

Deep Network Approximation: Achieving Arbitrary Accuracy with Fixed Number of Neurons

Zuwei Shen

*Department of Mathematics
National University of Singapore*

MATZUOWS@NUS.EDU.SG

Haizhao Yang

*Department of Mathematics
University of Maryland, College Park*

HZYANG@UMD.EDU

Shijun Zhang*

*Department of Mathematics
National University of Singapore*

ZHANGSHIJUN@U.NUS.EDU

Abstract

1 This paper develops simple feed-forward neural networks that achieve the universal
2 approximation property for all continuous functions with a fixed finite number of neurons.
3 These neural networks are simple because they are designed with a simple, computable,
4 and continuous activation function σ leveraging a triangular-wave function and the softsign
5 function. We first prove that σ -activated networks with width $36d(2d + 1)$ and depth 11
6 can approximate any continuous function on a d -dimensional hypercube within an arbi-
7 trarily small error. Hence, for supervised learning and its related regression problems, the
8 hypothesis space generated by these networks with a size not smaller than $36d(2d + 1) \times 11$
9 is dense in the continuous function space $C([a, b]^d)$ and therefore dense in the Lebesgue
10 spaces $L^p([a, b]^d)$ for $p \in [1, \infty)$. Furthermore, we show that classification functions arising
11 from image and signal classification are in the hypothesis space generated by σ -activated
12 networks with width $36d(2d + 1)$ and depth 12 when there exist pairwise disjoint bounded
13 closed subsets of \mathbb{R}^d such that the samples of the same class are located in the same subset.
14 Finally, we use numerical experimentation to show that replacing the rectified linear unit
15 (ReLU) activation function by ours would improve the experiment results.

16 **Keywords:** universal approximation property, fixed-size neural network, classification
17 function, periodic function, nonlinear approximation

18 1. Introduction

19 Deep neural networks have been widely used in data science and artificial intelligence. Their
20 tremendous successes in various applications have motivated extensive research to establish
21 the theoretical foundation of deep learning. Understanding the approximation capacity
22 of deep neural networks is one of the keys to revealing the power of deep learning. The
23 most basic layers of deep neural networks are nonlinear functions as the composition of
24 an affine linear transform and a nonlinear activation function. The composition of these
25 simple nonlinear functions can generate a complicated deep neural network with powerful
26 approximation capacity, which is the key difference from classic approximation tools. In
27 this paper, we show that the hypothesis space of deep neural networks generated from

* Corresponding author.

28 the composition of 11 such simple nonlinear functions is dense in the continuous function
 29 space $C([a, b]^d)$ when the affine linear transforms are parameterized with $\mathcal{O}(d^2)$ non-zero
 30 parameters in total and the nonlinear activation function is constructed from a simple
 31 triangular-wave function and the softsign function.

32 1.1 Main Results

33 One of the key elements of a neural network is its activation functions. Searching for simple
 34 activation functions enabling powerful approximation capacity of neural networks is an
 35 important mathematical problem that probably originated in the Kolmogorov superposition
 36 theorem (KST) (Kolmogorov, 1957) for Hilbert’s 13-th problem, where a two-hidden-layer
 37 neural network with $\mathcal{O}(d)$ neurons and complicated activation functions depending on the
 38 target functions are constructed to represent an arbitrary function in $C([0, 1]^d)$. Since then,
 39 whether simple and computable activation functions independent of the target function
 40 exist to make the space of neural networks with $\mathcal{O}(d)$ neurons dense in $C([0, 1]^d)$ or even
 41 equal to $C([0, 1]^d)$ has been an open problem. A function $\varrho : \mathbb{R} \rightarrow \mathbb{R}$ is said to be a universal
 42 activation function (UAF) if the function space generated by ϱ -activated networks with $C_{\varrho, d}$
 43 neurons is dense in $C([0, 1]^d)$, where $C_{\varrho, d}$ is a constant determined by ϱ and d . That is, if
 44 ϱ is a UAF, then ϱ -activated networks with $C_{\varrho, d}$ neurons can approximate any continuous
 45 function within an arbitrary error on $[0, 1]^d$ by only adjusting the parameters.

46 In this paper, we first construct a simple and computable example of UAFs. As a typical
 47 and simple UAF, this activation function is called elementary universal activation function
 48 (EUAF), and the corresponding networks are called EUAF networks. Then, we prove that
 49 the function space generated by EUAF networks with $\mathcal{O}(d^2)$ neurons is dense in $C([a, b]^d)$.
 50 Furthermore, it is shown that EUAF networks with $\mathcal{O}(d^2)$ neurons can exactly represent
 51 d -dimensional classification functions.

52 While a good activation function should be simple and numerically implementable, the
 53 neural network activated by it should be able to approximate continuous functions well
 54 with a manageable size. Considering these requirements and motivated by previous works
 55 (Yarotsky and Zhevnerchuk, 2020; Shen et al., 2021a,b), the activation function to be cho-
 56 sen should have appropriate nonlinearity, periodicity, and the capacity to reproduce step
 57 functions. It is challenging to find a single activation function with all these properties.
 58 Here, we propose an activation function with all required properties by using two simple
 59 functions σ_1 and σ_2 defined below.

60 Let σ_1 be the continuous triangular-wave function with period 2, i.e.,

$$61 \quad \sigma_1(x) := |x| \quad \text{for any } x \in [-1, 1]$$

62 and $\sigma_1(x + 2) = \sigma_1(x)$ for any $x \in \mathbb{R}$. Alternatively, σ_1 can also be written as:

$$63 \quad \sigma_1(x) = \left| x - 2 \left\lfloor \frac{x+1}{2} \right\rfloor \right| \quad \text{for any } x \in \mathbb{R}, \quad \text{where } \lfloor \cdot \rfloor \text{ is the floor function.}$$

64 Clearly, σ_1 is periodic and $x - \sigma_1(x)$ is a continuous variant of the floor function as desired.

65 To introduce high nonlinearity, let σ_2 be the softsign activation function commonly used
 66 in machine learning (Turian et al., 2009; Le and Zuidema, 2015):

$$67 \quad \sigma_2(x) := \frac{x}{|x| + 1} \quad \text{for any } x \in \mathbb{R}.$$

68 Then the activation function σ is defined as:

$$69 \quad \sigma(x) := \begin{cases} \sigma_1(x) & \text{for } x \in [0, \infty), \\ \sigma_2(x) & \text{for } x \in (-\infty, 0). \end{cases} \quad (1)$$

70 See an illustration of σ in Figure 1. This activation function σ is used to construct powerful
71 neural networks in this paper.

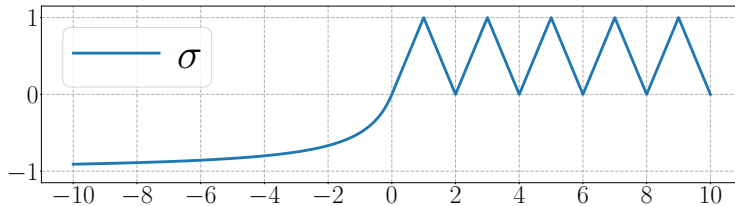


Figure 1: An illustration of σ on $[-10, 10]$.

72 As we shall see later, the periodicity of the triangular-wave function σ_1 and the (high)
73 nonlinearity of the softsign function σ_2 play crucial roles in the proofs of our main results.
74 One may find more details Section 2.2, which provides the ideas of proving our main results.
75 Observe that σ_1 is an even function and σ_2 is an odd function, i.e., $\sigma(x) = \sigma_1(x) = \sigma_1(-x)$
76 for any $x \geq 0$ and $-\sigma(-x) = -\sigma_2(-x) = \sigma_2(x)$ for any $x \geq 0$. This implies that $\sigma(x)$
77 and $-\sigma(-x)$ with $x \geq 0$ have both required periodicity and nonlinearity features and play
78 the same roles as $\sigma_1(x)$ and $\sigma_2(x)$, respectively. These requirements lead to our choice
79 of σ as the activation function. If allowed to be more complicated, one can design many
80 other UAFs satisfying stronger requirements for various applications. For example, the
81 idea of designing a C^s UAF is given in Section 4.1 and a sigmoidal UAF (see Figure 8) is
82 constructed in Section 4.2.

83 With the activation function σ in hand, let us introduce the network (architecture)
84 using σ as the activation function, called σ -activated network (architecture). To be precise,
85 a σ -activated network with a (vector) input $\mathbf{x} \in \mathbb{R}^d$, an output $\Phi(\mathbf{x}, \boldsymbol{\theta}) \in \mathbb{R}$, and $L \in \mathbb{N}^+$
86 hidden layers can be briefly described as follows:

$$87 \quad \mathbf{x} = \tilde{\mathbf{h}}_0 \xrightarrow{\mathbf{A}_0, \mathbf{b}_0} \mathbf{h}_1 \xrightarrow{\sigma} \tilde{\mathbf{h}}_1 \quad \cdots \quad \xrightarrow{\mathbf{A}_{L-1}, \mathbf{b}_{L-1}} \mathbf{h}_L \xrightarrow{\sigma} \tilde{\mathbf{h}}_L \xrightarrow{\mathbf{A}_L, \mathbf{b}_L} \mathbf{h}_{L+1} = \Phi(\mathbf{x}, \boldsymbol{\theta}), \quad (2)$$

88 where $N_0 = d \in \mathbb{N}^+$, $N_1, N_2, \dots, N_L \in \mathbb{N}^+$, $N_{L+1} = 1$, $\mathbf{A}_i \in \mathbb{R}^{N_{i+1} \times N_i}$ and $\mathbf{b}_i \in \mathbb{R}^{N_{i+1}}$ are
89 the weight matrix and the bias vector in the i -th affine linear transform \mathcal{L}_i , respectively,
90 i.e.,

$$91 \quad \mathbf{h}_{i+1} = \mathbf{A}_i \cdot \tilde{\mathbf{h}}_i + \mathbf{b}_i =: \mathcal{L}_i(\tilde{\mathbf{h}}_i) \quad \text{for } i = 0, 1, \dots, L$$

92 and

$$93 \quad \tilde{h}_{i,j} = \sigma(h_{i,j}) \quad \text{for } j = 1, 2, \dots, N_i \text{ and } i = 1, 2, \dots, L.$$

94 Here, $\tilde{h}_{i,j}$ and $h_{i,j}$ are the j -th entries of $\tilde{\mathbf{h}}_i$ and \mathbf{h}_i , respectively, for $j = 1, 2, \dots, N_i$ and $i =$
95 $1, 2, \dots, L$. $\boldsymbol{\theta}$ is a fattened vector consisting of all parameters in $\mathbf{A}_0, \mathbf{b}_0, \mathbf{A}_1, \mathbf{b}_1, \dots, \mathbf{A}_L, \mathbf{b}_L$.

96 With a slight abuse of notation, σ can be applied to a vector elementwisely, i.e., given
97 any $k \in \mathbb{N}^+$,

$$98 \quad \sigma(\mathbf{y}) = [\sigma(y_1), \sigma(y_2), \dots, \sigma(y_k)]^T \quad \text{for any } \mathbf{y} = [y_1, y_2, \dots, y_k]^T \in \mathbb{R}^k.$$

99 Then Φ can be represented in a form of function compositions as follows:

$$100 \quad \Phi(\mathbf{x}, \boldsymbol{\theta}) = \mathcal{L}_L \circ \sigma \circ \mathcal{L}_{L-1} \circ \cdots \circ \sigma \circ \mathcal{L}_1 \circ \sigma \circ \mathcal{L}_0(\mathbf{x}) \quad \text{for any } \mathbf{x} \in \mathbb{R}^d.$$

101 Given $N, L \in \mathbb{N}^+$, let $\Phi_{N,L}(\mathbf{x}, \boldsymbol{\theta})$ denote the σ -activated network architecture $\Phi(\mathbf{x}, \boldsymbol{\theta})$ in
 102 Equation (2) with $N_1 = N_2 = \cdots = N_L = N$. Let

$$103 \quad W = W_{d,N,L} = d \times N + N + (N \times N + N) \times (L - 1) + N \times 1 + 1 = \mathcal{O}(dN + N^2L)$$

104 be the total number of parameters in $\Phi_{N,L}(\mathbf{x}, \boldsymbol{\theta})$, i.e., $\boldsymbol{\theta} \in \mathbb{R}^W$.

105 Define the hypothesis space $\mathcal{H}_d(N, L)$ as the function space generated by d -input EUAF
 106 networks with width N and depth L , i.e.,

$$107 \quad \mathcal{H}_d(N, L) := \left\{ \phi : \phi(\mathbf{x}) = \Phi_{N,L}(\mathbf{x}, \boldsymbol{\theta}) \text{ for any } \mathbf{x} \in \mathbb{R}^d, \quad \boldsymbol{\theta} \in \mathbb{R}^W \right\}. \quad (3)$$

108 Let $C([a, b]^d)$ be the space of all continuous functions $f : [a, b]^d \rightarrow \mathbb{R}$ with the maximum
 109 norm. Our first main result, Theorem 1 below, shows that EUAF networks with a fixed
 110 size $\mathcal{O}(d^2)$ enjoy the universal approximation property by only adjusting their parameters.

111 **Theorem 1.** *Let $f \in C([a, b]^d)$ be a continuous function and $\mathcal{H}_d(N, L)$ be the hypothesis*
 112 *space defined in Equation (3) with $N = 36d(2d + 1)$ and $L = 11$. Then, for an arbitrary*
 113 *$\varepsilon > 0$, there exists $\phi \in \mathcal{H}_d(N, L)$ such that*

$$114 \quad |\phi(\mathbf{x}) - f(\mathbf{x})| < \varepsilon \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

115 To prove Theorem 1, we first summarize key proof ideas in Section 2.2 and then present
 116 the detailed proof later in Section 5.1.

117 **Remark.** *The network realizing ϕ in Theorem 1 has*

$$118 \quad d \times N + N + (N \times N + N) \times (L - 1) + N \times 1 + 1 \sim d^4$$

119 *parameters, where $N = 36d(2d + 1)$ and $L = 11$. However, as shown in our constructive*
 120 *proof of Theorem 1, it is enough to adjust $5437(d + 1)(2d + 1) = \mathcal{O}(d^2) \ll d^4$ parameters*
 121 *and set all the others to 0.*

122 Since for an arbitrary $M > 0$, $2M\sigma(\frac{x+M}{2M}) - M = x$ for all $x \in [-M, M]$, we can
 123 manually add hidden layers to EUAF networks without changing the output. This leads to
 124 the following immediate corollary of Theorem 1.

125 **Corollary 2.** *Assume $N \geq 36d(2d + 1)$ and $L \geq 11$. Then the hypothesis space $\mathcal{H}_d(N, L)$*
 126 *defined in Equation (3) is dense in $C([a, b]^d)$.*

127 The stable and accurate approximation of discontinuities has many real-world applica-
 128 tions and has been widely studied (Bernholdt et al., 2019; Beck et al., 2020; Gupta et al.,
 129 2020; Gedeon et al., 2021; Hu et al., 2021). Most of common discontinuous functions are in
 130 the Lebesgue spaces $L^p([a, b]^d)$ for $p \in [1, \infty)$. Let us consider the denseness of our hypoth-
 131 esis space in these function spaces. Since $C([a, b]^d)$ is dense in $L^p([a, b]^d)$ for $p \in [1, \infty)$,
 132 the hypothesis space in Corollary 2 is also dense in $L^p([a, b]^d)$ as shown in the following
 133 corollary.

134 **Corollary 3.** Assume $N \geq 36d(2d+1)$, $L \geq 11$, and $p \in [1, \infty)$. Then the hypothesis space
 135 $\mathcal{H}_d(N, L)$ defined in Equation (3) is dense in $L^p([a, b]^d)$.

136 This corollary implies that, for $f \in L^p([a, b]^d)$ and an arbitrary $\varepsilon > 0$, there exists
 137 $\phi \in \mathcal{H}_d(N, L)$ such that $\|\phi - f\|_{L^p([a, b]^d)} < \varepsilon$.

138 One can ask whether the arbitrary error $\varepsilon > 0$ in Theorem 1 can be further reduced to
 139 0. This is not true in general, but it is true for a class of interesting functions widely used in
 140 image classification. Given any pairwise disjoint bounded closed subsets $E_1, E_2, \dots, E_J \subseteq$
 141 \mathbb{R}^d , define “the classification function space” of these subsets as

$$142 \quad \mathcal{C}_d(E_1, E_2, \dots, E_J) := \left\{ f : f = \sum_{j=1}^J r_j \cdot \mathbb{1}_{E_j} \text{ for any } r_1, r_2, \dots, r_J \in \mathbb{Q} \right\},$$

143 where $\mathbb{1}_{E_n}$ is the indicator function of E_j for each j . Our second main result, Theorem 4
 144 below, shows that each element of $\mathcal{C}_d(E_1, E_2, \dots, E_J)$ can be exactly represented by a σ -
 145 activated network with $\mathcal{O}(d^2)$ neurons in $\bigcup_{j=1}^J E_j$.

146 **Theorem 4.** Let $E_1, E_2, \dots, E_J \subseteq \mathbb{R}^d$ be pairwise disjoint bounded closed subsets and
 147 $\mathcal{H}_d(N, L)$ be the hypothesis space defined in Equation (3) with $N = 36d(2d+1)$ and $L = 12$.
 148 Then, for an arbitrary $f \in \mathcal{C}_d(E_1, E_2, \dots, E_J)$, there exists $\phi \in \mathcal{H}_d(N, L)$ such that

$$149 \quad \phi(\mathbf{x}) = f(\mathbf{x}) \quad \text{for any } \mathbf{x} \in \bigcup_{j=1}^J E_j.$$

150 **Remark.** The network realizing ϕ in Theorem 4 has

$$151 \quad d \times N + N + (N \times N + N) \times (L - 1) + N \times 1 + 1 \sim d^4$$

152 parameters, where $N = 36d(2d+1)$ and $L = 12$. However, as shown in our constructive
 153 proof of Theorem 4 in Section 5.2, it is enough to adjust $5509(d+1)(2d+1) = \mathcal{O}(d^2) \ll d^4$
 154 parameters and set all the others to 0.

155 For a general function space \mathcal{F} , define $\mathcal{F}|_E := \{f|_E : f \in \mathcal{F}\}$, where $f|_E$ is the function
 156 achieved via limiting f on E . Then, we have a corollary of Theorem 4 as follows.

157 **Corollary 5.** Let $E_1, E_2, \dots, E_J \subseteq \mathbb{R}^d$ be pairwise disjoint bounded closed subsets and
 158 $\mathcal{H}_d(N, L)$ be the hypothesis space defined in Equation (3). If $N \geq 36d(2d+1)$ and $L \geq 12$,
 159 then

$$160 \quad \mathcal{C}_d(E_1, E_2, \dots, E_J)|_E \subseteq \mathcal{H}_d(N, L)|_E \quad \text{with } E = \bigcup_{j=1}^J E_j.$$

161 One of the most successful applications of deep learning is image and signal classifica-
 162 tion. In supervised classification problems, given a few samples and their labels (usually
 163 integers), the goal of the task is to learn how to assign a label to a new sample. For exam-
 164 ple, in binary classification via deep learning, a neural network is trained based on given
 165 samples (and labels) to approximate a classification function mapping one class of samples

166 to 0 and the other class of samples to 1. Theorem 4 (or Corollary 5) implies that the clas-
167 sification function can be exactly realized by an EUAF network with a size depending only
168 on the dimension of the problem domain via adjusting its parameters. This means that the
169 best approximation error of EUAF networks to classification functions in the classification
170 problem is 0.

171 We remark that, in the worst scenario, there might exist complicated high-dimensional
172 functions such that, the parameters of the EUAF network in Theorem 1 (or 4) require
173 high computer precision for storage, and the precision might be exponentially high in the
174 problem dimension. We refer to this as the curse of memory, which may make Theorem 1
175 and 4 less interesting in real-world applications, though the number of parameters can be
176 very small. The key question to be addressed is how rare the curse of memory would happen
177 in real-world applications. If the target functions in real-world applications typically have
178 no curse of memory with a high probability, then EUAF networks would be very useful in
179 real-world applications. In future work, we will explore the statistical characterization of
180 high-dimensional functions for the curse of memory of EUAF networks. Another approach
181 to reducing the memory requirement is to increase the network size. Our main result has
182 provided a network size $\mathcal{O}(d^2)$ to achieve an arbitrary error. If a larger network size is used,
183 the curse of memory can be lessened as we shall discuss in Section 1.4.

184 1.2 Related Work

185 In recent years, there has been an increasing amount of literature on the approximation
186 power of neural networks as a special case of nonlinear approximation (DeVore, 1998; Cohen
187 et al., 2022; Daubechies et al., 2022). In the early works of approximation theory for neural
188 networks, the universal approximation theorem (Cybenko, 1989; Hornik, 1991; Hornik et al.,
189 1989) without approximation errors showed that there exists a sufficiently large neural
190 network approximating a target function in a certain function space within any given error
191 $\varepsilon > 0$. There are also other versions of the universal approximation theorem. For example,
192 it was shown in (Lin and Jegelka, 2018) that the residual neural networks activated the
193 rectified linear unit (ReLU) with one neuron per hidden layer and a sufficiently large depth
194 are a universal approximator. The universal approximation property for general residual
195 neural networks was proved in (Li et al., to appear) via a dynamical system approach. In
196 all papers discussed above, the network size goes to infinity when the target approximation
197 error approaches 0. However, our result in Theorem 1 implies that EUAF networks with a
198 fixed size ($\mathcal{O}(d^2)$ neurons in total) can achieve an arbitrary small error for approximating
199 $f \in C([a, b]^d)$.

