approximation characterized by number of neurons. *Communications in Computational Physics*, 28(5):1768–1811, 2020.

- [54] Z. Shen, H. Yang, and S. Zhang. Deep network with approximation error being reciprocal of width to power of square root of depth. *arXiv e-prints*, page arXiv:2006.12231, Jun 2020.
- [55] T. Suzuki. Adaptivity of deep ReLU network for learning in Besov and mixed smooth Besov spaces: optimal rate and curse of dimensionality. In *International Conference on Learning Representations*, 2019.
- [56] P. Wang, Q. Hu, Y. Zhang, C. Zhang, Y. Liu, and J. Cheng. Two-step quantization for low-bit neural networks. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4376–4384, 2018.
- [57] T. F. Xie and F. L. Cao. The rate of approximation of gaussian radial basis neural networks in continuous function space. *Acta Mathematica Sinica, English Series*, 29(2):295–302, Feb 2013.
- [58] D. Yarotsky. Error bounds for approximations with deep ReLU networks. *Neural Networks*, 94:103–114, 2017.
- [59] D. Yarotsky. Optimal approximation of continuous functions by very deep ReLU networks. In S. Bubeck, V. Perchet, and P. Rigollet, editors, *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, pages 639–649. PMLR, Jul 2018.
- [60] D. Yarotsky and A. Zhevnerchuk. The phase diagram of approximation rates for deep neural networks. *arXiv e-prints*, page arXiv:1906.09477, Jun 2019.
- [61] P. Yin, J. Lyu, S. Zhang, S. Osher, Y. Qi, and J. Xin. Understanding straight-through estimator in training activation quantized neural nets. *arXiv e-prints*, page arXiv:1903.05662, Mar 2019.