200 The approximation errors in terms of the total number of parameters of ReLU networks
201 are well studied for basic function spaces with (nearly) optimal approximation errors, e.g.,
202 (nearly) optimal asymptotic errors for continuous functions (Yarotsky, 2018), C^s functions
203 (Yarotsky and Zhevnerchuk, 2020), piecewise smooth functions (Petersen and Voigtlaender,
204 2018), solutions of special PDEs (Elbrächter et al., 2022; Beck et al., 2020), functions that
205 can be optimally approximated by affine systems (Bölcskei et al., 2019), and Sobolev spaces
206 (Yang et al., 2022; Hon and Yang, 2021). Approximation errors in terms of width and
207 depth would be more useful than those in terms of the total number of nonzero parameters
208 in practice, because width and depth are two essential hyper-parameters in every numerical

209 algorithm instead of the number of nonzero parameters. This motivated the works on the
 210 (nearly) optimal non-asymptotic errors in terms of width and depth with explicit pre-factors
 211 for approximating continuous functions in (Shen et al., 2020, 2022; Zhang, 2020) and for
 212 C^s functions in (Lu et al., 2021; Zhang, 2020). As the errors are (nearly) optimal, there are
 213 two possible directions to improve the approximation error in order to reduce the effect of
 214 the curse of dimensionality. The first one is to consider smaller target function spaces, e.g.,
 215 analytic functions (E and Wang, 2018; Bonito et al., 2021), Barron spaces (Barron, 1993; E
 216 et al., 2019b; E and Wojtowytsch, 2022; Siegel and Xu, 2021), and band-limited functions
 217 (Chen and Wu, 2019; Montanelli et al., 2021).

218 Another direction is to design advanced activation functions, where one can use mul-
 219 tiple activation functions, to enhance the power of neural networks, especially to conquer
 220 the curse of dimensionality in network approximation. There have been several papers de-
 221 signing activation functions to achieve good approximation errors. The results in (Yarotsky
 222 and Zhevnerchuk, 2020) imply that (sin, ReLU)-activated neural networks (i.e., the acti-
 223 vation function of a neuron can be chosen from either sin or ReLU) with W parameters
 224 can approximate Lipschitz continuous functions with an asymptotic approximation error
 225 $\mathcal{O}(e^{-c_d\sqrt{W}})$, where c_d is a constant depending on d . In (Shen et al., 2021a), it was shown
 226 that (Floor, ReLU)-activated neural networks with width $\mathcal{O}(N)$ and depth $\mathcal{O}(L)$ admit an
 227 quantitative approximation error $\mathcal{O}(\sqrt{d}N^{-\sqrt{L}})$ for Lipschitz continuous functions, conquer-
 228 ing the curse of dimensionality in approximation with a root-exponentially small error in
 229 depth L .¹ In (Shen et al., 2021b), it was shown that, even if the depth is as small as 3, neu-
 230 ral networks with width N and $\mathcal{O}(d + N)$ nonzero parameters can approximate Lipschitz
 231 continuous functions with an exponentially small error $\mathcal{O}(\sqrt{d}2^{-N})$, if the floor function
 232 $\lfloor x \rfloor$, the exponential function 2^x , and the step function $\mathbf{1}_{\{x \geq 0\}}$ are used as activation func-
 233 tions. Recently in (Jiao et al., 2021), the results in (Yarotsky and Zhevnerchuk, 2020; Shen
 234 et al., 2021b) were combined to avoid the curse of dimensionality using ReLU, sin, and 2^x
 235 activation functions. Corollary 2 implies that the hypothesis space of EUAF networks acti-
 236 vated by a single activation function with $\mathcal{O}(d^2)$ neurons is dense in $C([a, b]^d)$. Particularly,
 237 all continuous functions can be arbitrarily approximated by fixed-size EUAF networks with
 238 width N and depth L on a d -dimensional hypercube whenever $N \geq 36d(2d + 1)$ and $L \geq 11$.

239 There is another research line for the approximation error of neural networks: applying
 240 KST (Kolmogorov, 1957) or its variants to explore new activation functions for a fixed-
 241 size network to achieve an arbitrary error. The original KST shows that any multivariate
 242 function $f \in C([0, 1]^d)$ can be represented as $f(\mathbf{x}) = \sum_{i=0}^{2d} g_i(\sum_{j=1}^d h_{i,j}(x_j))$ for any $\mathbf{x} =$
 243 $[x_1, x_2, \dots, x_d]^T \in [0, 1]^d$, where g_i and $h_{i,j}$ are univariate continuous functions. In fact,
 244 the composition architecture of KST can be regarded as a special neural network with
 245 (complicated) activation functions depending on the target function, which results in the
 246 failure of KST in practice. To alleviate this issue, a single activation function independent
 247 of the target function is designed in (Maiorov and Pinkus, 1999) to construct networks
 248 with a fixed size ($\mathcal{O}(d)$ neurons) to achieve an arbitrary error for approximating functions
 249 in $C([-1, 1]^d)$. However, the activation function in (Maiorov and Pinkus, 1999) has no

1. Although there is no curse of dimensionality in network approximation, the construction requires exponentially many data samples of the target function and computer memory. Hence, there would be a curse of dimensionality in inferring a target function from its finite samples when standard learning techniques are applied to a computer.

250 closed form and is hardly computable. See Section 2.2 for a detailed discussion of the
 251 construction in (Maierov and Pinkus, 1999). The computability issue of activation functions
 252 was addressed recently in (Yarotsky, 2021). It was shown in (Yarotsky, 2021) that, for an
 253 arbitrary $\varepsilon > 0$ and any function f in $C([0, 1]^d)$, there exists a network of size only depending
 254 on d constructed with multiple activation functions either (sin & arcsin) or ($[\cdot]$ & a non-
 255 polynomial analytic function) to approximate f within an error ε . To the best of our
 256 knowledge, there is no explicit characterization of the size dependence on d in (Yarotsky,
 257 2021). For example, a very important question is whether the dependence can be mild,
 258 e.g., only a polynomial of d , or has to be severe, e.g., exponentially in d . The results of
 259 the current paper provide positive answers to all the issues discussed above: We show that
 260 EUAF networks with a simple and computable activation function, width $36d(2d + 1)$, and
 261 depth 11 can approximate functions in $C([a, b]^d)$ within an arbitrary pre-specified error
 262 $\varepsilon > 0$.

263 In summary, this paper aims to design a simple and computable activation function
 264 σ to construct fixed-size neural networks with the universal approximation property. The
 265 network width and depth are explicitly characterized, depending only on the dimension
 266 d . The fixed-size neural network is designed to approximate any continuous functions on a
 267 hypercube within an arbitrary error by only adjusting $\mathcal{O}(d^2)$ network parameters. Moreover,
 268 we prove that an arbitrary classification function can be exactly represented by such a
 269 fixed-size network architecture via only adjusting $\mathcal{O}(d^2)$ network parameters. The main
 270 contribution of this paper is to develop a rigorous mathematical analysis for the universal
 271 approximation property of fixed-size neural networks. The mathematical analysis developed
 272 here would provide a deeper understanding for other neural networks and the approximation
 273 results discussed here can be applied to the full error analysis of deep learning in the next
 274 subsection.

275 1.3 Error Analysis

276 We will briefly discuss the full error analysis of deep neural networks. Let $\Phi(\mathbf{x}, \boldsymbol{\theta})$ denote a
 277 function of $\mathbf{x} \in \mathcal{X}$ generated by a network architecture parameterized with $\boldsymbol{\theta} \in \mathbb{R}^W$. Given
 278 a target function f defined on \mathcal{X} , the final goal is to find the expected risk minimizer

$$279 \quad \boldsymbol{\theta}_{\mathcal{D}} \in \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^W} R_{\mathcal{D}}(\boldsymbol{\theta}), \quad \text{where } R_{\mathcal{D}}(\boldsymbol{\theta}) := \mathbb{E}_{\mathbf{x} \sim U(\mathcal{X})} [\ell(\Phi(\mathbf{x}, \boldsymbol{\theta}), f(\mathbf{x}))]$$

280 with an unknown data distribution $U(\mathcal{X})$ over \mathcal{X} and a loss function $\ell(\cdot, \cdot)$ typically taken
 281 as $\ell(y_1, y_2) = \frac{1}{2}|y_1 - y_2|^2$. Note that $\boldsymbol{\theta}_{\mathcal{D}}$ may not be always achievable. For any pre-specified
 282 $\eta > 0$, one can always identify $\boldsymbol{\theta}_{\mathcal{D}, \eta} \in \mathbb{R}^W$ instead of $\boldsymbol{\theta}_{\mathcal{D}}$ such that

$$283 \quad R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}, \eta}) \leq \inf_{\boldsymbol{\theta} \in \mathbb{R}^W} R_{\mathcal{D}}(\boldsymbol{\theta}) + \eta/2. \quad (4)$$

284 Since the expected risk $R_{\mathcal{D}}(\boldsymbol{\theta})$ is not available in practice, we use the empirical risk $R_{\mathcal{S}}(\boldsymbol{\theta})$
 285 to approximate $R_{\mathcal{D}}(\boldsymbol{\theta})$ for given samples $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n$ and our goal is to identify the
 286 empirical risk minimizer

$$287 \quad \boldsymbol{\theta}_{\mathcal{S}} \in \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^W} R_{\mathcal{S}}(\boldsymbol{\theta}), \quad \text{where } R_{\mathcal{S}}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \ell(\Phi(\mathbf{x}_i, \boldsymbol{\theta}), f(\mathbf{x}_i)).$$

288 Similarly, $\boldsymbol{\theta}_S$ is not always achievable. For any pre-specified $\eta > 0$, one can always identify
 289 $\boldsymbol{\theta}_{S,\eta} \in \mathbb{R}^W$ instead of $\boldsymbol{\theta}_S$ such that

$$290 \quad R_S(\boldsymbol{\theta}_{S,\eta}) \leq \inf_{\boldsymbol{\theta} \in \mathbb{R}^W} R_S(\boldsymbol{\theta}) + \eta/2. \quad (5)$$

291 In practical implementation, only a numerical minimizer $\boldsymbol{\theta}_N$ of $R_S(\boldsymbol{\theta})$ can be achieved via
 292 a numerical optimization method. The discrepancy between the learned function $\Phi(\mathbf{x}, \boldsymbol{\theta}_N)$
 293 and the target function f is measured by $R_{\mathcal{D}}(\boldsymbol{\theta}_N)$, which is bounded by

$$294 \quad \begin{aligned} R_{\mathcal{D}}(\boldsymbol{\theta}_N) &= \underbrace{[R_{\mathcal{D}}(\boldsymbol{\theta}_N) - R_S(\boldsymbol{\theta}_N)]}_{\text{GE}} + \underbrace{[R_S(\boldsymbol{\theta}_N) - R_S(\boldsymbol{\theta}_{S,\eta})]}_{\text{OE}} + \underbrace{[R_S(\boldsymbol{\theta}_{S,\eta}) - R_S(\boldsymbol{\theta}_{\mathcal{D},\eta})]}_{\leq \eta/2 \text{ by (5)}} + \underbrace{[R_S(\boldsymbol{\theta}_{\mathcal{D},\eta}) - R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D},\eta})]}_{\text{GE}} + \underbrace{R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D},\eta})}_{\leq \inf_{\boldsymbol{\theta} \in \mathbb{R}^W} R_{\mathcal{D}}(\boldsymbol{\theta}) + \eta/2 \text{ by (4)}} \\ &\leq \underbrace{\eta}_{\text{perturbation}} + \underbrace{\inf_{\boldsymbol{\theta} \in \mathbb{R}^W} R_{\mathcal{D}}(\boldsymbol{\theta})}_{\text{approximation error}} + \underbrace{[R_S(\boldsymbol{\theta}_N) - \inf_{\boldsymbol{\theta} \in \mathbb{R}^W} R_S(\boldsymbol{\theta})]}_{\text{optimization error (OE)}} + \underbrace{[R_{\mathcal{D}}(\boldsymbol{\theta}_N) - R_S(\boldsymbol{\theta}_N)] + [R_S(\boldsymbol{\theta}_{\mathcal{D},\eta}) - R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D},\eta})]}_{\text{generalization error (GE)}}. \end{aligned}$$

295 If $\Phi(\mathbf{x}, \boldsymbol{\theta})$ is realized by EUAF networks, then Theorem 1 implies

$$296 \quad \inf_{\boldsymbol{\theta} \in \mathbb{R}^W} \|\Phi(\cdot, \boldsymbol{\theta}) - f(\cdot)\|_{L^\infty(\mathcal{X})} = 0 \quad \text{for all } f \in C(\mathcal{X}) \text{ with } \mathcal{X} = [a, b]^d.$$

297 It follows that

$$298 \quad \inf_{\boldsymbol{\theta} \in \mathbb{R}^W} R_{\mathcal{D}}(\boldsymbol{\theta}) = \inf_{\boldsymbol{\theta} \in \mathbb{R}^W} \mathbb{E}_{\mathbf{x} \sim U(\mathcal{X})} [\ell(\Phi(\mathbf{x}, \boldsymbol{\theta}), f(\mathbf{x}))] = 0.$$

299 Since the pre-specified hyper-parameter η can be arbitrarily small, the full error analysis
 300 can be reduced to the analysis of the optimization and generalization errors, which de-
 301 pends on data samples, optimization algorithms, etc. One could refer to (Neyshabur et al.,
 302 2019; E et al., 2019a,b; E and Wojtowysch, 2020; Kawaguchi, 2016; Nguyen and Hein,
 303 2017; Kawaguchi and Bengio, 2019; He et al., 2020; Li et al., 2019) for the analysis of the
 304 generalization and optimization errors.

305 1.4 Computability

306 The EUAF network is simple and computable in the sense that the output and subgradient
 307 of EUAF networks can be efficiently evaluated. The computability of EUAF implies that
 308 we can numerically implement the optimization algorithm to find a numerical minimizer
 309 of the empirical risk. Therefore, EUAF can be directly applied to existing deep learning
 310 software in the same way as other popular activation functions (such as ReLU or Sigmoid).
 311 For further discussion on the computability of EUAF, one may refer to Section 3, which
 312 provides experiments to explore the numerical properties of EUAF. As opposed to the
 313 computability of EUAF, the activation function proposed in (Maiorov and Pinkus, 1999)
 314 is not computable in the sense that there is no numerical algorithm to evaluate the output
 315 and subgradient of the corresponding network.

316 As we shall see later in the proof of Theorem 1, our EUAF network may require suffi-
 317 ciently large parameters to achieve an arbitrarily small error. The magnitude of network
 318 parameters in Theorem 1 can be dramatically reduced by increasing the network size. In
 319 particular, if we replace each elemental block like Figure 2(a) by a block like Figure 2(b),
 320 then the magnitude of parameters can be roughly reduced to its square root. By repeatedly
 321 applying this idea, it is easy to prove that the magnitude of parameters can be exponentially

322 reduced as the network size increases linearly. If we fix the size of these larger networks and
 323 only tune their parameters, they can still approximate high-dimensional continuous func-
 324 tions within an arbitrarily small error. How to fix a network size to balance the number
 325 of parameters and their memory depends on both the computer hardware and software.
 326 The goal of this paper is to demonstrate the existence of a simple network with a fixed size
 327 achieving an arbitrary error in spite of the magnitude of parameters and we have shown
 328 that the network size can be as small as $\mathcal{O}(d^2)$. It is interesting to investigate the balance
 329 between the network size and the memory requirement in the future.

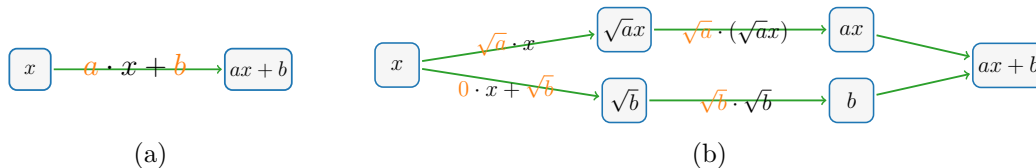


Figure 2: Illustrations of the magnitude reduction of parameters for a sub-network. The parameters are marked in orange. Without loss of generality, $a \gg 1$ and $b \gg 1$. (a) Return $ax + b$ via two large parameters a and b . (b) Return $ax + b$ via several small parameters bounded by $\max\{\sqrt{a}, \sqrt{b}\}$.

330 In real-world applications, the parameters of the EUAF network are learned from the
 331 samples of the target function, which involves sophisticated numerical optimization. We
 332 refer to the learnability of network parameters as the existence of a numerical optimization
 333 algorithm that can identify network parameters to achieve a target approximation error.
 334 The computability of the EUAF networks does not imply learnability, which involves ap-
 335 proximation, optimization, and generalization error analyses. The result in this paper shows
 336 that there exist computable EUAF networks achieving an arbitrarily small approximation
 337 error. This means the learnability of the best approximation is reduced to achieving small
 338 generalization and optimization errors, which depend on the given data, the empirical risk
 339 model, and the optimization algorithm. Therefore, whether or not EUAF networks would
 340 be useful in real-world applications also depends on optimization and generalization, which
 341 is out of the scope of this paper. The optimization and generalization error analyses of
 342 practical deep neural networks including EUAF networks is a challenging problem. To the
 343 best of our knowledge, there is no complete error analysis to address the learnability of
 344 neural networks with nonlinear activation functions.

345 The rest of this paper is organized as follows. In Section 2, we first summarize notations
 346 used in this paper and then discuss the ideas of proving Theorem 1. Section 3 focuses
 347 on numerical experimentation of EUAF, which acts as a proof of concept to explore the
 348 numerical properties of EUAF. Next, several UAFs with better properties are proposed in
 349 Section 4. After that, we use several sections to present the complete proofs of Theorems 1
 350 and 4. In Section 5, by assuming Theorem 6 is true, we give the detailed proofs of Theo-
 351 rems 1 and 4. Theorem 6 is proved in Section 6 based on Proposition 7, the proof of which
 352 can be found in Section 7. Finally, Section 8 concludes this paper with a short discussion.

353 2. Notations and Proof Ideas

354 In this section, we first summarize notations used in this paper and then discuss the ideas
 355 of proving Theorem 1.

356 2.1 Notations

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

358 • Let \mathbb{R} , \mathbb{Q} , and \mathbb{Z} denote the set of real numbers, rational numbers, and integers,
 359 respectively.

360 • Let \mathbb{N} and \mathbb{N}^+ denote the set of natural numbers and positive natural numbers, re-
 361 spectively. That is, $\mathbb{N}^+ = \{1, 2, 3, \dots\}$ and $\mathbb{N} = \mathbb{N}^+ \cup \{0\}$.

362 • For any $x \in \mathbb{R}$, let $\lfloor x \rfloor := \max\{n : n \leq x, n \in \mathbb{Z}\}$ and $\lceil x \rceil := \min\{n : n \geq x, n \in \mathbb{Z}\}$.

363 • Let $\mathbf{1}_S$ be the indicator (characteristic) function of a set S , i.e., $\mathbf{1}_S$ is equal to 1 on S
 364 and 0 outside S .

365 • The set difference of two sets A and B is denoted by $A \setminus B := \{x : x \in A, x \notin B\}$.

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

368 bold lowercase letters. For example, $\mathbf{v} = [v_1, v_2, \dots, v_d]^T = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_d \end{bmatrix} \in \mathbb{R}^d$ is a column

369 vector. Besides, “[” and “]” are used to partition matrices (vectors) into blocks, e.g.,

370 $\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}$.

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

373
$$\|\mathbf{x}\|_p = \|\mathbf{x}\|_{\ell^p} := (|x_1|^p + |x_2|^p + \dots + |x_d|^p)^{1/p}.$$

374 In the case $p = \infty$,

375
$$\|\mathbf{x}\|_\infty = \|\mathbf{x}\|_{\ell^\infty} := \max\{|x_i| : i = 1, 2, \dots, d\}.$$

376 • For any $a_1, a_2, \dots, a_J \in \mathbb{R}$, we say a_1, a_2, \dots, a_J are **rationally independent** if
 377 they are linearly independent over the rational numbers \mathbb{Q} . That is, if there exist
 378 $\lambda_1, \lambda_2, \dots, \lambda_J \in \mathbb{Q}$ such that $\sum_{j=1}^J \lambda_j \cdot a_j = 0$, then $\lambda_1 = \lambda_2 = \dots = \lambda_J = 0$. For a
 379 simple example, $1, \sqrt{2}$, and $\sqrt{3}$ are rationally independent.

380 • An **algebraic** number is any complex number (including real numbers) that is a root
 381 of a polynomial equation with rational coefficients, i.e., α is an algebraic number if
 382 and only if there exist $\lambda_0, \lambda_1, \dots, \lambda_J \in \mathbb{Q}$ with $\sum_{j=0}^J \lambda_j \alpha^j = 0$.² Denote the set of all
 383 algebraic numbers by \mathbb{A} . We say a complex number is **transcendental** if it is not

2. For simplicity, we denote $1 = x^0$ for any $x \in \mathbb{R}$, including the case 0^0 .

384 in \mathbb{A} . The set \mathbb{A} is countable, and, therefore, almost all numbers are transcendental.
 385 The best known transcendental numbers are π (the ratio of a circle’s circumference
 386 to its diameter) and e (the natural logarithmic base).

- 387 • The expression “a network (architecture) with width N and depth L ” means
 - 388 – The number of neurons in each **hidden** layer of this network (architecture) is no
 389 more than N .
 - 390 – The number of **hidden** layers of this network (architecture) is no more than L .

391 2.2 Key Ideas of Proving Theorem 1

392 The proof of Theorem 1 has two main steps: 1) prove the one-dimensional case; 2) reduce
 393 the d -dimensional approximation to the one-dimensional case via KST (Kolmogorov, 1957).
 394 In fact, in the case of $d = 1$, the size of the network in Theorem 1 can be further reduced as
 395 shown in Theorem 6 below. Theorem 6 is actually an enhanced version of Theorem 1 and
 396 hence implies Theorem 1 in the case $d = 1$.

397 **Theorem 6.** *Let $f \in C([a, b])$ be a continuous function. Then, for an arbitrary $\varepsilon > 0$,
 398 there exists a function ϕ generated by an EUAF network with width 36 and depth 5 such
 399 that*

$$400 \quad |\phi(x) - f(x)| < \varepsilon \quad \text{for any } x \in [a, b] \subseteq \mathbb{R}.$$

401 The detailed proof of Theorem 6 can be found in Section 6. The main ideas of proving
 402 Theorem 6 are developed from some ideas of our early works (Shen et al., 2021a,b). Roughly
 403 speaking, we eventually convert a function approximation problem in an interval (e.g.,
 404 $[0, 1)$) to a point-fitting problem via the composition architecture of neural networks in the
 405 following three main steps.³

- 406 • Divide $[0, 1)$ into small intervals $\mathcal{I}_k = [\frac{k-1}{K}, \frac{k}{K})$ with a left endpoint x_k for $k \in$
 407 $\{1, 2, \dots, K\}$, where K is an integer determined by the given error and the target
 408 function f .
- 409 • Construct a sub-network to generate a function ϕ_1 mapping the whole interval \mathcal{I}_k to k
 410 for each k . The floor function $\lfloor \cdot \rfloor$ is a good choice to implement this step. Precisely, we
 411 can define $\phi_1(x) = \lfloor Kx \rfloor$. The floor function is not continuous and has zero-derivative
 412 almost everywhere. As we shall see later, σ_1 (or σ) can be a continuous alternative to
 413 implement this step, but the construction is more complicated.
- 414 • The final step is to design another sub-network to generate a function ϕ_2 mapping k
 415 approximately to $f(x_k)$ for each k . Then $\phi_2 \circ \phi_1(x) = \phi_2(k) \approx f(x_k) \approx f(x)$ for any
 416 $x \in \mathcal{I}_k$ and $k \in \{1, 2, \dots, K\}$, which implies $\phi_2 \circ \phi_1 \approx f$ on $[0, 1)$. After the above two
 417 steps, we simplify the approximation problem to a point-fitting problem, where k is
 418 approximately mapped to $f(k)$. This step is the bottleneck of the construction in our
 419 previous papers (Shen et al., 2021a,b). Roughly speaking, the final approximation
 420 error is essentially determined by how many points we can fit using a neural network.

3. The goal of a point-fitting problem is to identify a function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ in a given hypothesis space (e.g., the space of functions realized by neural networks) such that $|\phi(\mathbf{x}_i) - y_i| < \varepsilon$ for $i = 1, 2, \dots, n$ and a pre-specified error $\varepsilon > 0$, where $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \subseteq \mathbb{R}^{d+1}$ are given samples.

421 For the second step, the capacity to generate step functions with sufficiently many
 422 “steps” via a sub-network with a limited number of neurons plays an important role. The
 423 reproduced step functions can be considered as a continuous version of the floor function
 424 ($\lfloor \cdot \rfloor$) in (Shen et al., 2021a,b), which is a perfect step function with infinite “steps” that
 425 improves the approximation power of networks as shown in (Shen et al., 2021a,b). The key
 426 ingredient in the third step of the proof of Theorem 6 is essentially a point-fitting problem
 427 with arbitrarily many points. This requires the following proposition motivated by the well-
 428 known fact that an irrational winding on the torus is dense. See Figure 3 for illustrations
 429 of such a fact. Here, we propose a new point-fitting technique that can fit arbitrarily many
 430 points within an arbitrary error using fixed-size neural networks.

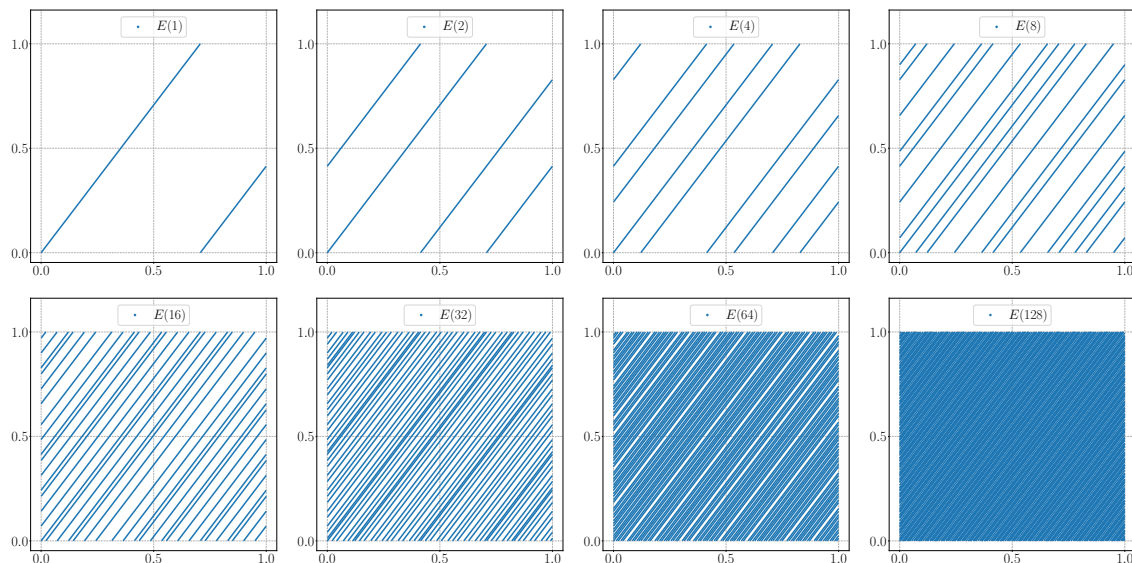


Figure 3: Illustrations of the denseness of $E(\infty)$ in $[0, 1]^2$, where $E(r)$ is a winding of an “irrational” direction $[1, \sqrt{2}]^T$ on $[0, r)$, i.e., $E(r) = \{[\tau(t), \tau(\sqrt{2}t)]^T : t \in [0, r)\}$ with $\tau(t) = t - \lfloor t \rfloor$.

431 **Proposition 7.** *For any $K \in \mathbb{N}^+$, the following point set*

432
$$\left\{ \left[\sigma_1\left(\frac{w}{\pi+1}\right), \sigma_1\left(\frac{w}{\pi+2}\right), \dots, \sigma_1\left(\frac{w}{\pi+K}\right) \right]^T : w \in \mathbb{R} \right\} \subseteq [0, 1]^K$$

433 *is dense in $[0, 1]^K$, where π is the ratio of a circle’s circumference to its diameter.*

434 The proof of Proposition 7 can be found in Section 7. To prove the denseness in Propo-
 435 sition 7, we borrow some ideas from transcendental number theory and Diophantine ap-
 436 proximations in number theory. The number π used in Proposition 7 is transcendental. It
 437 can be replaced by any other transcendental number.

438 Proposition 7 implies that for any given sample points $(k, y_k) \in \mathbb{R}^2$ with $y_k \in [0, 1]$ for
 439 $k = 1, 2, \dots, K$ and any $K \in \mathbb{N}^+$, there exists $w_0 \in \mathbb{R}$ such that the function $x \mapsto \sigma_1\left(\frac{w_0}{\pi+x}\right)$
 440 fit the points $(k, y_k) \in \mathbb{R}^2$ for $k = 1, 2, \dots, K$ within an arbitrary pre-specified error $\varepsilon > 0$.
 441 To put it another way, for any $\varepsilon > 0$, there exists $w_0 \in \mathbb{R}$ such that $|\sigma_1\left(\frac{w_0}{\pi+k}\right) - y_k| < \varepsilon$ for
 442 all k .

443 As we shall see later in the proof of Proposition 7, the key point is the periodicity of the
 444 outer function σ_1 . Of course, the inner function $x \mapsto \frac{w_0}{\pi+x}$ is also necessary since it helps
 445 to adjust sample points for $x = 1, 2, \dots, K$. In fact, the inner function $x \mapsto \frac{w_0}{\pi+x}$ can be
 446 regarded as a variant of σ_2 via scaling and shifting. The periodicity has been explored to
 447 improve neural network approximation in the literature, e.g. the sine function in (Yarotsky
 448 and Zhevnerchuk, 2020) is periodic and the floor function ($\lfloor \cdot \rfloor$) in (Shen et al., 2021a,b) is
 449 implicitly periodic because $x - \lfloor x \rfloor$ is periodic. We remark that a similar result holds if we
 450 replace σ_1 by a non-trivial periodic function and replace the sample locations $x = 1, 2, \dots, K$
 451 by distinct rational numbers $r_1, r_2, \dots, r_K \in \mathbb{Q}$. See Section 7 for a further discussion.

452 Theorem 6 essentially proves Theorem 1 for the univariate case. To prove the general
 453 case, we need the Kolmogorov superposition theorem (KST) (Kolmogorov, 1957) given
 454 below to reduce a multivariate problem to a one-dimensional case.

455 **Theorem 8 (KST).** *There exist continuous functions $h_{i,j} \in C([0, 1])$ for $i = 0, 1, \dots, 2d$
 456 and $j = 1, 2, \dots, d$ such that any continuous function $f \in C([0, 1]^d)$ can be represented as*

$$457 \quad f(\mathbf{x}) = \sum_{i=0}^{2d} g_i \left(\sum_{j=1}^d h_{i,j}(x_j) \right) \quad \text{for any } \mathbf{x} = [x_1, x_2, \dots, x_d]^T \in [0, 1]^d,$$

458 where $g_i : \mathbb{R} \rightarrow \mathbb{R}$ is a continuous function for each $i \in \{0, 1, \dots, 2d\}$.

459 KST is often used to reduce a multidimensional problem to a one-dimensional one. In
 460 fact, the compositional representation in KST can be regarded as a special neural network
 461 with (complicated) activation functions depending on the target function, which makes
 462 KST useless in practical computation. To avoid this dependency, an activation function
 463 was designed in (Maierov and Pinkus, 1999) to construct neural network representations
 464 with $\mathcal{O}(d)$ neurons that can approximate functions in $C([-1, 1]^d)$ within an arbitrary error.
 465 Let us briefly summarize the main ideas in (Maierov and Pinkus, 1999): 1) Identify a dense
 466 and countable subset $\{u_k\}_{k=1}^\infty$ of $C([-1, 1])$, e.g., polynomials with rational coefficients. 2)
 467 Construct an activation function ϱ to encode all $u_k(x)$ for $x \in [-1, 1]$. In fact, for each
 468 k , $u_k|_{[-1, 1]}$ is “stored” in ϱ on $[4k, 4k + 2]$, and the values of ϱ on $[4k + 2, 4k + 4]$ are
 469 properly assigned to make ϱ a smooth and monotonically increasing function. That is, let
 470 $\varrho(x + 4k + 1) = a_k + b_k x + c_k u_k(x)$ for any $x \in [-1, 1]$ with carefully chosen constants
 471 a_k , b_k , and $c_k \neq 0$ such that $\varrho(x)$ can be a sigmoidal function. 3) For any $g \in C([-1, 1])$,
 472 there exists a one-hidden-layer ϱ -activated network with width 3 approximating g within
 473 an arbitrary error $\delta > 0$, i.e., there exists k such that $g(x) \stackrel{\delta}{\approx} u_k(x) = \frac{\varrho(x+4k+1) - a_k - b_k x}{c_k}$
 474 for any $x \in [-1, 1]$. 4) Replace the inner and outer functions in KST with these one-
 475 hidden-layer networks to achieve a two-hidden-layer ϱ -activated network with width $\mathcal{O}(d)$
 476 to approximate $f \in C([-1, 1]^d)$ within an arbitrary error $\varepsilon > 0$. As we can see, the key
 477 point of the construction in (Maierov and Pinkus, 1999) is to encode a dense and countable
 478 subset of the target function space in an activation function.

479 Note that both (Maierov and Pinkus, 1999) and this paper use KST to reduce dimension.
 480 However, the activation function of (Maierov and Pinkus, 1999) is complicated without any
 481 closed form and there is no efficient numerical algorithm to evaluate it. After encoding
 482 a dense subset of continuous function into a single but complicated activation function,
 483 one only needs to construct affine linear transformations to select appropriate functions

484 of this dense subset from this complicated activation function to construct approximation.
 485 Hence, such a complicated activation function simplifies the proof of the denseness, since
 486 the denseness is encoded in the activation function. As a contrast, we design a simple
 487 activation function with efficient numerical implementation (see Figure 1 for an illustration)
 488 achieving the universal approximation property with fixed-size networks, because simple and
 489 implementable activation functions are a basic requirement for a neural network to be used
 490 in applications. However, the proof of the denseness of a neural network generated by such
 491 a simple activation function becomes difficult. A sophisticated analysis will be developed
 492 in the rest of this paper to overcome the difficulties.

493 3. Experimentation

494 In this section, we will conduct two simple experiments as a proof of concept to explore
 495 the numerical performances of the EUAF activation function. Let us first discuss the
 496 numerical implementation of EUAF in `PyTorch`. To enable the automatic differentiation
 497 feature for EUAF, we need to implement EUAF based on PyTorch built-in functions. With
 498 the following four built-in functions $\text{abs}(x) = |x|$, $\text{floor}(x) = \lfloor x \rfloor$,

$$499 \quad \text{softsign}(x) = \frac{x}{|x| + 1}, \quad \text{and} \quad \text{sign}(x) = \begin{cases} 1 & \text{if } x > 0, \\ 0 & \text{if } x = 0, \\ -1 & \text{if } x < 0, \end{cases}$$

500 we can represent EUAF as

$$\begin{aligned} \text{EUAF}(x) &= \begin{cases} \text{softsign}(x) & \text{if } x < 0, \\ |x - 2\lfloor \frac{x+1}{2} \rfloor| & \text{if } x \geq 0 \end{cases} \\ 501 \quad &= \text{softsign}(x) \cdot \frac{1 - \text{sign}(x)}{2} + \left| x - 2\lfloor \frac{x+1}{2} \rfloor \right| \cdot \frac{1 + \text{sign}(x)}{2} \\ &= \text{softsign}(x) \cdot \frac{1 - \text{sign}(x)}{2} + \text{abs}\left(x - 2 \cdot \text{floor}\left(\frac{x+1}{2}\right)\right) \cdot \frac{1 + \text{sign}(x)}{2}. \end{aligned}$$

502 Thus, it is numerically cheap to compute EUAF and its subgradient. We believe the EUAF
 503 activation function can achieve good results in some real-world applications if proper op-
 504 timization algorithms are developed for EUAF. In this paper, we only conduct two simple
 505 experiments: a function approximation experiment in Section 3.1 and a classification ex-
 506 periment in Section 3.2.

507 Next, let us briefly discuss when our EUAF activation function would outperform the
 508 practically used ones (e.g., ReLU, Sigmoid, and Softsign), which is based on full error
 509 analysis in Section 1.3. In our discussion, we take the ReLU activation function as an
 510 example and suppose the optimization error is well-controlled. Clearly, replacing ReLU by
 511 EUAF can reduce the approximation error, but would result in a large generalization error.
 512 Thus, we would expect that EUAF achieves better results than ReLU if the approximation
 513 error is larger than the generalization error. That means EUAF would outperform ReLU
 514 in the following two cases.

- 515 • The approximation error is pretty large (e.g., the target function is sufficiently com-
 516 plicated).

517 • The generalization error is well-controlled (e.g., there are sufficiently many samples).

518 If a given problem does not belong to these two cases, one may consider replacing only
 519 a small number of ReLUs by EUAFs. In the function approximation experiment in Sec-
 520 tion 3.1, we first choose a complicated target function and then generate sufficiently many
 521 samples to reduce the generalization error. In the classification experiment in Section 3.2,
 522 we control the generalization error via three common methods: keeping network parameters
 523 small via L2 regularization, dropout (Hinton et al., 2012; Srivastava et al., 2014), and batch
 524 normalization (Ioffe and Szegedy, 2015).

525 3.1 Function Approximation

526 We will design fully connected neural network (FCNN) architectures activated by ReLU
 527 or EUAF to solve a function approximation problem. To better compare the approxima-
 528 tion power of ReLU and EUAF activation functions, we choose a complicated (oscillatory)
 529 function f as the target function, where f is defined as

$$530 \quad f(x_1, x_2) := 0.6 \sin(8x_1) + 0.4 \sin(16x_2) \quad \text{for any } (x_1, x_2) \in [0, 1]^2.$$

531 To compare the numerical performances of ReLU and EUAF activation functions, we
 532 design two FCNN architectures with different activation functions. Both of them have 4
 533 hidden layers and each hidden layer has 80 neurons. For simplicity, we denote them as
 534 FCNN1 and FCNN2. See illustrations of them in Figure 4. FCNN1 is a standard fully
 535 connected ReLU network and FCNN2 can be regarded as a variant of FCNN1 by replacing
 536 ReLU by EUAF.

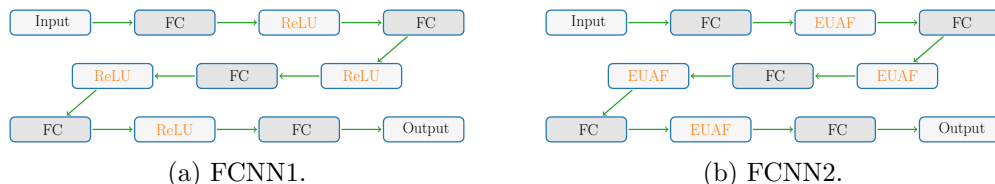


Figure 4: Illustrations of FCNN1 and FCNN2. FC represents a fully connected layer.

537 Before presenting the numerical results, let us present the hyper-parameters for training
 538 FCNN1 and FCNN2. We randomly choose 10^6 training samples and 10^5 test samples in
 539 $[0, 1]^2$. The number of epochs and the batch size are set to 500 and 256, respectively.
 540 We adopt RAdam (Liu et al., 2020) as the optimization method and the learning rate is
 541 $0.002 \times 0.9^{i-1}$ in epochs $5(i-1) + 1$ to $5i$ for $i = 1, 2, \dots, 100$. Several loss functions are
 542 used to estimate the training and test losses, including the mean squared error (MSE), the
 543 mean absolute error (MAE), and the maximum (MAX) loss functions. To illustrate MSE,
 544 MAE and MAX losses, we denote ϕ as the network-generated function and $\mathbf{x}_1, \dots, \mathbf{x}_m$ as
 545 the test samples ($m = 10^5$ in our setting). Then, the MSE loss is given by $\frac{1}{m} \sum_{i=1}^m (\phi(x_i) -$
 546 $f(x_i))^2$, the MAE loss is given by $\frac{1}{m} \sum_{i=1}^m |\phi(x_i) - f(x_i)|$, and the MAX loss is given by
 547 $\max \{|\phi(x_i) - f(x_i)| : i = 1, 2, \dots, m\}$. The MSE loss is used in our training process. In
 548 the settings above, we repeat the experiment 12 times and discard 2 top-performing and 2
 549 bottom-performing trials by using the average of test losses (MSE) in the last 100 epochs

550 as the performance criterion. For each epoch, we adopt the average of training (test) losses
 551 in the rest 8 trials as the target training (test) loss.

552 Next, let us present the experiment results to compare the numerical performances of
 553 ReLU and EUAF activation functions. Training and test losses (MSE) over epochs for
 554 FCNN1 and FCNN2 are summarized in Figure 5.

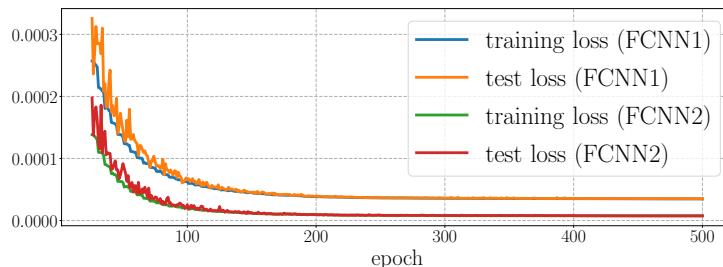


Figure 5: Training and test losses (MSE) in epochs 25-500 for FCNN1 and FCNN2.

555 In Table 1, we present a comparison of FCNN1 and FCNN2 for the average of the test
 556 losses in the last 100 epochs measured in several loss functions. As we can see from Figure 5
 557 and Table 1, FCNN2 performs better than FCNN1. That means replacing ReLU by EUAF
 558 would improve experiment results.

activation function		test loss		
		MSE	MAE	MAX
FCNN1	ReLU	3.53×10^{-5}	4.57×10^{-3}	3.69×10^{-2}
FCNN2	EUAF	7.56×10^{-6}	2.13×10^{-3}	1.48×10^{-2}

Table 1: Test loss comparison.

559 3.2 Classification

560 The goal of a classification problem with $J \in \mathbb{N}^+$ classes is to identify a classification
 561 function f defined by

$$562 \quad f(\mathbf{x}) = j \quad \text{for any } \mathbf{x} \in E_j \text{ and } j = 0, 1, \dots, J - 1,$$

563 where E_0, E_1, \dots, E_{J-1} are pairwise disjoint bounded closed subsets of \mathbb{R}^d and all samples
 564 with a label j are contained in E_j for each j . Such a classification function f can be
 565 continuously extended to \mathbb{R}^d , which means a classification problem can also be regarded as
 566 a continuous function approximation problem. We take the case $J = 2$ as an example to
 567 illustrate the extension. The multiclass case is similar. By defining

$$568 \quad \text{dist}(\mathbf{x}, E_i) := \inf_{\mathbf{y} \in E_i} \|\mathbf{x} - \mathbf{y}\|_2 \quad \text{for any } \mathbf{x} \in \mathbb{R}^d \text{ and } i = 0, 1,$$

569 we have $\text{dist}(\mathbf{x}, E_0) + \text{dist}(\mathbf{x}, E_1) > 0$ for any $\mathbf{x} \in \mathbb{R}^d$. Thus, we can define

$$570 \quad \tilde{f}(\mathbf{x}) := \frac{\text{dist}(\mathbf{x}, E_0)}{\text{dist}(\mathbf{x}, E_0) + \text{dist}(\mathbf{x}, E_1)} \quad \text{for any } \mathbf{x} \in \mathbb{R}^d.$$

571 It is easy to verify that \tilde{f} is continuous on \mathbb{R}^d and

572
$$\tilde{f}(\mathbf{x}) = \begin{cases} 0 & \text{if } \mathbf{x} \in E_0, \\ 1 & \text{if } \mathbf{x} \in E_1 \end{cases} = f(\mathbf{x}) \quad \text{for any } \mathbf{x} \in E_0 \cup E_1.$$

573 That means \tilde{f} is a continuous extension of f . That means we can apply our theory to
574 classification problems.

575 We will design convolutional neural network (CNN) architectures activated by ReLU or
576 EUAF to solve a classification problem corresponding to a standard benchmark data set
577 Fashion-MNIST (Xiao et al., 2017). This data set consists of a training set of 60000 samples
578 and a test set of 10000 samples. Each sample is a 28×28 grayscale image, associated with a
579 label from 10 classes. To compare the numerical performances of ReLU and EUAF activa-
580 tion functions, we design two small CNN architectures with different activation functions.
581 Both of them have two convolutional layers and two fully connected layers. For simplicity,
582 we denote them as CNN1 and CNN2. See illustrations of them in Figure 6. We present
583 more details of CNN1 and CNN2 in Table 2.

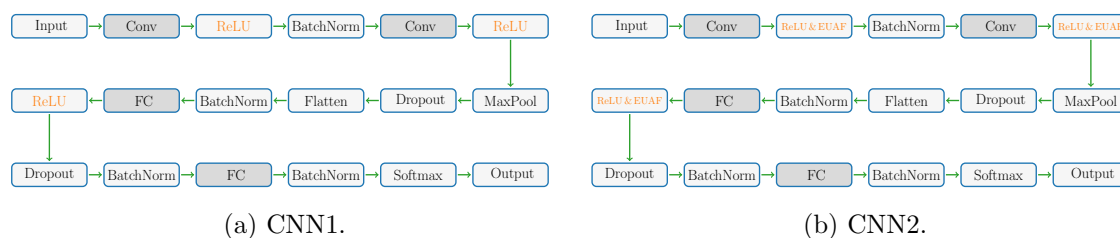


Figure 6: Illustrations of CNN1 and CNN2. Conv and FC represent convolutional and fully connected layers, respectively.

layers	activation function		output size of each layer	dropout	batch normalization
	CNN1	CNN2			
input $\in \mathbb{R}^{28 \times 28}$			28×28		
Conv-1: $1 \times (3 \times 3)$, 24	ReLU	EUAF, $1 \times (26 \times 26)$ ReLU, $23 \times (26 \times 26)$	$24 \times (26 \times 26)$		yes
Conv-2: $24 \times (3 \times 3)$, 24	ReLU	EUAF, $1 \times (24 \times 24)$ ReLU, $23 \times (24 \times 24)$	3456 (MaxPool & Flatten)	0.25	yes
FC-1: 3456, 48	ReLU	EUAF, 1 ReLU, 47	48	0.5	yes
FC-2: 48, 10			10 (Softmax)		yes
output $\in \mathbb{R}^{10}$					

Table 2: Details of CNN1 and CNN2.

584 CNN1 is activated by ReLU, while CNN2 is activated by ReLU and EUAF. In CNN2,
585 only one channel (neuron) of a convolutional (fully connected) hidden layer is activated
586 by EUAF. CNN2 can be regarded as a variant of CNN1 by replacing a small number of
587 ReLUs by EUAFs. This follows a natural question: Why do we not make all (or most)
588 neurons (channels) of CNN2 activated by EUAF? We use only a few EUAFs in CNN2 for
589 two reasons listed below.

- Since the number of available training samples is limited, using too many EUAF activation functions would lead to a large generalization error.
- The key difference of EUAF to the practical used activation functions (e.g., ReLU, Sigmoid, and Softsign) is the periodic part on $[0, \infty)$. As we shall see later in the proof of our main theorem, only a small number of neurons in the constructed network require the periodic property. Thus, we would expect that neural networks activated by the practical used activation functions and a few EUAFs are super expressive.

Next, let us discuss why we choose relatively small network architectures. Since the Fashion-MNIST classification problem is simple, the expressive power of a relatively large ReLU CNN architecture is enough. That means there is no need to introduce EUAF if the network architecture is relatively large. We believe EUAF would be useful for complicated classification problems.

We remark that we use CNNs to approximate an equivalent variant \hat{f} of the original classification function f mentioned previously, where \hat{f} is given by

$$\hat{f}(\mathbf{x}) = \mathbf{e}_j \quad \text{for any } \mathbf{x} \in E_j \text{ and } j = 0, 1, \dots, J - 1,$$

where $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_J\}$ is the standard basis of \mathbb{R}^J , i.e., $\mathbf{e}_j \in \mathbb{R}^J$ denotes the vector with a 1 in the j -th coordinate and 0's elsewhere.

Before presenting the numerical results, let us present the hyper-parameters for training two CNN architectures above. We use the cross-entropy loss function to evaluate the loss. The number of epochs and the batch size are set to 500 and 128, respectively. We adopt RAdam (Liu et al., 2020) as the optimization method. The weight decay of the optimizer is 0.0001 and the learning rate is $0.002 \times 0.9^{i-1}$ in epochs $5(i-1) + 1$ to $5i$ for $i = 1, 2, \dots, 100$. All training (test) samples in the Fashion-MNIST data set are standardized in our experiment, i.e., we rescale all training (test) samples to have a mean of 0 and a standard deviation of 1. In the settings above, we repeat the experiment 48 times and discard 8 top-performing and 8 bottom-performing trials by using the average of test accuracy in the last 100 epochs as the performance criterion. For each epoch, we adopt the average of test accuracies in the rest 32 trials as the target test accuracy.

Let us present the experiment results to compare the numerical performances of CNN1 and CNN2. The test accuracy comparison of CNN1 and CNN2 is summarized in Table 3.

	activation function	largest accuracy	average of largest 100 accuracies	average accuracy in last 100 epochs
CNN1	ReLU	0.933066	0.932852	0.932698
CNN2	ReLU and EUAF	0.933922	0.933685	0.933508

Table 3: Test accuracy comparison.

For each of CNN1 and CNN2, we present the largest test accuracy, the average of largest 100 test accuracies over epochs, and the average of test accuracies in the last 100 epochs. For an intuitive comparison, we also provide illustrations of the test accuracy over epochs for CNN1 and CNN2 in Figure 7. As we can see from Table 3 and Figure 7, CNN2 performs better than CNN1. That means replacing a small number of ReLUs by EUAFs would improve the experiment results.

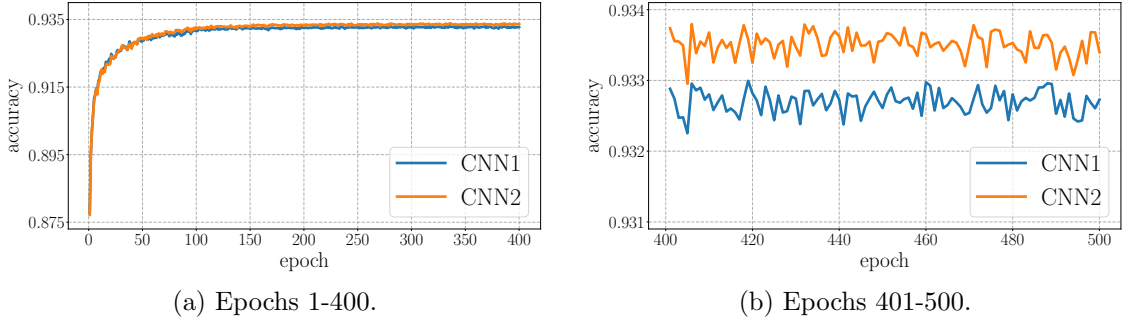


Figure 7: Test accuracy over epochs.

626 **4. Other Examples of UAFs**

627 This section aims at designing new UAFs with additional properties such as smooth or
 628 sigmoidal functions. As discussed in the introduction and shown in the proof of our main
 629 theorem, the construction of UAFs mainly relies on three properties: high nonlinearity,
 630 periodicity, and the capacity to reproduce step functions. The EUAF σ defined in Equa-
 631 tion (1) is a simple and typical example of UAFs satisfying these three properties. Indeed,
 632 having these properties plays an important role in our proof and is a necessary but not
 633 sufficient condition for designing a UAF. In other words, these properties are important,
 634 but cannot guarantee the successful construction of UAFs.

635 Here, we present another idea to design new UAFs, which mainly relies on the following
 636 observation: If a UAF ϱ can be approximated by a fixed-size network activated by a new
 637 activation function $\tilde{\varrho}$ within an arbitrary error on any bounded interval, then $\tilde{\varrho}$ is also a
 638 UAF. Such an observation is a direct result of the lemma below.

639 **Lemma 9.** *Let $\varrho, \tilde{\varrho} : \mathbb{R} \rightarrow \mathbb{R}$ be two functions with $\varrho \in C(\mathbb{R})$. For an arbitrary given
 640 function $f \in [a, b]^d \rightarrow \mathbb{R}$ and any $\varepsilon > 0$, suppose that the following two conditions hold:*

- 641 • *There exists a function ϕ_ϱ realized by a ϱ -activated network with width N and depth
 642 L such that*

$$643 \quad |\phi_\varrho(\mathbf{x}) - f(\mathbf{x})| < \varepsilon/2 \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

- 644 • *For any $M > 0$ and each $\delta \in (0, 1)$, there exists a function ϱ_δ realized by a $\tilde{\varrho}$ -activated
 645 network with width \tilde{N} and depth \tilde{L} such that*

$$646 \quad \varrho_\delta(t) \rightrightarrows \varrho(t) \quad \text{as } \delta \rightarrow 0^+ \quad \text{for any } t \in [-M, M],$$

647 *where \rightrightarrows denotes the uniform convergence.*

648 *Then, there exists a function $\phi = \phi_{\tilde{\varrho}}$ generated by a $\tilde{\varrho}$ -activated network with width $N \cdot \tilde{N}$
 649 and depth $L \cdot \tilde{L}$ such that*

$$650 \quad |\phi(\mathbf{x}) - f(\mathbf{x})| < \varepsilon \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

651 The proof of Lemma 9 is placed in Section 4.3. Based on Lemma 9, we will propose
 652 two UAFs with better mathematical properties. That is, the idea of designing a C^s UAF
 653 is given in Section 4.1 and a sigmoidal UAF is constructed in detail in Section 4.2.

654 **4.1 Smooth UAF**

655 The smoothness of a function is one of the most desired properties in mathematical modeling
 656 and computation. The EUAF σ is continuous but not smooth. So we will show how to
 657 construct a C^s UAF based on an existing one. The key point is the fact that the indefinite
 658 integral of a continuous function is continuously differentiable.

659 Suppose ϱ is a continuous UAF. Define

$$660 \quad \tilde{\varrho}(x) := \int_0^x \varrho(t) dt \quad \text{for any } x \in \mathbb{R}.$$

661 For any $M > 0$, it holds that

$$662 \quad \frac{\tilde{\varrho}(x + \delta) - \tilde{\varrho}(x)}{\delta} = \frac{1}{\delta} \int_x^{x+\delta} \varrho(t) dt \rightrightarrows \varrho(x) \quad \text{as } \delta \rightarrow 0^+ \quad \text{for any } x \in [-M, M].$$

663 This means ϱ can be approximated by a one-hidden-layer $\tilde{\varrho}$ -activated network with width
 664 2 arbitrarily well on any bounded interval. It follows that $\tilde{\varrho}$ is also a UAF. By repeated
 665 applications of the above idea, one could easily construct a C^s UAF.

666 In particular, set $\varrho_0 = \sigma$ and define $\varrho_1, \varrho_2, \dots, \varrho_s$ by induction as follows.

$$667 \quad \varrho_{i+1}(x) := \int_0^x \varrho_i(t) dt \quad \text{for any } x \in \mathbb{R} \text{ and } i \in \{0, 1, \dots, s-1\}. \quad (6)$$

668 Then ϱ_s is a C^s UAF as shown in the following theorem.

669 **Theorem 10.** *Let $\varrho_s \in C^s(\mathbb{R})$ be the function defined in Equation (6) for any $s \in \mathbb{N}^+$.
 670 Then, for any $f \in C([a, b]^d)$ and any $\varepsilon > 0$, there exists a function ϕ generated by a
 671 ϱ_s -activated network with width $72sd(2d+1)$ and depth 11 such that*

$$672 \quad |\phi(\mathbf{x}) - f(\mathbf{x})| < \varepsilon \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

673 *Proof.* For any $i \in \{0, 1, \dots, s-1\}$ and any $M > 0$, it is easy to verify that

$$674 \quad \frac{\varrho_{i+1}(x + \delta) - \varrho_{i+1}(x)}{\delta} = \frac{1}{\delta} \int_x^{x+\delta} \varrho_i(t) dt \rightrightarrows \varrho_i(x) \quad \text{as } \delta \rightarrow 0^+ \quad \text{for any } x \in [-M, M].$$

675 This means ϱ_i can be approximated by a one-hidden-layer ϱ_{i+1} -activated network with
 676 width 2 arbitrarily well on any bounded interval. By induction, one could easily prove that
 677 $\varrho_0 = \sigma$ can be approximated by a one-hidden-layer ϱ_s -activated network with width $2s$
 678 arbitrarily well on any bounded interval. That is, for each $\delta \in (0, 1)$, there exists a function
 679 $\sigma_{s,\delta}$ realized by a ϱ_s -activated network with width $2s$ and depth 1 such that

$$680 \quad \sigma_{s,\delta}(t) \rightrightarrows \sigma(t) \quad \text{as } \delta \rightarrow 0^+ \quad \text{for any } t \in [-M, M].$$

681 By Theorem 1, there exists a function ϕ_σ generated by a σ -activated network with width
 682 $36d(2d+1)$ and depth 11 such that

$$683 \quad |\phi_\sigma(\mathbf{x}) - f(\mathbf{x})| < \varepsilon/2 \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

684 Then, by Lemma 9, there exists another function $\phi = \phi_{\varrho_s}$ realized by a ϱ_s -activated network
 685 with width $2s \times 36d(2d+1) = 72sd(2d+1)$ and depth $1 \times 11 = 11$ such that

$$686 \quad |\phi(\mathbf{x}) - f(\mathbf{x})| < \varepsilon \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

687 So we finish the proof. ■

688 **4.2 Sigmoidal UAF**

689 Many activation functions used in real-world applications are sigmoidal functions. Gener-
 690 ally, we say a function $g : \mathbb{R} \rightarrow \mathbb{R}$ is sigmoidal (or sigmoid, e.g., see (Han and Moraga,
 691 1995)) if it satisfies the following conditions.

- 692 • Bounded: $\lim_{x \rightarrow \infty} g(x) = 1$ and $\lim_{x \rightarrow -\infty} g(x) = -1$ (or 0).
- 693 • Differentiable: $g'(x)$ exists and continuous for all $x \in \mathbb{R}$.
- 694 • Increasing: $g'(x)$ is non-negative for all $x \in \mathbb{R}$.

695 Our goal is to construct a sigmoidal UAF. To this end, we need to design a new function
 696 $\tilde{\sigma}$ based on σ such that σ can be reproduced/approximated by a $\tilde{\sigma}$ -activated network with
 697 a fixed size. Making $\tilde{\sigma}$ bounded and increasing is not difficult. The key is to make $\tilde{\sigma}$
 698 continuously differentiable, which can be implemented by the fact that the indefinite integral
 699 of a continuous function is continuously differentiable. To be exact, we can define $\tilde{\sigma}$ as
 700 follows.

- 701 • For $x \in (-\infty, 0]$, define $\tilde{\sigma}(x) := \sigma(x) = \frac{x}{-x+1}$.
- 702 • For $x \in (0, \infty)$, define

$$703 \quad \tilde{\sigma}(x) := \int_0^x \frac{c\sigma(t) + 1}{(2t + 1)^2} dt, \quad \text{where } c = \frac{1}{2 \int_0^\infty \frac{\sigma(t)}{(2t+1)^2} dt} \approx 2.554.$$

704 We remark that there are many possible choices for the integrand in the above definition
 705 of $\tilde{\sigma}(x)$ for $x \in (0, \infty)$. Here, we just give a simple example. See an illustration of $\tilde{\sigma}$ in
 706 Figure 8.

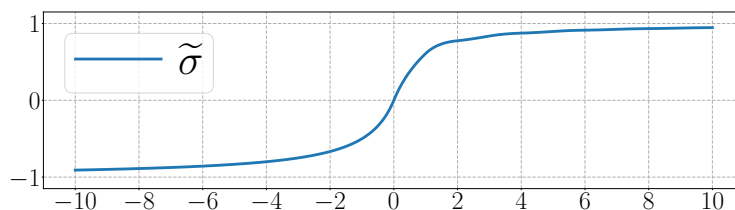


Figure 8: An illustration of $\tilde{\sigma}$ on $[-10, 10]$.

707 Then $\tilde{\sigma}$ is a sigmoidal function as verified below.

- 708 • Clearly, $\lim_{x \rightarrow -\infty} \tilde{\sigma}(x) = \lim_{x \rightarrow -\infty} \frac{x}{-x+1} = -1$. Moreover,

$$709 \quad \lim_{x \rightarrow \infty} \tilde{\sigma}(x) = \int_0^\infty \frac{c\sigma(t) + 1}{(2t + 1)^2} dt = \frac{1}{2} + \int_0^\infty \frac{1}{(2t + 1)^2} dt = 1.$$

- 710 • Obviously, $\tilde{\sigma}$ is continuously differentiable on $(-\infty, 0)$ and $(0, \infty)$. Meanwhile, we
 711 have $\tilde{\sigma}'(0) = 1$ and $\lim_{x \rightarrow 0} \tilde{\sigma}'(x) = 1$. Therefore, we have $\tilde{\sigma} \in C^1(\mathbb{R})$ as desired.

712 • For $x \in (-\infty, 0)$, $\tilde{\sigma}'(x) = \frac{1}{(-x+1)^2} > 0$. For $x = 0$, $\tilde{\sigma}'(x) = 1 > 0$. For $x \in (0, \infty)$,
713 $\tilde{\sigma}'(x) = \frac{c\sigma(x)+1}{(2x+1)^2} > 0$. Therefore, $\tilde{\sigma}'(x) > 0$ for all $x \in \mathbb{R}$.

714 Based on Theorem 1 corresponding to σ , we establish a similar theorem for $\tilde{\sigma}$, Theo-
715 rem 11 below, showing that fixed-size $\tilde{\sigma}$ -activated networks can also approximate continuous
716 functions within an arbitrary error on a hypercube.

717 **Theorem 11.** For any $f \in C([a, b]^d)$ and any $\varepsilon > 0$, there exists a function ϕ generated by
718 a $\tilde{\sigma}$ -activated network with width $1800d(2d + 1)$ and depth 66 such that

$$719 \quad |\phi(\mathbf{x}) - f(\mathbf{x})| < \varepsilon \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

720 To prove this theorem based on Theorem 1, we only need to show σ can be approximated
721 by a fixed-size $\tilde{\sigma}$ -activated network within an arbitrary error on any pre-specified interval
722 as presented in the following lemma.

723 **Lemma 12.** For any $\varepsilon > 0$ and any $M > 0$, there exists a function ϕ realized by a $\tilde{\sigma}$ -
724 activated network with width 50 and depth 6 such that

$$725 \quad |\phi(x) - \sigma(x)| < \varepsilon \quad \text{for any } x \in [-M, M].$$

726 The proof of Lemma 12 can be found later. By assuming Lemma 12 is true, we can give
727 the proof of Theorem 11.

728 *Proof of Theorem 11.* By Theorem 1, there exists a function ϕ_σ generated by a σ -activated
729 network with width $36d(2d + 1)$ and depth 11 such that

$$730 \quad |\phi_\sigma(\mathbf{x}) - f(\mathbf{x})| < \varepsilon/2 \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

731 By Lemma 12, for any $M > 0$ and each $\delta \in (0, 1)$, there exists a function σ_δ realized by a
732 $\tilde{\sigma}$ -activated network with width 50 and depth 6 such that

$$733 \quad \sigma_\delta(t) \rightrightarrows \sigma(t) \quad \text{as } \delta \rightarrow 0^+ \quad \text{for any } t \in [-M, M].$$

734 Then, by Lemma 9, there exists another function $\phi = \phi_{\tilde{\sigma}}$ realized by a $\tilde{\sigma}$ -activated network
735 with width $50 \times 36d(2d + 1) = 1800d(2d + 1)$ and depth $6 \times 11 = 66$ such that

$$736 \quad |\phi(\mathbf{x}) - f(\mathbf{x})| < \varepsilon \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

737 So we finish the proof. ■

738 Finally, let us present the detailed proof of Lemma 12.

739 *Proof of Lemma 12.* Since $1 = \tilde{\sigma}'(0) = \lim_{x \rightarrow 0} \frac{\tilde{\sigma}(x)}{x}$, it is easy to show: For any $\mathcal{E} > 0$ and
740 any $R > 0$, there exists a sufficiently small $w > 0$ such that

$$741 \quad |\tilde{\sigma}(wx)/w - x| < \mathcal{E} \quad \text{for any } x \in [-R, R].$$

742 Thus, we may assume the identity map is allowed to be the activation function in $\tilde{\sigma}$ -activated
743 networks. Without loss of generality, we may assume $M \geq 2$ because $\widehat{M} = \max\{2, M\}$
744 implies $\widehat{M} \geq 2$ and $[-M, M] \subseteq [-\widehat{M}, \widehat{M}]$.

745 For simplicity, we denote $\widehat{\mathcal{H}}(N, L)$ as the (hypothesis) space of functions generated by
746 $\tilde{\sigma}$ -activated networks with width N and depth L . Then the proof can be roughly divided
747 into three steps as follows.

748 (1) Design $\Gamma \in \widetilde{\mathcal{H}}(9, 2)$ to reproduce xy on $[-4\widetilde{M}, 4\widetilde{M}]^2$, where $\widetilde{M} = (M + 1)^2$.

749 (2) Design $\psi_\delta \in \widetilde{\mathcal{H}}(9, 4)$ based on the first step to approximate σ well on $[0, M]$.

750 (3) Design $\phi \in \widetilde{\mathcal{H}}(50, 6)$ based on the previous two steps to approximate σ well on $[-M, M]$.

751 The details of these three steps can be found below.

752 **Step 1:** Design $\Gamma \in \widetilde{\mathcal{H}}(9, 2)$ to reproduce xy on $[-4\widetilde{M}, 4\widetilde{M}]^2$.

753 Observe that

$$754 \quad \widetilde{\sigma}(y) + 1 = \frac{y}{|y| + 1} + 1 = \frac{y}{-y + 1} + 1 = \frac{1}{-y + 1} \quad \text{for any } y \leq 0.$$

755 For any $x \in [-4, 4]$, we have $-x - 4 \leq 0$ and $-x - 5 \leq 0$, implying

$$\begin{aligned} 756 \quad \widetilde{\sigma}(-x - 4) - \widetilde{\sigma}(-x - 5) &= \left(\widetilde{\sigma}(-x - 4) + 1 \right) - \left(\widetilde{\sigma}(-x - 5) + 1 \right) \\ &= \frac{1}{-(-x - 4) + 1} - \frac{1}{-(-x - 5) + 1} \\ &= \frac{1}{x + 5} - \frac{1}{x + 6} = \frac{1}{(x + 5)(x + 6)}. \end{aligned}$$

757 It follows from $1 - \frac{90}{(x+5)(x+6)} \leq 0$ for any $x \in [-4, 4]$ that

$$758 \quad \widetilde{\sigma}\left(1 - \frac{90}{(x + 5)(x + 6)}\right) + 1 = \frac{1}{-\left(1 - \frac{90}{(x+5)(x+6)}\right) + 1} = \frac{x^2 + 11x + 30}{90},$$

759 implying

$$\begin{aligned} 760 \quad x^2 &= 90\widetilde{\sigma}\left(1 - \frac{90}{(x + 5)(x + 6)}\right) + 90 - (11x + 30) \\ &= 90\widetilde{\sigma}\left(1 - 90(\widetilde{\sigma}(-x - 4) - \widetilde{\sigma}(-x - 5))\right) - 11x + 60 \\ &= 90\widetilde{\sigma}\left(1 - 90\widetilde{\sigma}(-x - 4) + 90\widetilde{\sigma}(-x - 5)\right) - 11x + 60. \end{aligned}$$

761 Thus, x^2 can be realized by a $\widetilde{\sigma}$ -activated network with width 3 and depth 2 on $[-4, 4]$. Set
762 $\widetilde{M} = (M + 1)^2$. Then, for any $x, y \in [-4\widetilde{M}, 4\widetilde{M}]$, we have $\frac{x}{2\widetilde{M}}, \frac{y}{2\widetilde{M}}, \frac{x+y}{2\widetilde{M}} \in [-4, 4]$. Recall
763 the fact

$$764 \quad xy = 2\widetilde{M}^2 \left(\left(\frac{x+y}{2\widetilde{M}} \right)^2 - \left(\frac{x}{2\widetilde{M}} \right)^2 - \left(\frac{y}{2\widetilde{M}} \right)^2 \right).$$

765 Therefore, xy can be realized by a $\widetilde{\sigma}$ -activated network with width 9 and depth 2 for any
766 $x, y \in [-4\widetilde{M}, 4\widetilde{M}]$. That is, there exists $\Gamma \in \widetilde{\mathcal{H}}(9, 2)$ such that $\Gamma(x, y) = xy$ on $[-4\widetilde{M}, 4\widetilde{M}]^2$.

767 **Step 2:** Design $\psi_\delta \in \widetilde{\mathcal{H}}(9, 4)$ to approximate σ well on $[0, M]$.

768 Recall that x^2 can be realized by a $\widetilde{\sigma}$ -activated network with width 3 and depth 2 on
769 $[-4, 4]$. There exists $\psi_1 \in \widetilde{\mathcal{H}}(3, 2)$ such that

$$770 \quad \psi_1(x) = \frac{(2x + 1)^2}{(2M + 1)^2} \quad \text{for any } x \in [-M, M].$$

771 For any small $\delta > 0$, we define

$$772 \quad \psi_{2,\delta}(x) := \frac{\tilde{\sigma}(x + \delta) - \tilde{\sigma}(x)}{\delta} \quad \text{for any } x \in \mathbb{R}.$$

773 Then, we have $\psi_{2,\delta} \in \widetilde{\mathcal{H}}(2, 1)$ and

$$774 \quad \psi_{2,\delta}(x) := \frac{\tilde{\sigma}(x + \delta) - \tilde{\sigma}(x)}{\delta} \rightrightarrows \frac{d}{dx} \tilde{\sigma}(x) = \frac{c\sigma(x) + 1}{(2x + 1)^2} \quad \text{as } \delta \rightarrow 0^+$$

775 for any $x \in [0, M]$, where c is a constant given by

$$776 \quad c = \frac{1}{2 \int_0^\infty \frac{\sigma(t)}{(2t+1)^2} dt} \approx 2.554.$$

777 For any small $\delta > 0$, we define

$$778 \quad \psi_\delta(x) := \frac{(2M+1)^2}{c} \Gamma(\psi_1(x), \psi_{2,\delta}(x)) - \frac{1}{c} \quad \text{for any } x \in \mathbb{R}.$$

779 Since $\Gamma \in \widetilde{\mathcal{H}}(9, 2)$, $\psi_1 \in \widetilde{\mathcal{H}}(3, 2)$, and $\psi_{2,\delta} \in \widetilde{\mathcal{H}}(2, 1)$, we have $\psi_\delta \in \widetilde{\mathcal{H}}(9, 4)$.

780 Clearly, for any $x \in [0, M]$, we have $\psi_1(x) = \frac{(2x+1)^2}{(2M+1)^2} \in [0, 1]$ and $\psi_{2,\delta}(x) \rightrightarrows \frac{c\sigma(x)+1}{(2x+1)^2} \in$
 781 $[0, c+1] \subseteq [0, 3.6]$, implying $\psi_1(x), \psi_{2,\delta}(x) \in [-4, 4] \subseteq [-4\widetilde{M}, 4\widetilde{M}]$ for any small $\delta > 0$.
 782 Thus, for any $x \in [0, M]$, as δ goes to 0^+ , we have

$$783 \quad \begin{aligned} \psi_\delta(x) &= \frac{(2M+1)^2}{c} \Gamma(\psi_1(x), \psi_{2,\delta}(x)) - \frac{1}{c} = \frac{(2M+1)^2}{c} \cdot \psi_1(x) \cdot \psi_{2,\delta}(x) - \frac{1}{c} \\ &\rightrightarrows \frac{(2M+1)^2}{c} \cdot \frac{(2x+1)^2}{(2M+1)^2} \cdot \frac{c\sigma(x)+1}{(2x+1)^2} - \frac{1}{c} = \sigma(x). \end{aligned}$$

784 That is, for any $x \in [0, M]$,

$$785 \quad \psi_\delta(x) \rightrightarrows \sigma(x) \quad \text{as } \delta \rightarrow 0^+.$$

786 **Step 3:** Design $\phi \in \widetilde{\mathcal{H}}(50, 6)$ to approximate σ well on $[-M, M]$.

787 Note that $\tilde{\sigma}(x) = \sigma(x)$ for all $x \in [-M, 0)$ and $\psi_\delta(x)$ approximates $\sigma(x)$ well for all
 788 $x \in [0, M]$. Then, we have

$$789 \quad \psi_\delta(x) \cdot \mathbf{1}_{\{x \in [0, M]\}} + \tilde{\sigma}(x) \cdot \mathbf{1}_{\{x \in [-M, 0)\}}$$

790 approximates $\sigma(x)$ well for all $x \in [-M, M]$. However, it is impossible to approximate
 791 $\mathbf{1}_{\{x \in [0, M]\}}$ well by a $\tilde{\sigma}$ -activated network due to the continuity of $\tilde{\sigma}$. To address this gap,
 792 we will construct a continuous function g to replace $\mathbf{1}_{\{x \in [0, M]\}}$ such that

$$793 \quad \psi_\delta(x) \cdot g(x) + \tilde{\sigma}(x) \cdot (1 - g(x)) \tag{7}$$

794 can also approximate $\sigma(x)$ well for all $x \in [-M, M]$.

795 By the continuity of $\tilde{\sigma}$ and σ , there exists a small $\eta_0 \in (0, 1)$ such that

$$796 \quad |\tilde{\sigma}(x)| < \varepsilon/6 \quad \text{and} \quad |\sigma(x)| < \varepsilon/6 \quad \text{for any } x \in [0, \eta_0]. \tag{8}$$

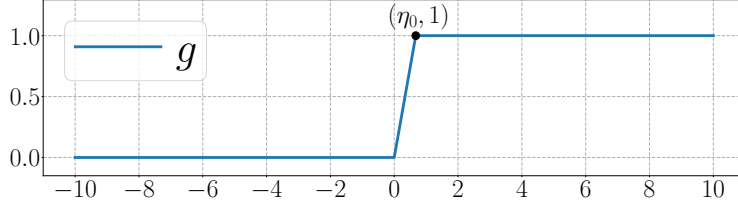


Figure 9: An illustration of g on $[-10, 10]$.

797 Then we define

$$798 \quad g(x) := \frac{\text{ReLU}(x) - \text{ReLU}(x - \eta_0)}{\eta_0}, \quad \text{where } \text{ReLU}(x) = \max\{0, x\} \quad \text{for any } x \in \mathbb{R}.$$

799 See Figure 9 for an illustration of g .

800 We will construct a $\tilde{\sigma}$ -activated network to approximate g well. To this end, we first
 801 design a $\tilde{\sigma}$ -activated network to approximate the ReLU function well. For any $x \in [-M -$
 802 $1, M + 1]$, we have $\frac{x}{M+1} + 1 \in [0, 2] \subseteq [0, M]$, implying

$$803 \quad 1 - \psi_\delta\left(\frac{x}{M+1} + 1\right) \rightrightarrows 1 - \sigma\left(\frac{x}{M+1} + 1\right) = \left|\frac{x}{M+1}\right| \quad \text{as } \delta \rightarrow 0^+,$$

804 where the last equality comes from $1 - \sigma(y) = |y - 1|$ for any $y \in [0, 2]$. Recall that

$$805 \quad \text{ReLU}(x) = \frac{x}{2} + \frac{|x|}{2} = \frac{x}{2} + \frac{M+1}{2} \cdot \left|\frac{x}{M+1}\right|$$

806 for any $x \in [-M - 1, M + 1]$. For any small $\delta > 0$, we define

$$807 \quad \tilde{g}_\delta(x) := \frac{x}{2} + \frac{M+1}{2} \left(1 - \psi_\delta\left(\frac{x}{M+1} + 1\right)\right) \quad \text{for any } x \in \mathbb{R}.$$

808 Then, $\psi_\delta \in \widetilde{\mathcal{H}}(9, 4)$ implies $\tilde{g}_\delta \in \widetilde{\mathcal{H}}(10, 4)$. Moreover, for any $x \in [-M - 1, M + 1]$,

$$809 \quad \tilde{g}_\delta(x) \rightrightarrows \frac{x}{2} + \frac{M+1}{2} \cdot \left|\frac{x}{M+1}\right| = \text{ReLU}(x) \quad \text{as } \delta \rightarrow 0^+.$$

810 Define

$$811 \quad g_\delta(x) := \frac{\tilde{g}_\delta(x) - \tilde{g}_\delta(x - \eta_0)}{\eta_0} \quad \text{for any } x \in \mathbb{R}.$$

812 Clearly, $\tilde{g}_\delta \in \widetilde{\mathcal{H}}(10, 4)$ implies $g_\delta \in \widetilde{\mathcal{H}}(20, 4)$. For any $x \in [-M, M]$, we have $x, x - \eta_0 \in$
 813 $[-M - 1, M + 1]$, implying

$$814 \quad g_\delta(x) = \frac{\tilde{g}_\delta(x) - \tilde{g}_\delta(x - \eta_0)}{\eta_0} \rightrightarrows \frac{\text{ReLU}(x) - \text{ReLU}(x - \eta_0)}{\eta_0} = g(x) \quad \text{as } \delta \rightarrow 0^+.$$

815 Next, motivated by Equation (7), we can define ϕ_δ to approximate σ well on $[-M, M]$.
 816 The definition of ϕ_δ is given by

$$817 \quad \phi_\delta(x) := \Gamma\left(\psi_\delta(x), g_\delta(x)\right) + \Gamma\left(\tilde{\sigma}(x), 1 - g_\delta(x)\right) \quad \text{for any } x \in \mathbb{R}.$$

818 Since $\Gamma \in \widetilde{\mathcal{H}}(9, 2)$, $\psi_\delta \in \widetilde{\mathcal{H}}(9, 4)$, and $g_\delta, 1 - g_\delta \in \widetilde{\mathcal{H}}(20, 4)$, we have

819
$$\phi_\delta \in \widetilde{\mathcal{H}}(9 + 20 + 1 + 20, 4 + 2) = \widetilde{\mathcal{H}}(50, 6).$$

820 Clearly, $\tilde{\sigma}(x)$, $g_\delta(x)$, and $1 - g_\delta(x)$ are all in $[-4\widetilde{M}, 4\widetilde{M}]$ for any small $\delta > 0$ and all
 821 $x \in [-M, M]$. We will show $\psi_\delta(x) \in [-4\widetilde{M}, 4\widetilde{M}]$ for any small $\delta > 0$ and all $x \in [-M, M]$
 822 via two cases as follows.

823 • For any $x \in [0, M]$, $\psi_\delta(x) \rightrightarrows \sigma(x)$ implies $\psi_\delta(x) \in [-4\widetilde{M}, 4\widetilde{M}]$ for any small $\delta > 0$.

824 • For any $x \in [-M, 0)$, we have $\psi_1(x) = \frac{(2x+1)^2}{(2M+1)^2} \in [0, 1]$ and

825
$$\psi_{2,\delta}(x) = \frac{\tilde{\sigma}(x+\delta) - \tilde{\sigma}(x)}{\delta} \rightrightarrows \frac{d}{dx} \tilde{\sigma}(x) = \frac{1}{(-x+1)^2} \quad \text{as } \delta \rightarrow 0^+.$$

826 Thus, for any $x \in [-M, 0)$, as δ goes to 0^+ , we get

827
$$\begin{aligned} \psi_\delta(x) &= \frac{(2M+1)^2}{c} \Gamma\left(\psi_1(x), \psi_{2,\delta}(x)\right) - \frac{1}{c} = \frac{(2M+1)^2}{c} \cdot \psi_1(x) \cdot \psi_{2,\delta}(x) - \frac{1}{c} \\ &\rightrightarrows \frac{(2M+1)^2}{c} \cdot \frac{(2x+1)^2}{(2M+1)^2} \cdot \frac{1}{(-x+1)^2} - \frac{1}{c} = \frac{(2x+1)^2 - 1}{c(-x+1)^2}. \end{aligned}$$

828 For all $x \in [-M, 0)$, we have $c(-x+1)^2 \geq 1$, implying $\frac{(2x+1)^2 - 1}{c(-x+1)^2} \geq \frac{-1}{c(-x+1)^2} \geq -1$ and

829
$$\begin{aligned} \frac{(2x+1)^2 - 1}{c(-x+1)^2} &\leq \frac{(2|x|+1)^2 - 1}{c(-x+1)^2} \leq (2|x| + 1)^2 - 1 = 4(|x| + 1/2)^2 - 1 \\ &\leq 4(M + 1)^2 - 1 = 4\widetilde{M} - 1. \end{aligned}$$

830 That is, $\frac{(2x+1)^2 - 1}{c(-x+1)^2} \in [-1, 4\widetilde{M} - 1]$ for all $x \in [-M, 0)$, implying $\psi_\delta(x) \in [-4\widetilde{M}, 4\widetilde{M}]$
 831 for any small $\delta > 0$.

832 Hence, for any $x \in [\eta_0, M]$, we have $1 - g(x) = 0$, implying

833
$$\phi_\delta(x) = \psi_\delta(x) \cdot g_\delta(x) + \tilde{\sigma}(x) \cdot (1 - g_\delta(x)) \rightrightarrows \sigma(x) \cdot g(x) + 0 = \sigma(x) \quad \text{as } \delta \rightarrow 0^+.$$

834 Similarly, for any $x \in [-M, 0]$, we have $g(x) = 0$, implying

835
$$\phi_\delta(x) = \psi_\delta(x) \cdot g_\delta(x) + \tilde{\sigma}(x) \cdot (1 - g_\delta(x)) \rightrightarrows 0 + \tilde{\sigma}(x) \cdot (1 - g(x)) = \sigma(x) \quad \text{as } \delta \rightarrow 0^+.$$

836 Therefore, there exists a small $\delta_0 > 0$ such that

837
$$|\phi_{\delta_0}(x) - \sigma(x)| < \varepsilon \quad \text{for any } x \in [-M, 0] \cup [\eta_0, M],$$

838
$$\|g_{\delta_0}\|_{L^\infty([0, \eta_0])} \leq 2, \quad \|1 - g_{\delta_0}\|_{L^\infty([0, \eta_0])} \leq 2, \quad \text{and}$$

839
$$\|\psi_{\delta_0}\|_{L^\infty([0, \eta_0])} \leq \|\sigma\|_{L^\infty([0, \eta_0])} + \varepsilon/12,$$

840 where the above inequality comes from the fact $\psi_\delta(x)$ uniformly converges to $\sigma(x)$ for any
 841 $x \in [0, \eta_0] \subseteq [0, M]$.

842 Clearly, for any $x \in [0, \eta_0]$, by Equation (8), we have

$$\begin{aligned}
|\phi_{\delta_0}(x) - \sigma(x)| &\leq |\phi_{\delta_0}(x)| + |\sigma(x)| < \left| \psi_{\delta_0}(x) \cdot g_{\delta_0}(x) + \tilde{\sigma}(x) \cdot (1 - g_{\delta_0}(x)) \right| + \varepsilon/6 \\
&\leq |\psi_{\delta_0}(x)| \cdot |g_{\delta_0}(x)| + |\tilde{\sigma}(x)| \cdot |1 - g_{\delta_0}(x)| + \varepsilon/6 \\
843 \quad &\leq (\|\sigma\|_{L^\infty([0, \eta_0])} + \frac{\varepsilon}{12}) \cdot 2 + \frac{\varepsilon}{6} \cdot 2 + \frac{\varepsilon}{6} \\
&\leq \left(\frac{\varepsilon}{6} + \frac{\varepsilon}{12}\right) \cdot 2 + \frac{\varepsilon}{6} \cdot 2 + \frac{\varepsilon}{6} = \varepsilon.
\end{aligned}$$

844 By setting $\phi = \phi_{\delta_0}$, we have $\phi = \phi_{\delta_0} \in \widetilde{\mathcal{H}}(50, 6)$ and

$$845 \quad |\phi(x) - \sigma(x)| = |\phi_{\delta_0}(x) - \sigma(x)| < \varepsilon \quad \text{for any } x \in [-M, M].$$

846 So we finish the proof. ■

847 4.3 Proof of Lemma 9

848 Let the activation function be applied to a vector elementwisely. Then ϕ_ϱ can be represented
849 in a form of function compositions as follows:

$$850 \quad \phi_\varrho(\mathbf{x}) = \mathcal{L}_L \circ \varrho \circ \mathcal{L}_{L-1} \circ \cdots \circ \varrho \circ \mathcal{L}_1 \circ \varrho \circ \mathcal{L}_0(\mathbf{x}) \quad \text{for any } \mathbf{x} \in \mathbb{R}^d,$$

851 where $N_0 = d$, $N_1, N_2, \dots, N_L \in \mathbb{N}^+$, $N_{L+1} = 1$, $\mathbf{A}_\ell \in \mathbb{R}^{N_{\ell+1} \times N_\ell}$ and $\mathbf{b}_\ell \in \mathbb{R}^{N_{\ell+1}}$ are the
852 weight matrix and the bias vector in the ℓ -th affine linear transform $\mathcal{L}_\ell : \mathbf{y} \mapsto \mathbf{A}_\ell \mathbf{y} + \mathbf{b}_\ell$ for
853 each $\ell \in \{0, 1, \dots, L\}$. Define

$$854 \quad \phi_{\varrho_\delta}(\mathbf{x}) := \mathcal{L}_L \circ \varrho_\delta \circ \mathcal{L}_{L-1} \circ \cdots \circ \varrho_\delta \circ \mathcal{L}_1 \circ \varrho_\delta \circ \mathcal{L}_0(\mathbf{x}) \quad \text{for any } \mathbf{x} \in \mathbb{R}^d.$$

855 Recall that ϱ_δ can be realized by a $\tilde{\varrho}$ -activated network with width \tilde{N} and depth \tilde{L} . Thus,
856 ϕ_{ϱ_δ} can be realized by a $\tilde{\varrho}$ -activated network with width $N \cdot \tilde{N}$ and depth $L \cdot \tilde{L}$. We will
857 prove

$$858 \quad \phi_{\varrho_\delta}(\mathbf{x}) \rightrightarrows \phi_\varrho(\mathbf{x}) \quad \text{as } \delta \rightarrow 0^+ \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

859 For any $\mathbf{x} \in \mathbb{R}^d$ and each $\ell \in \{1, 2, \dots, L+1\}$, define

$$860 \quad \mathbf{h}_\ell(\mathbf{x}) := \mathcal{L}_{\ell-1} \circ \varrho \circ \mathcal{L}_{\ell-2} \circ \cdots \circ \varrho \circ \mathcal{L}_1 \circ \varrho \circ \mathcal{L}_0(\mathbf{x})$$

861 and

$$862 \quad \mathbf{h}_{\ell, \delta}(\mathbf{x}) := \mathcal{L}_{\ell-1} \circ \varrho_\delta \circ \mathcal{L}_{\ell-2} \circ \cdots \circ \varrho_\delta \circ \mathcal{L}_1 \circ \varrho_\delta \circ \mathcal{L}_0(\mathbf{x}).$$

863 Note that \mathbf{h}_ℓ and $\mathbf{h}_{\ell, \delta}$ are two maps from \mathbb{R}^d to \mathbb{R}^{N_ℓ} for each ℓ .

864 We will prove by induction that

$$865 \quad \mathbf{h}_{\ell, \delta}(\mathbf{x}) \rightrightarrows \mathbf{h}_\ell(\mathbf{x}) \quad \text{as } \delta \rightarrow 0^+ \tag{9}$$

866 for any $\mathbf{x} \in [a, b]^d$ and each $\ell \in \{1, 2, \dots, L+1\}$.

867 First, we consider the case $\ell = 1$. Clearly,

$$868 \quad \mathbf{h}_{1, \delta}(\mathbf{x}) = \mathcal{L}_0(\mathbf{x}) = \mathbf{h}_1(\mathbf{x}) \quad \text{as } \delta \rightarrow 0^+ \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

869 This means Equation (9) holds for $\ell = 1$.

870 Next, suppose Equation (9) holds for $\ell = i \in \{1, 2, \dots, L\}$. Our goal is to prove that it
871 also holds for $\ell = i + 1$. Determine $M > 0$ by defining

$$872 \quad M := \sup \left\{ \|\mathbf{h}_j(\mathbf{x})\|_\infty + 1 : \mathbf{x} \in [a, b]^d, \quad j = 1, 2, \dots, L + 1 \right\},$$

873 where the continuity of ϱ guarantees the above supremum is finite, i.e., $M \in (1, \infty)$. By
874 the induction hypothesis, we have

$$875 \quad \mathbf{h}_{i,\delta}(\mathbf{x}) \rightrightarrows \mathbf{h}_i(\mathbf{x}) \quad \text{as } \delta \rightarrow 0^+ \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

876 Clearly, for any $\mathbf{x} \in [a, b]^d$, we have $\|\mathbf{h}_i(\mathbf{x})\|_\infty \leq M$ and $\|\mathbf{h}_{i,\delta}(\mathbf{x})\|_\infty \leq \|\mathbf{h}_i(\mathbf{x})\|_\infty + 1 \leq M$
877 for any small $\delta > 0$.

878 Recall the fact $\varrho_\delta(t) \rightrightarrows \varrho(t)$ as $\delta \rightarrow 0^+$ for any $t \in [-M, M]$. Then, we have

$$879 \quad \varrho_\delta \circ \mathbf{h}_{i,\delta}(\mathbf{x}) - \varrho \circ \mathbf{h}_{i,\delta}(\mathbf{x}) \rightrightarrows \mathbf{0} \quad \text{as } \delta \rightarrow 0^+ \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

880 The continuity of ϱ implies the uniform continuity of ϱ on $[-M, M]$, from which we deduce

$$881 \quad \varrho \circ \mathbf{h}_{i,\delta}(\mathbf{x}) - \varrho \circ \mathbf{h}_i(\mathbf{x}) \rightrightarrows \mathbf{0} \quad \text{as } \delta \rightarrow 0^+ \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

882 Therefore, for any $\mathbf{x} \in [a, b]^d$, as $\delta \rightarrow 0^+$, we have

$$883 \quad \varrho_\delta \circ \mathbf{h}_{i,\delta}(\mathbf{x}) - \varrho \circ \mathbf{h}_i(\mathbf{x}) = \underbrace{\varrho_\delta \circ \mathbf{h}_{i,\delta}(\mathbf{x}) - \varrho \circ \mathbf{h}_{i,\delta}(\mathbf{x})}_{\rightrightarrows \mathbf{0}} + \underbrace{\varrho \circ \mathbf{h}_{i,\delta}(\mathbf{x}) - \varrho \circ \mathbf{h}_i(\mathbf{x})}_{\rightrightarrows \mathbf{0}} \rightrightarrows \mathbf{0},$$

884 implying

$$885 \quad \mathbf{h}_{i+1,\delta}(\mathbf{x}) = \mathcal{L}_i \circ \varrho_\delta \circ \mathbf{h}_{i,\delta}(\mathbf{x}) \rightrightarrows \mathcal{L}_i \circ \varrho \circ \mathbf{h}_i(\mathbf{x}) = \mathbf{h}_{i+1}(\mathbf{x}).$$

886 This means Equation (9) holds for $\ell = i + 1$. So we complete the inductive step.

887 By the principle of induction, we have

$$888 \quad \phi_{\varrho_\delta}(\mathbf{x}) = \mathbf{h}_{L+1,\delta}(\mathbf{x}) \rightrightarrows \mathbf{h}_{L+1}(\mathbf{x}) = \phi_\varrho(\mathbf{x}) \quad \text{as } \delta \rightarrow 0^+ \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

889 There exists a small $\delta_0 > 0$ such that

$$890 \quad |\phi_{\varrho_{\delta_0}}(\mathbf{x}) - \phi_\varrho(\mathbf{x})| < \varepsilon/2 \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

891 By defining $\phi := \phi_{\varrho_{\delta_0}}$, we have

$$892 \quad |\phi(\mathbf{x}) - f(\mathbf{x})| \leq |\phi_{\varrho_{\delta_0}}(\mathbf{x}) - \phi_\varrho(\mathbf{x})| + |\phi_\varrho(\mathbf{x}) - f(\mathbf{x})| < \varepsilon/2 + \varepsilon/2 = \varepsilon$$

893 for any $\mathbf{x} \in [a, b]^d$. Moreover, $\phi = \phi_{\varrho_{\delta_0}}$ can be generated by a $\tilde{\varrho}$ -activated network with
894 width $N \cdot \tilde{N}$ and depth $L \cdot \tilde{L}$. So we finish the proof.

895 5. Detailed Proofs of Theorems 1 and 4

896 In this section, we will give the detailed proofs of Theorems 1 and 4. First, we prove The-
897 orem 1 based on Theorem 6, which will be proved in Section 6. Next, we apply Theorem 1
898 to prove Theorem 4.

899 **5.1 Proof of Theorem 1**

900 The detailed proof of Theorem 1 converts the above ideas mentioned in Section 2.2 to
 901 implementations using neural networks with fixed sizes. The whole construction procedure
 902 can be divided into three steps.

903 (1) Apply KST to reduce dimension, i.e., represent $f \in C([a, b]^d)$ by the compositions and
 904 combinations of univariate continuous functions.

905 (2) Apply Theorem 6 to design sub-networks to approximate the univariate continuous
 906 functions in the previous step within the desired error.

907 (3) Integrate the sub-networks to form the final network and estimate its size.

908 The details of these three steps can be found below.

909 **Step 1:** Apply KST to reduce dimension.

910 To apply KST, we define a linear function $\mathcal{L}_1(t) = (b-a)t + a$ for any $t \in [0, 1]$. Clearly,
 911 \mathcal{L}_1 is a bijection from $[0, 1]$ to $[a, b]$. Define

$$912 \quad \tilde{f}(\mathbf{y}) := f(\mathcal{L}_1(y_1), \mathcal{L}_1(y_2), \dots, \mathcal{L}_1(y_d)) \quad \text{for any } \mathbf{y} = [y_1, y_2, \dots, y_d]^T \in [0, 1]^d.$$

913 Then, $\tilde{f} : [0, 1]^d \rightarrow \mathbb{R}$ is a continuous function since $f \in C([a, b]^d)$. By Theorem 8, there
 914 exists $\tilde{h}_{i,j} \in C([0, 1])$ and $\tilde{g}_i \in C(\mathbb{R})$ for $i = 0, 1, \dots, 2d$ and $j = 1, 2, \dots, d$ such that

$$915 \quad \tilde{f}(\mathbf{y}) = \sum_{i=0}^{2d} \tilde{g}_i \left(\sum_{j=1}^d \tilde{h}_{i,j}(y_j) \right) \quad \text{for any } \mathbf{y} = [y_1, y_2, \dots, y_d]^T \in [0, 1]^d.$$

916 Let $\tilde{\mathcal{L}}_1$ be the inverse of \mathcal{L}_1 , i.e., $\tilde{\mathcal{L}}_1(t) = (t-a)/(b-a)$ for any $t \in [a, b]$. Then, for any
 917 $x_j \in [a, b]$, there exists a unique $y_j \in [0, 1]$ such that $\mathcal{L}_1(y_j) = x_j$ and $y_j = \tilde{\mathcal{L}}_1(x_j)$ for any
 918 $j = 1, 2, \dots, d$, which implies

$$919 \quad \begin{aligned} f(\mathbf{x}) &= f(x_1, x_2, \dots, x_d) = f(\mathcal{L}_1(y_1), \mathcal{L}_1(y_2), \dots, \mathcal{L}_1(y_d)) = \tilde{f}(\mathbf{y}) \\ &= \sum_{i=0}^{2d} \tilde{g}_i \left(\sum_{j=1}^d \tilde{h}_{i,j}(y_j) \right) = \sum_{i=0}^{2d} \tilde{g}_i \left(\sum_{j=1}^d \tilde{h}_{i,j}(\tilde{\mathcal{L}}_1(x_j)) \right) = \sum_{i=0}^{2d} \tilde{g}_i \left(\sum_{j=1}^d \tilde{h}_{i,j} \circ \tilde{\mathcal{L}}_1(x_j) \right). \end{aligned}$$

920 It follows that

$$921 \quad f(\mathbf{x}) = \sum_{i=0}^{2d} \tilde{g}_i \left(\sum_{j=1}^d \tilde{h}_{i,j} \circ \tilde{\mathcal{L}}_1(x_j) \right) = \sum_{i=0}^{2d} \tilde{g}_i \circ \hat{h}_i(\mathbf{x}) \quad \text{for any } \mathbf{x} \in [a, b]^d,$$

922 where

$$923 \quad \hat{h}_i(\mathbf{x}) = \sum_{j=1}^d \tilde{h}_{i,j} \circ \tilde{\mathcal{L}}_1(x_j) \quad \text{for any } \mathbf{x} = [x_1, x_2, \dots, x_d]^T \in [a, b]^d. \quad (10)$$

924 Set

$$925 \quad M = \max_{i \in \{0, 1, \dots, 2d\}} \|\hat{h}_i\|_{L^\infty([a, b]^d)} + 1 > 0.$$

926 Define $\mathcal{L}_2(t) = (t + 2M)/4M$ and $\tilde{\mathcal{L}}_2(t) = 4Mt - 2M$ for any $t \in \mathbb{R}$. Then, \mathcal{L}_2 is a
 927 bijection from $[-M, M]$ to $[\frac{1}{4}, \frac{3}{4}]$ and $\tilde{\mathcal{L}}_2$ is the inverse of \mathcal{L}_2 . Clearly, $\tilde{\mathcal{L}}_2 \circ \mathcal{L}_2(t) = t$ for any
 928 $t \in [-M, M]$, which implies $\hat{h}_i(\mathbf{x}) = \tilde{\mathcal{L}}_2 \circ \mathcal{L}_2 \circ \hat{h}_i(\mathbf{x})$ for any $\mathbf{x} \in [a, b]^d$. Therefore, for any
 929 $\mathbf{x} \in [a, b]^d$, we have

$$930 \quad f(\mathbf{x}) = \sum_{i=0}^{2d} \tilde{g}_i \circ \hat{h}_i(\mathbf{x}) = \sum_{i=0}^{2d} \tilde{g}_i \circ \tilde{\mathcal{L}}_2 \circ \mathcal{L}_2 \circ \hat{h}_i(\mathbf{x}) = \sum_{i=0}^{2d} g_i \circ h_i(\mathbf{x}),$$

931 where

$$932 \quad g_i = \tilde{g}_i \circ \tilde{\mathcal{L}}_2 \quad \text{and} \quad h_i = \mathcal{L}_2 \circ \hat{h}_i \quad \text{for } i = 0, 1, \dots, 2d. \quad (11)$$

933 Clearly, $\mathcal{L}_2(t) \in [\frac{1}{4}, \frac{3}{4}]$ for any $t \in [-M, M]$, which implies

$$934 \quad h_i(\mathbf{x}) = \mathcal{L}_2 \circ \hat{h}_i(\mathbf{x}) \in [\frac{1}{4}, \frac{3}{4}] \quad \text{for any } \mathbf{x} \in [a, b]^d \text{ and } i = 0, 1, \dots, 2d.$$

935 **Step 2:** Design sub-networks to approximate g_i and h_i .

936 Next, we will construct sub-networks to approximate g_i and h_i for each i . Obviously,
 937 $g_i = \tilde{g}_i \circ \tilde{\mathcal{L}}_2$ is continuous on \mathbb{R} and hence uniformly continuous on $[0, 1]$ for each i . Thus,
 938 for $i = 0, 1, \dots, 2d$, there exists $\delta_i > 0$ such that

$$939 \quad |g_i(z_1) - g_i(z_2)| < \varepsilon/(4d + 2) \quad \text{for any } z_1, z_2 \in [0, 1] \text{ with } |z_1 - z_2| < \delta_i.$$

940 Set $\delta = \min(\{\delta_i : i = 0, 1, \dots, 2d\} \cup \{\frac{1}{4}\})$. Then, for $i = 0, 1, \dots, 2d$, we have

$$941 \quad |g_i(z_1) - g_i(z_2)| < \varepsilon/(4d + 2) \quad \text{for any } z_1, z_2 \in [0, 1] \text{ with } |z_1 - z_2| < \delta. \quad (12)$$

942 For each $i \in \{0, 1, \dots, 2d\}$, by Theorem 6, there exists a function ϕ_i generated by an
 943 EUAF network with width 36 and depth 5 such that

$$944 \quad |g_i(z) - \phi_i(z)| < \varepsilon/(4d + 2) \quad \text{for any } z \in [0, 1]. \quad (13)$$

945 Fix $i \in \{0, 1, \dots, 2d\}$, we will design an EUAF network to generate a function $\psi_i : [a, b]^d \rightarrow \mathbb{R}$ satisfying

$$946 \quad |h_i(\mathbf{x}) - \psi_i(\mathbf{x})| < \delta \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

948 For any $\mathbf{x} = [x_1, x_2, \dots, x_d]^T \in [a, b]^d$, by Equations (10) and (11), we have

$$949 \quad \begin{aligned} h_i(\mathbf{x}) &= \mathcal{L}_2 \circ \hat{h}_i(\mathbf{x}) = \mathcal{L}_2 \left(\sum_{j=1}^d \tilde{h}_{i,j} \circ \tilde{\mathcal{L}}_1(x_j) \right) = \frac{\left(\sum_{j=1}^d \tilde{h}_{i,j} \circ \tilde{\mathcal{L}}_1(x_j) \right) + 2M}{4M} \\ &= \sum_{j=1}^d \left(\frac{\tilde{h}_{i,j} \circ \tilde{\mathcal{L}}_1(x_j)}{4M} + \frac{1}{2d} \right) = \sum_{j=1}^d h_{i,j}(x_j), \end{aligned}$$

950 where

$$951 \quad h_{i,j}(t) := \frac{\tilde{h}_{i,j} \circ \tilde{\mathcal{L}}_1(t)}{4M} + \frac{1}{2d} \quad \text{for any } t \in [a, b], i = 0, 1, \dots, 2d, \text{ and } j = 1, 2, \dots, d.$$

952 It is easy to verify that $h_{i,j} \in C([a, b]^d)$ each $i \in \{0, 1, \dots, 2d\}$ and each $j \in \{1, 2, \dots, d\}$,
 953 from which we deduce by Theorem 6 that there exists a function $\psi_{i,j}$ generated by an EUAF
 954 network with width 36 and depth 5 such that

$$955 \quad |h_{i,j}(t) - \psi_{i,j}(t)| < \delta/d \quad \text{for any } t \in [a, b].$$

956 For each $i \in \{0, 1, \dots, 2d\}$, we define

$$957 \quad \psi_i(\mathbf{x}) := \sum_{j=1}^d \psi_{i,j}(x_j) \quad \text{for any } \mathbf{x} = [x_1, x_2, \dots, x_d]^T \in [a, b]^d.$$

958 Then, for any $\mathbf{x} = [x_1, x_2, \dots, x_d]^T \in [a, b]^d$ and each $i \in \{0, 1, \dots, 2d\}$, we have

$$959 \quad |h_i(\mathbf{x}) - \psi_i(\mathbf{x})| = \left| \sum_{j=1}^d h_{i,j}(x_j) - \sum_{j=1}^d \psi_{i,j}(x_j) \right| = \sum_{j=1}^d |h_{i,j}(x_j) - \psi_{i,j}(x_j)| < \sum_{j=1}^d \delta/d = \delta.$$

960 **Step 3:** Integrate sub-networks.

961 Finally, we build an integrated network with the desired size to approximate the target
 962 function f . The desired function ϕ can be defined as

$$963 \quad \phi(\mathbf{x}) := \sum_{i=0}^{2d} \phi_i \circ \psi_i(\mathbf{x}) = \sum_{i=0}^{2d} \phi_i \left(\sum_{j=1}^d \psi_{i,j}(x_j) \right) \quad \text{for any } \mathbf{x} = [x_1, x_2, \dots, x_d]^T \in [a, b]^d.$$

964 Let us first estimate the approximation error and then determine the size of the target
 965 network realizing ϕ . See Figure 10 for an illustration of the target network realizing ϕ for
 966 the case $d = 2$.

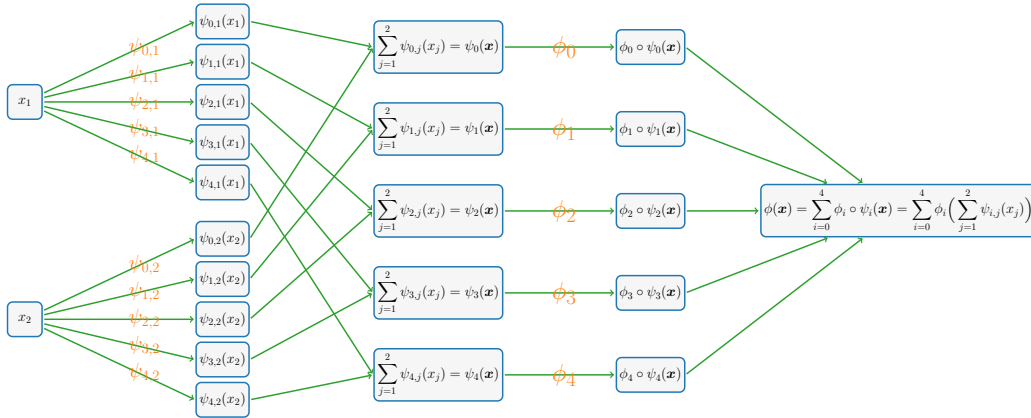


Figure 10: An illustration of the target network realizing ϕ for any $\mathbf{x} \in [a, b]^d$ in the case of $d = 2$. This network contains $(2d + 1)d + (2d + 1) = (d + 1)(2d + 1)$ sub-networks that realize $\psi_{i,j}$ and ϕ_i for $i = 0, 1, \dots, 2d$ and $j = 1, 2, \dots, d$.

967 Fix $\mathbf{x} \in [a, b]^d$ and $i \in \{0, 1, \dots, 2d\}$. Recall that $h_i(\mathbf{x}) \in [\frac{1}{4}, \frac{3}{4}]$ and

$$968 \quad |h_i(\mathbf{x}) - \psi_i(\mathbf{x})| < \delta \leq \frac{1}{4},$$

969 implying $\psi_i(\mathbf{x}) \in [0, 1]$. Then, by Equation (12) (set $z_1 = h_i(\mathbf{x})$ and $z_2 = \psi_i(\mathbf{x})$ therein),
 970 we have

$$971 \quad \left| g_i \circ h_i(\mathbf{x}) - g_i \circ \psi_i(\mathbf{x}) \right| = \left| g_i(h_i(\mathbf{x})) - g_i(\psi_i(\mathbf{x})) \right| < \varepsilon/(4d + 2).$$

972 By Equation (13) (set $z = \psi_i(\mathbf{x}) \in [0, 1]$ therein), we have

$$973 \quad \left| g_i \circ \psi_i(\mathbf{x}) - \phi_i \circ \psi_i(\mathbf{x}) \right| = \left| g_i(\psi_i(\mathbf{x})) - \phi_i(\psi_i(\mathbf{x})) \right| < \varepsilon/(4d + 2).$$

974 Therefore, for any $\mathbf{x} \in [a, b]^d$, we have

$$\begin{aligned} |f(\mathbf{x}) - \phi(\mathbf{x})| &= \left| \sum_{i=0}^{2d} g_i \circ h_i(\mathbf{x}) - \sum_{i=0}^{2d} \phi_i \circ \psi_i(\mathbf{x}) \right| = \sum_{i=0}^{2d} \left| g_i \circ h_i(\mathbf{x}) - \phi_i \circ \psi_i(\mathbf{x}) \right| \\ 975 \quad &\leq \sum_{i=0}^{2d} \left(\left| g_i \circ h_i(\mathbf{x}) - g_i \circ \psi_i(\mathbf{x}) \right| + \left| g_i \circ \psi_i(\mathbf{x}) - \phi_i \circ \psi_i(\mathbf{x}) \right| \right) \\ &< \sum_{i=0}^{2d} \left(\varepsilon/(4d + 2) + \varepsilon/(4d + 2) \right) = \varepsilon. \end{aligned}$$

976 It remains to show ϕ can be generated by an EUAF network with the desired size. Recall
 977 that, for each $i \in \{0, 1, \dots, 2d\}$ and each $j \in \{1, 2, \dots, d\}$, $\psi_{i,j}$ can be generated by an EUAF
 978 network with width 36, depth 5, and therefore at most

$$979 \quad (1 \times 36 + 36) + (36 \times 36 + 36) \times 4 + (36 \times 1 + 1) = 5437$$

980 nonzero parameters. Hence, for each $i \in \{0, 1, \dots, 2d\}$, ψ_i , given by $\psi_i(\mathbf{x}) = \sum_{j=1}^d \psi_{i,j}(x_j)$,
 981 can be generated by an EUAF network with width $36d$, depth 5, and at most $5437d$ nonzero
 982 parameters.

983 Since $\psi_i(\mathbf{x}) \in [0, 1]$ for any $\mathbf{x} \in [a, b]^d$ and $i = 0, 1, \dots, 2d$, we have $\sigma(\psi_i(\mathbf{x})) = \psi_i(\mathbf{x})$
 984 for any $\mathbf{x} \in [a, b]^d$. Hence, $\phi_i \circ \psi_i$ can be generated by an EUAF network as visualized in
 985 Figure 11.

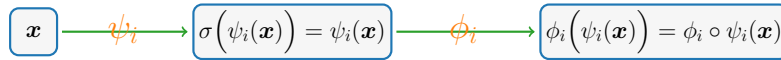


Figure 11: An illustration of the target EUAF network generating $\phi_i \circ \psi_i(\mathbf{x})$ for any $\mathbf{x} \in [a, b]^d$ and $i = 0, 1, \dots, 2d$.

986 Recall that ϕ_i can be generated by an EUAF network with width 36 and depth 5.
 987 Hence, the network generating ϕ_i has at most 5437 nonzero parameters. As we can see from
 988 Figure 11, $\phi_i \circ \psi_i$ can be generated by an EUAF network with width $\max\{36d, 36\} = 36d$,
 989 depth $5 + 1 + 5 = 11$, and at most $5437d + 5437 = 5437(d + 1)$ nonzero parameters. This
 990 means $\phi = \sum_{i=0}^{2d} \phi_i \circ \psi_i$ can be generated by an EUAF network with width $36d(2d + 1)$,
 991 depth 11, and therefore at most $5437(d + 1)(2d + 1)$ nonzero parameters as desired. So we
 992 finish the proof.

993 **5.2 Proof of Theorem 4**

994 The proof of Theorem 4 relies on a basic result of real analysis given in the following lemma.

995 **Lemma 13.** *Suppose $A, B \subseteq \mathbb{R}^d$ are two disjoint bounded closed sets. Then, there exists*
 996 *a continuous function $f \in C(\mathbb{R}^d)$ such that $f(\mathbf{x}) = 1$ for any $\mathbf{x} \in A$ and $f(\mathbf{y}) = 0$ for any*
 997 *$\mathbf{y} \in B$.*

998 *Proof.* Define $\text{dist}(\mathbf{x}, A) = \inf\{\|\mathbf{x} - \mathbf{y}\|_2 : \mathbf{y} \in A\}$ and $\text{dist}(\mathbf{x}, B) = \inf\{\|\mathbf{x} - \mathbf{y}\|_2 : \mathbf{y} \in B\}$
 999 for any $\mathbf{x} \in \mathbb{R}^d$. It is easy to verify that $\text{dist}(\mathbf{x}, A)$ and $\text{dist}(\mathbf{x}, B)$ are continuous in $\mathbf{x} \in \mathbb{R}^d$.
 1000 Since $A, B \subseteq \mathbb{R}^d$ are two disjoint bounded closed subsets, we have $\text{dist}(\mathbf{x}, A) + \text{dist}(\mathbf{x}, B) > 0$
 1001 for any $\mathbf{x} \in \mathbb{R}^d$. Finally, define

$$1002 \quad f(\mathbf{x}) := \frac{\text{dist}(\mathbf{x}, B)}{\text{dist}(\mathbf{x}, A) + \text{dist}(\mathbf{x}, B)} \quad \text{for any } \mathbf{x} \in \mathbb{R}^d.$$

1003 Then f meets the requirements. So we finish the proof. ■

1004 With Lemma 13, we can prove Theorem 4.

1005 *Proof of Theorem 4.* For any $f = \sum_{j=1}^J r_j \cdot \mathbf{1}_{E_j} \in \mathcal{C}_d(E_1, E_2, \dots, E_J)$, our goal is to construct
 1006 a function ϕ generated by a σ -activated network such that $\phi(\mathbf{x}) = f(\mathbf{x})$ for any $\mathbf{x} \in$
 1007 $\bigcup_{j=1}^J E_j$, where E_1, E_2, \dots, E_J are pairwise disjoint bounded closed subsets of \mathbb{R}^d . Define
 1008 $E := \bigcup_{j=1}^J E_j$ and choose $a, b \in \mathbb{R}$ properly such that $E \subseteq [a, b]^d$.

1009 For each $j \in \{1, 2, \dots, J\}$, E_j and $\tilde{E}_j := E \setminus E_j$ are two disjoint bounded closed subsets.
 1010 Then, for each j , by Lemma 13, there exists $g_j \in C(\mathbb{R}^d)$ such that $g_j(\mathbf{x}) = 1$ for any $\mathbf{x} \in E_j$
 1011 and $g_j(\mathbf{y}) = 0$ for any $\mathbf{y} \in \tilde{E}_j = E \setminus E_j$. By defining $g := \sum_{j=1}^J r_j \cdot g_j \in C(\mathbb{R}^d)$, we have
 1012 $g(\mathbf{x}) = \sum_{j=1}^J r_j \cdot \mathbf{1}_{E_j}(\mathbf{x}) = f(\mathbf{x})$ for any $\mathbf{x} \in E = \bigcup_{j=1}^J E_j$.

1013 Since r_1, r_2, \dots, r_J are rational numbers and $g : [a, b]^d \rightarrow \mathbb{R}$ is continuous, there exist
 1014 $n_1, n_2 \in \mathbb{Z} \setminus \{0\}$ such that

- 1015 • $n_1 \cdot r_j + n_2 \in \mathbb{N}^+$ for $j = 1, 2, \dots, J$;
- 1016 • $n_1 \cdot g(\mathbf{x}) + n_2 \geq 0$ for any $\mathbf{x} \in [a, b]^d$.

1017 By applying Theorem 1 to $2(n_1 \cdot g + n_2) + 1 \in C([a, b]^d)$, there exists a function ϕ_1
 1018 generated by an EUAF network with width $36d(2d + 1)$, depth 11, and at most $5437(d +$
 1019 $1)(2d + 1)$ nonzero parameters such that

$$1020 \quad \left| 2(n_1 \cdot g(\mathbf{x}) + n_2) + 1 - \phi_1(\mathbf{x}) \right| \leq 1/2 \quad \text{for any } \mathbf{x} \in [a, b]^d. \quad (14)$$

1021 It follows that

$$1022 \quad \left| 2\left(n_1 \cdot \sum_{j=1}^J r_j \cdot \mathbf{1}_{E_j}(\mathbf{x}) + n_2\right) + 1 - \phi_1(\mathbf{x}) \right| \leq 1/2 \quad \text{for any } \mathbf{x} \in E = \bigcup_{j=1}^J E_j.$$

1023 Since E_1, E_2, \dots, E_J are pairwise disjoint, we have

$$1024 \quad \left| 2(n_1 \cdot r_j + n_2) + 1 - \phi_1(\mathbf{x}) \right| \leq 1/2 \quad \text{for any } \mathbf{x} \in E_j \text{ and each } j \in \{1, 2, \dots, J\}. \quad (15)$$

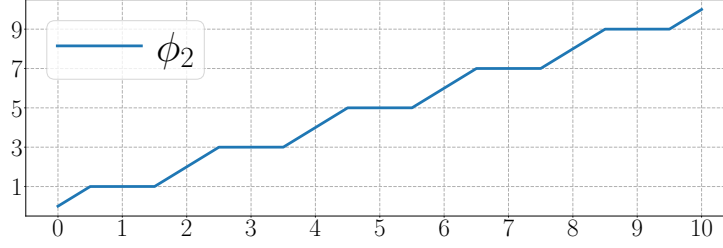


Figure 12: An illustration of ϕ_2 on $[0, 10]$.

1025 Define $\phi_2(x) = x + 1/2 - \sigma(x + 3/2)$ for any $x \in \mathbb{R}$. See Figure 12 for an illustration.

1026 It is easy to verify that

$$1027 \quad \phi_2(y) = 2k + 1 \quad \text{for any } y \text{ and } k \in \mathbb{N}^+ \text{ with } |2k + 1 - y| \leq 1/2. \quad (16)$$

1028 Then, by Equations (15) and (16) (set $y = \phi_1(\mathbf{x})$ and $k = n_1 \cdot r_j + n_2$ therein), we have

$$1029 \quad \phi_2 \circ \phi_1(\mathbf{x}) = \phi_2(y) = 2k + 1 = 2(n_1 \cdot r_j + n_2) + 1$$

1030 for any $\mathbf{x} \in E_j$ and any $j \in \{1, 2, \dots, J\}$, which implies

$$1031 \quad \frac{\phi_2 \circ \phi_1(\mathbf{x}) - 2n_2 - 1}{2n_1} = r_j \quad \text{for any } \mathbf{x} \in E_j \text{ and any } j \in \{1, 2, \dots, J\}.$$

1032 Define

$$1033 \quad \phi(\mathbf{x}) := \frac{\phi_2 \circ \phi_1(\mathbf{x}) - 2n_2 - 1}{2n_1} \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

1034 Clearly, we have $\phi(\mathbf{x}) = r_j$ for any $\mathbf{x} \in E_j$ and each $j \in \{1, 2, \dots, J\}$, which implies

$$1035 \quad \phi(\mathbf{x}) = \sum_{j=1}^J r_j \cdot \mathbb{1}_{E_j}(\mathbf{x}) = f(\mathbf{x}) \quad \text{for any } \mathbf{x} \in E = \bigcup_{j=1}^J E_j.$$

1036 It remains to show that ϕ can be generated by an EUAF network with the desired size.

1037 Set $M = 2\|n_1 g + n_2\|_{L^\infty([a, b]^d)} + 3/2 > 0$. By Equation (14) and the fact $n_1 \cdot g(\mathbf{x}) + n_2 \geq 0$

1038 for any $\mathbf{x} \in [a, b]^d$, we have

$$1039 \quad \phi_1(\mathbf{x}) \in \left[1/2, 2\|n_1 g + n_2\|_{L^\infty([a, b]^d)} + 1 + 1/2\right] \subseteq [0, M] \quad \text{for any } \mathbf{x} \in [a, b]^d.$$

1040 Then, for any $\mathbf{x} \in [a, b]^d$, we have

$$1041 \quad \begin{aligned} \phi_2 \circ \phi_1(\mathbf{x}) &= \phi_1(\mathbf{x}) + 1/2 - \sigma(\phi_1(\mathbf{x}) + 3/2) \\ &= M\sigma(\phi_1(\mathbf{x})/M) + 1/2 - \sigma(\phi_1(\mathbf{x}) + 3/2). \end{aligned}$$

1042 It follows that

$$1043 \quad \phi(\mathbf{x}) = \frac{\phi_2 \circ \phi_1(\mathbf{x}) - 2n_2 - 1}{2n_1} = \frac{M\sigma(\phi_1(\mathbf{x})/M) - \sigma(\phi_1(\mathbf{x}) + 3/2) - 2n_2 - 1/2}{2n_1},$$

1044 for any $\mathbf{x} \in [a, b]^d$. That means the network realizing ϕ has just one more hidden layer with
 1045 2 neurons, compared to the network realizing ϕ_1 . Recall that ϕ_1 can be generated by an
 1046 EUAF network with width $36d(2d+1)$, depth 11, and at most $5437(d+1)(2d+1)$ nonzero
 1047 parameters. Therefore, ϕ , limited on $[a, b]^d$, can be generated by an EUAF network with
 1048 width $36d(2d+1)$, depth 12, and at most

$$1049 \quad 5437(d+1)(2d+1) + \underbrace{2 \times 36d(2d+1) + 2 + 2 \times 1 + 1}_{\text{all possible new parameters}} \leq 5509(d+1)(2d+1)$$

1050 nonzero parameters. So we finish the proof. ■

1051 6. Proof of Theorem 6

1052 To prove Theorem 6, we need to introduce two auxiliary theorems, Theorems 14 and 15,
 1053 which serve as two important intermediate steps.

1054 **Theorem 14.** *Let $f \in C([0, 1])$ be a continuous function. Given any $\varepsilon > 0$, if K is a*
 1055 *positive integer satisfying*

$$1056 \quad |f(x_1) - f(x_2)| < \varepsilon/2 \quad \text{for any } x_1, x_2 \in [0, 1] \text{ with } |x_1 - x_2| < 1/K, \quad (17)$$

1057 *then there exists a function ϕ generated by an EUAF network with width 2 and depth 3 such*
 1058 *that $\|\phi\|_{L^\infty([0,1])} \leq \|f\|_{L^\infty([0,1])} + 1$ and*

$$1059 \quad |\phi(x) - f(x)| < \varepsilon \quad \text{for any } x \in \bigcup_{k=0}^{K-1} \left[\frac{2k}{2K}, \frac{2k+1}{2K} \right].$$

1060 **Theorem 15.** *Let $f \in C([0, 1])$ be a continuous function. Then, for any $\varepsilon > 0$, there exists*
 1061 *a function ϕ generated by an EUAF network with width 36 and depth 5 such that*

$$1062 \quad |\phi(x) - f(x)| < \varepsilon \quad \text{for any } x \in [0, \frac{9}{10}].$$

1063 To prove Theorem 14, we only need to care about the approximation on one “half” of
 1064 $[0, 1]$. It is not necessary to care about the approximation on the other “half” of $[0, 1]$. Such
 1065 an idea is similar to the “trifling region” in (Lu et al., 2021; Zhang, 2020). As we shall
 1066 see later, the proof of Theorem 14 can eventually be converted to a point-fitting problem,
 1067 which can be solved by applying Proposition 7.

1068 The key idea to prove Theorem 15 is to apply Theorem 14 to several horizontally shifted
 1069 variants of the target function. Then a good approximation can be constructed via the
 1070 combinations and multiplications of these variants, similar to the idea of (Lu et al., 2021;
 1071 Zhang, 2020) to obtain an error estimation with the L^∞ -norm from a result with the L^p -
 1072 norm for $p \in [1, \infty)$.

1073 The proofs of Theorems 14 and 15 will be presented in Sections 6.1 and 6.2, respectively.
 1074 Let us first prove Theorem 6 by assuming Theorem 15 is true.

1075 *Proof of Theorem 6.* Define a linear function \mathcal{L} by $\mathcal{L}(x) = a + \frac{10(b-a)}{9}x$ for any $x \in [0, \frac{9}{10}]$.
 1076 Then \mathcal{L} is a bijection from $[0, \frac{9}{10}]$ to $[a, b]$. It follows that $f \circ \mathcal{L}$ is a continuous function

1077 on $[0, \frac{9}{10}]$. By Theorem 15, there exists a function $\tilde{\phi}$ generated by an EUAF network with
 1078 width 36 and depth 5 such that

1079
$$|f \circ \mathcal{L}(x) - \tilde{\phi}(x)| < \varepsilon \quad \text{for any } x \in [0, \frac{9}{10}].$$

1080 Define $\tilde{\mathcal{L}}(y) = \frac{9(y-a)}{10(b-a)}$ for any $y \in [a, b]$. Clearly, it is the inverse of \mathcal{L} , i.e., $\mathcal{L} \circ \tilde{\mathcal{L}}(y) = y$
 1081 for any $y \in [a, b]$. Therefore, for any $y \in [a, b]$, we have $x = \tilde{\mathcal{L}}(y) \in [0, \frac{9}{10}]$, which implies

1082
$$\begin{aligned} |f(y) - \tilde{\phi} \circ \tilde{\mathcal{L}}(y)| &= |f \circ \mathcal{L} \circ \tilde{\mathcal{L}}(y) - \tilde{\phi} \circ \tilde{\mathcal{L}}(y)| \\ &= |f \circ \mathcal{L}(\tilde{\mathcal{L}}(y)) - \tilde{\phi}(\tilde{\mathcal{L}}(y))| = |f \circ \mathcal{L}(x) - \tilde{\phi}(x)| < \varepsilon. \end{aligned}$$

1083 By defining $\phi := \tilde{\phi} \circ \tilde{\mathcal{L}}$, we have $|f(y) - \phi(y)| < \varepsilon$ for any $y \in [a, b]$ as desired.

1084 Note that $\tilde{\phi}$ can be realized by an EUAF network with width 36 and depth 5. We can
 1085 compose $\tilde{\mathcal{L}}$ and the affine linear map of the network $\tilde{\phi}$ that connects the input layer and
 1086 the first hidden layer. Therefore, $\phi = \tilde{\phi} \circ \tilde{\mathcal{L}}$ can also be realized by an EUAF network with
 1087 width 36 and depth 5. So we finish the proof. ■

1088 6.1 Proof of Theorem 14

1089 Partition $[0, 1]$ into $2K$ small intervals \mathcal{I}_k and $\tilde{\mathcal{I}}_k$ for $k = 1, 2, \dots, K$, i.e.,

1090
$$\mathcal{I}_k = [\frac{2k-2}{2K}, \frac{2k-1}{2K}] \quad \text{and} \quad \tilde{\mathcal{I}}_k = [\frac{2k-1}{2K}, \frac{2k}{2K}].$$

1091 Clearly, $[0, 1] = \bigcup_{k=1}^K (\mathcal{I}_k \cup \tilde{\mathcal{I}}_k)$. Let x_k be the right endpoint of \mathcal{I}_k , i.e., $x_k = \frac{2k-1}{2K}$ for
 1092 $k = 1, 2, \dots, K$. See an illustration of \mathcal{I}_k , $\tilde{\mathcal{I}}_k$, and x_k in Figure 13 for the case $K = 5$.

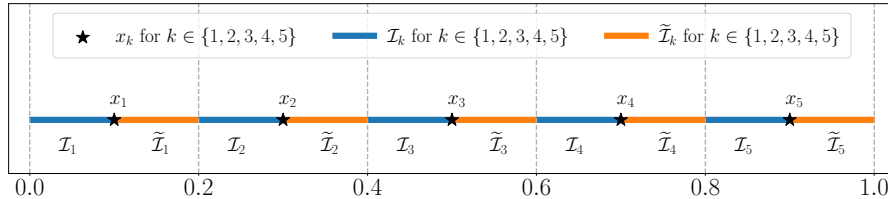


Figure 13: An illustration of \mathcal{I}_k and $\tilde{\mathcal{I}}_k$ for $k \in \{1, 2, \dots, K\}$ with $K = 5$.

1093 Our goal is to construct a function ϕ generated by an EUAF network with the desired
 1094 size to approximate f well on \mathcal{I}_k for $k = 1, 2, \dots, K$. It is not necessary to care about the
 1095 values of ϕ on $\tilde{\mathcal{I}}_k$ for all k . In other words, we only need to care about the approximation
 1096 on one “half” of $[0, 1]$, which is the key for our proof.

1097 Define $\psi(x) := x - \sigma(x)$ for any $x \in \mathbb{R}$, where σ is defined in Equation (1). See Figure 14
 1098 for an illustration of ψ .

1099 It is easy to verify that

1100
$$\psi(y) = 2k - 2 \quad \text{for any } y \in [2k - 2, 2k - 1] \text{ and each } k \in \{1, 2, \dots, K\}.$$

1101 It follows that

1102
$$\psi(2Kx)/2 + 1 = k \quad \text{for any } x \in [\frac{2k-2}{2K}, \frac{2k-1}{2K}] = \mathcal{I}_k \text{ and each } k \in \{1, 2, \dots, K\}. \quad (18)$$

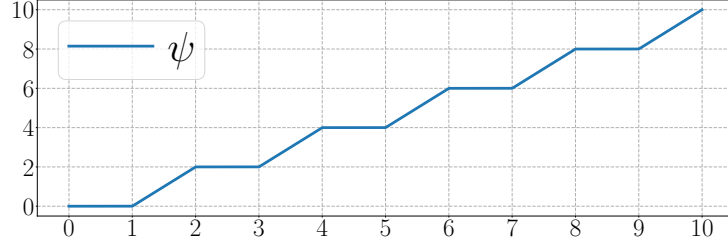


Figure 14: An illustration of ψ on $[0, 10]$.

1103 Recall that x_k is the right endpoint of \mathcal{I}_k for $k = 1, 2, \dots, K$. Set $M = \|f\|_{L^\infty([0,1])} + 1$
 1104 and define

$$1105 \quad \xi_k := \frac{f(x_k) + M}{2M} \in [0, 1] \quad \text{for } k = 1, 2, \dots, K.$$

1106 Then $[\xi_1, \xi_2, \dots, \xi_K]^T$ is in $[0, 1]^K$. By Proposition 7, there exists $w_0 \in \mathbb{R}$ such that

$$1107 \quad \left| \sigma_1\left(\frac{w_0}{\pi+k}\right) - \xi_k \right| < \varepsilon/(4M) \quad \text{for } k = 1, 2, \dots, K.$$

1108 Let m_0 be an integer larger than $|w_0|$, e.g., set $m_0 = \lfloor |w_0| \rfloor + 1$. It is easy to verify that

$$1109 \quad \frac{w_0}{\pi+k} + 2m_0 \geq 0 \quad \text{for any } x \in [0, 1].$$

1110 Since $\sigma(x) = \sigma_1(x)$ for any $x \geq 0$ and σ_1 is periodic with period 2, we have

$$1111 \quad \left| \sigma\left(\frac{w_0}{\pi+k} + 2m_0\right) - \xi_k \right| = \left| \sigma_1\left(\frac{w_0}{\pi+k} + 2m_0\right) - \xi_k \right| = \left| \sigma_1\left(\frac{w_0}{\pi+k}\right) - \xi_k \right| < \varepsilon/(4M),$$

1112 for $k = 1, 2, \dots, K$. It follows that

$$1113 \quad \begin{aligned} \left| 2M\sigma\left(\frac{w_0}{\pi+k} + 2m_0\right) - M - f(x_k) \right| &= \left| 2M\sigma\left(\frac{w_0}{\pi+k} + 2m_0\right) - M - (2M\xi_k - M) \right| \\ &= 2M \left| \sigma\left(\frac{w_0}{\pi+k} + 2m_0\right) - \xi_k \right| < 2M \cdot \frac{\varepsilon}{4M} = \varepsilon/2, \end{aligned} \quad (19)$$

1114 for $k = 1, 2, \dots, K$.

1115 The desired ϕ is defined as

$$1116 \quad \phi(x) := 2M\sigma\left(\frac{w_0}{\pi+\psi(2Kx)/2+1} + 2m_0\right) - M \quad \text{for any } x \in [0, 1].$$

1117 Recall that $m_0 \geq |w_0|$ and $\psi(x) \geq 0$ for any $x \geq 0$, which implies

$$1118 \quad \frac{w_0}{\pi+\psi(2Kx)/2+1} + 2m_0 \geq 0 \quad \text{for any } x \in [0, 1].$$

1119 It follows that $\|\phi\|_{L^\infty([0,1])} \leq M = \|f\|_{L^\infty([0,1])} + 1$ since $0 \leq \sigma(y) \leq 1$ for any $y \geq 0$.

1120 For any $x \in \mathcal{I}_k$ and each $k \in \{1, 2, \dots, K\}$, by Equation (18), we have $\psi(2Kx)/2+1 = k$,
 1121 which implies

$$1122 \quad \phi(x) = 2M\sigma\left(\frac{w_0}{\pi+\psi(2Kx)/2+1} + 2m_0\right) - M = 2M\sigma\left(\frac{w_0}{\pi+k} + 2m_0\right) - M.$$

1123 Clearly, for any $x \in \mathcal{I}_k$ and each $k \in \{1, 2, \dots, K\}$, we have $|x_k - x| < 1/K$. Then, by
 1124 Equation (17), we get

$$1125 \quad |f(x_k) - f(x)| < \varepsilon/2 \quad \text{for any } x \in \mathcal{I}_k \text{ and each } k \in \{1, 2, \dots, K\}.$$

1126 Therefore, by Equation (19), we have

$$\begin{aligned}
1127 \quad |\phi(x) - f(x)| &= \left| 2M\sigma\left(\frac{w_0}{\pi+k} + 2m_0\right) - M - f(x) \right| \\
&\leq \left| 2M\sigma\left(\frac{w_0}{\pi+k} + 2m_0\right) - M - f(x_k) \right| + |f(x_k) - f(x)| < \varepsilon/2 + \varepsilon/2 = \varepsilon
\end{aligned}$$

1128 for any $x \in \mathcal{I}_k$ and each $k \in \{1, 2, \dots, K\}$. It follows that

$$1129 \quad |\phi(x) - f(x)| < \varepsilon \quad \text{for any } x \in \bigcup_{j=1}^K \mathcal{I}_j = \bigcup_{j=1}^K \left[\frac{2j-2}{2K}, \frac{2j-1}{2K} \right] = \bigcup_{k=0}^{K-1} \left[\frac{2k}{2K}, \frac{2k+1}{2K} \right].$$

1130 It remains to show that ϕ can be generated by an EUAF network with the desired size.
1131 Observe that

$$1132 \quad \sigma(y) + 1 = \frac{y}{|y| + 1} + 1 = \frac{y}{-y + 1} + 1 = \frac{1}{-y + 1} \quad \text{for any } y \leq 0.$$

1133 By setting $y = -\pi - \psi(2Kx)/2 \leq 0$ for any $x \in [0, 1]$, we have

$$\begin{aligned}
1134 \quad \frac{1}{\pi + \psi(2Kx)/2 + 1} &= \frac{1}{-y + 1} = \sigma(y) + 1 = \sigma(-\pi - \psi(2Kx)/2) + 1 \\
&= \sigma\left(-\pi - (2Kx - \sigma(2Kx))/2\right) + 1 \\
&= \sigma\left(-\pi - Kx + \sigma(2Kx)/2\right) + 1,
\end{aligned}$$

1135 where the second-to-last equality comes from $\psi(z) = z - \sigma(z)$ for any $z \in \mathbb{R}$. Therefore, we
1136 get

$$\begin{aligned}
1137 \quad \phi(x) &= 2M\sigma\left(\frac{w_0}{\pi + \psi(2Kx)/2 + 1} + 2m_0\right) - M \\
&= 2M\sigma\left(w_0\sigma\left(-\pi - Kx + \sigma(2Kx)/2\right) + w_0 + 2m_0\right) - M.
\end{aligned} \tag{20}$$

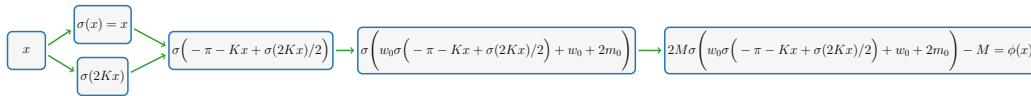


Figure 15: An illustration of the target EUAF network realizing $\phi(x)$ for $x \in [0, 1]$ based on Equation (20).

1138 Thus, the desired EUAF network realizing ϕ is shown in Figure 15. Clearly, the network
1139 in Figure 15 has width 2 and depth 3 as desired. It is easy to verify that the network
1140 architecture corresponding ϕ is independent of the target function f and the desired error
1141 ε . That is, we can fix the architecture and only adjust parameters to achieve the desired
1142 approximation error. So we finish the proof.

1143 **6.2 Proof of Theorem 15**

1144 The key idea of proving Theorem 15 is to apply Theorem 14 to several horizontally shifted
 1145 variants of the target function. Then a good approximation can be expected via combina-
 1146 tions and multiplications of these variants. Thus, we need to reproduce $f(x, y) = xy$ locally
 1147 via an EUAF network as shown in the following lemma.

1148 **Lemma 16.** *For any $M > 0$, there exists a function ϕ generated by an EUAF network with*
 1149 *width 9 and depth 2 such that*

1150
$$\phi(x, y) = xy \quad \text{for any } x, y \in [-M, M].$$

1151 The proof of this lemma is given in Section 6.3. Now let us first prove Theorem 15 by
 1152 assuming this lemma is true.

1153 *Proof of Theorem 15.* Set $\tilde{\varepsilon} = \varepsilon/4$ and extend f from $[0, 1]$ to $[-1, 1]$ by defining $f(x) = f(0)$
 1154 for any $x \in [-1, 0)$. Then f is continuous on $[-1, 1]$ and therefore uniformly continuous.
 1155 Thus, there exists $K = K(f, \varepsilon) \in \mathbb{N}^+$ with $K \geq 10$ such that

1156
$$|f(x_1) - f(x_2)| < \tilde{\varepsilon}/2 \quad \text{for any } x_1, x_2 \in [-1, 1] \text{ with } |x_1 - x_2| < 1/K.$$

1157 For $i = 1, 2, 3, 4$, define

1158
$$f_i(x) := f\left(x - \frac{i}{4K}\right) \quad \text{for any } x \in [0, 1].$$

1159 For each $i \in \{1, 2, 3, 4\}$ and any $x_1, x_2 \in [0, 1]$ with $|x_1 - x_2| < 1/K$, we have $x_1 -$
 1160 $\frac{i}{4K}, x_2 - \frac{i}{4K} \in [-1, 1]$ and $\left|(x_1 - \frac{i}{4K}) - (x_2 - \frac{i}{4K})\right| = |x_1 - x_2| < 1/K$, which implies

1161
$$|f_i(x_1) - f_i(x_2)| = \left|f\left(x_1 - \frac{i}{4K}\right) - f\left(x_2 - \frac{i}{4K}\right)\right| < \tilde{\varepsilon}/2.$$

1162 That is, for $i = 1, 2, 3, 4$, we have

1163
$$|f_i(x_1) - f_i(x_2)| < \tilde{\varepsilon}/2 \quad \text{for any } x_1, x_2 \in [0, 1] \text{ with } |x_1 - x_2| < 1/K,$$

1164 which means we can apply Theorem 14 to $f_i \in C([0, 1])$. For each $i \in \{1, 2, 3, 4\}$, by
 1165 Theorem 14, there exists a function ϕ_i generated by an EUAF network with width 2 and
 1166 depth 3 such that

1167
$$\|\phi_i\|_{L^\infty([0,1])} \leq \|f_i\|_{L^\infty([0,1])} + 1 \leq \|f\|_{L^\infty([-1,1])} + 1$$

1168 and

1169
$$|\phi_i(x) - f_i(x)| < \tilde{\varepsilon} = \varepsilon/4 \quad \text{for any } x \in \bigcup_{k=0}^{K-1} \left[\frac{2k}{2K}, \frac{2k+1}{2K}\right].$$

1170 Define

1171
$$\psi(x) = \sigma(x + 1 - \sigma(x + 1)) \quad \text{for any } x \in \mathbb{R}.$$

1172 See an illustration of ψ on $[0, 2K]$ for $K = 5$ in Figure 16.

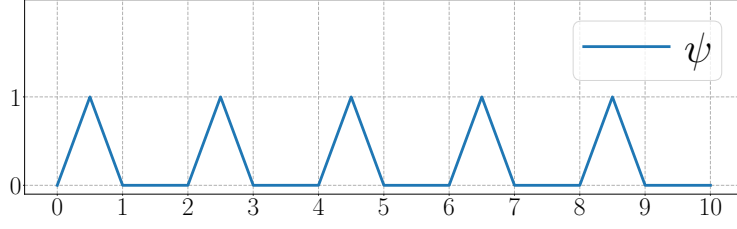


Figure 16: An illustration of ψ on $[0, 2K]$ for $K = 5$.

1173 Clearly, $0 \leq \psi(2Kx) \leq 1$ for any $x \in [0, 1]$, from which we deduce

$$1174 \quad \left| (\phi_i(x) - f_i(x))\psi(2Kx) \right| \leq |\phi_i(x) - f_i(x)| < \varepsilon/4 \quad \text{for any } x \in \bigcup_{k=0}^{K-1} \left[\frac{2k}{2K}, \frac{2k+1}{2K} \right].$$

1175 Observe that $\psi(y) = 0$ for $y \in \bigcup_{k=0}^{K-1} [2k+1, 2k+2]$, which implies

$$1176 \quad \psi(2Kx) = 0 \quad \text{for any } x \in \bigcup_{k=0}^{K-1} \left[\frac{2k+1}{2K}, \frac{2k+2}{2K} \right] \supseteq [0, 1] \setminus \bigcup_{k=0}^{K-1} \left[\frac{2k}{2K}, \frac{2k+1}{2K} \right].$$

1177 It follows that

$$1178 \quad \left| (\phi_i(x) - f_i(x))\psi(2Kx) \right| < \varepsilon/4 \quad \text{for any } x \in [0, 1] \text{ and } i = 1, 2, 3, 4. \quad (21)$$

1179 For each $i \in \{1, 2, 3, 4\}$ and any $z \in [0, \frac{9}{10}] \subseteq [0, 1 - \frac{1}{K}] \subseteq [0, 1 - \frac{i}{4K}]$, we have

$$1180 \quad y_i = z + \frac{i}{4K} \in [\frac{i}{4K}, 1] \subseteq [0, 1].$$

1181 Therefore, by bringing $x = y_i \in [0, 1]$ into Equation (21), we have

$$1182 \quad \begin{aligned} \varepsilon/4 &> \left| (\phi_i(y_i) - f_i(y_i))\psi(2Ky_i) \right| = \left| \phi_i(y_i)\psi(2Ky_i) - f_i(y_i)\psi(2Ky_i) \right| \\ &= \left| \phi_i\left(z + \frac{i}{4K}\right)\psi\left(2K\left(z + \frac{i}{4K}\right)\right) - f_i\left(z + \frac{i}{4K}\right)\psi\left(2K\left(z + \frac{i}{4K}\right)\right) \right| \\ &= \left| \phi_i\left(z + \frac{i}{4K}\right)\psi\left(2Kz + \frac{i}{2}\right) - f\left(z\right)\psi\left(2Kz + \frac{i}{2}\right) \right| \end{aligned} \quad (22)$$

1183 for any $z \in [0, \frac{9}{10}]$, where the last equality comes from the fact that $f_i(x) = f(x - \frac{i}{4K})$ for
1184 any $x \in [0, 1] \supseteq [\frac{i}{4K}, 1]$. The desired ϕ is defined as

$$1185 \quad \phi(x) := \sum_{i=1}^4 \phi_i\left(x + \frac{i}{4K}\right)\psi\left(2Kx + \frac{i}{2}\right) \quad \text{for any } x \in [0, \frac{9}{10}].$$

1186 It is easy to verify that $\sum_{i=1}^4 \psi\left(x + \frac{i}{2}\right) = 1$ for any $x \geq 0$ based on the definition of ψ .
1187 See Figure 17 for illustrations. It follows that $\sum_{i=1}^4 \psi\left(2Kz + \frac{i}{2}\right) = 1$ for any $z \in [0, \frac{9}{10}]$.

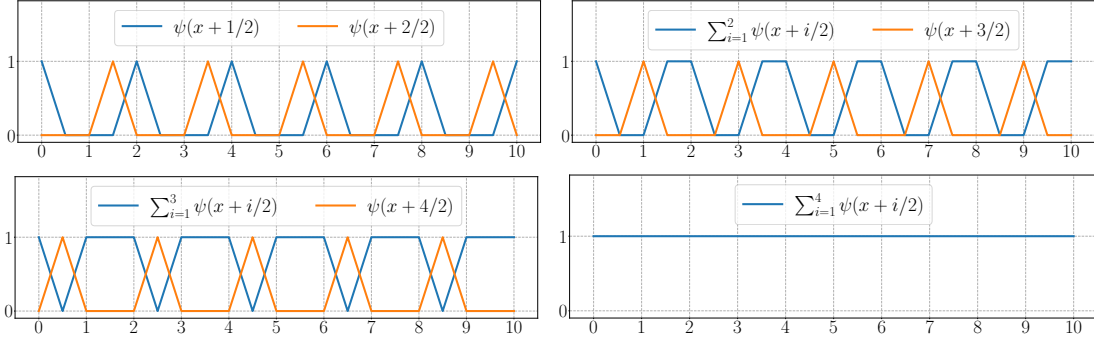


Figure 17: Illustrations of $\sum_{i=1}^4 \psi(x + i/2) = 1$ for any $x \in [0, 10]$.

1188 Hence, for any $z \in [0, \frac{9}{10}]$, by Equation (22), we have

$$\begin{aligned}
 1189 \quad |\phi(z) - f(z)| &= \left| \sum_{i=1}^4 \phi_i(z + \frac{i}{4K}) \psi(2Kz + \frac{i}{2}) - f(z) \sum_{i=1}^4 \psi(2Kz + \frac{i}{2}) \right| \\
 &\leq \sum_{i=1}^4 \left| \phi_i(z + \frac{i}{4K}) \psi(2Kz + \frac{i}{2}) - f(z) \psi(2Kz + \frac{i}{2}) \right| < 4 \cdot \frac{\varepsilon}{4} = \varepsilon.
 \end{aligned}$$

1190 That is, $|\phi(x) - f(x)| < \varepsilon$ for any $x \in [0, \frac{9}{10}]$ as desired. It remains to show that ϕ , limited
 1191 on $[0, \frac{9}{10}]$, can be generated by an EUAF network with the desired size.

1192 Note that $x + 1 = (2K + 1)\sigma(\frac{x+1}{2K+1})$ for any $x \in [0, 2K]$, which implies

$$1193 \quad \psi(x) = \sigma(x + 1 - \sigma(x + 1)) = \sigma\left((2K + 1)\sigma\left(\frac{x+1}{2K+1}\right) - \sigma(x + 1)\right).$$

1194 This means ψ , limited on $[0, 2K]$, can be generated by an EUAF network with width 2 and
 1195 depth 2. Since $0 \leq 2Kx + \frac{i}{2} \leq 2K \cdot \frac{9}{10} + 2 = 2K(\frac{9}{10} + \frac{1}{K}) \leq 2K$ for any $x \in [0, \frac{9}{10}]$, $\psi(2K \cdot + \frac{i}{2})$,
 1196 limited on $[0, \frac{9}{10}]$, can also be generated by an EUAF network with width 2 and depth 2.

1197 Note that ϕ_i , limited on $[0, 1]$, can also be generated by an EUAF network with width 2
 1198 and depth 3. Clearly, $x + \frac{i}{4K} \in [0, 1]$ for any $x \in [0, \frac{9}{10}]$, and, therefore, $\phi_i(\cdot + \frac{i}{4K})$, limited
 1199 on $[0, \frac{9}{10}]$, can also be generated by an EUAF network with width 2 and depth 3.

1200 Recall that $\|\phi_i\|_{L^\infty([0,1])} \leq \|f\|_{L^\infty([-1,1])} + 1 =: M$. Thus, $|\phi_i(x + \frac{i}{4K})| \leq M$ and
 1201 $|\psi(2Kx + \frac{i}{2})| \leq 1 \leq M$ for any $x \in [0, \frac{9}{10}]$ and $i = 1, 2, 3, 4$. By Lemma 16, there exists a
 1202 function Γ generated by an EUAF network with width 9 and depth 2 such that

$$1203 \quad \Gamma(x, y) = xy \quad \text{for any } x, y \in [-M, M].$$

1204 It follows that

$$1205 \quad \Gamma\left(\phi_i(x + \frac{i}{4K}), \psi(2Kx + \frac{i}{2})\right) = \phi_i(x + \frac{i}{4K})\psi(2Kx + \frac{i}{2}) \quad \text{for } i = 1, 2, 3, 4.$$

1206 Therefore, each component of $\phi(x)$, $\phi_i(x + \frac{i}{4K})\psi(2Kx + \frac{i}{2})$ for each $i \in \{1, 2, 3, 4\}$, can
 1207 be generated by the network in Figure 18 for any $x \in [0, \frac{9}{10}]$. Clearly, such a network has
 1208 width 9 and depth 6. Since the 4-th hidden layer of the network in Figure 18 uses the

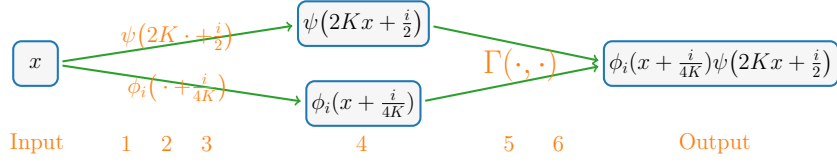


Figure 18: An illustration of the target EUAF network realizing each component of $\phi(x)$, $\phi_i(x + \frac{i}{4K})\psi(2Kx + \frac{i}{2})$, for any $x \in [0, \frac{9}{10}]$ and each $i \in \{1, 2, 3, 4\}$. The networks realizing $\phi_i(\cdot + \frac{i}{4K})$ and $\psi(2K\cdot + \frac{i}{2})$ can be placed in parallel since we can manually add a hidden layers to ψ since $\sigma \circ \psi(2Kx + \frac{i}{2}) = \psi(2Kx + \frac{i}{2})$ for any $x \in [0, \frac{9}{10}]$.

1209 identity map as an activation function for each neuron in this hidden layer, we can reduce
 1210 the depth by 1 via composing two adjacent affine linear maps to generate a new one. Thus,
 1211 the network in Figure 18 can be interpreted as an EUAF network with width 9 and depth
 1212 5.

1213 Note that ϕ is the sum of its four components, namely,

$$\phi(x) = \sum_{i=1}^4 \phi_i(x + \frac{i}{4K})\psi(2Kx + \frac{i}{2}) \quad \text{for any } x \in [0, \frac{9}{10}].$$

1215 Therefore, ϕ , limited on $[0, \frac{9}{10}]$, can be generated by an EUAF network with width $9 \times 4 = 36$
 1216 and depth 5 as desired. It is easy to verify that the designed network architecture is
 1217 independent of the target function f and the desired error ε . That is, we can fix the
 1218 architecture and only adjust parameters to achieve an arbitrarily small approximation error.
 1219 So we finish the proof. ■

1220 6.3 Proof of Lemma 16

1221 The key idea of proving Lemma 16 is the polarization identity $2xy = (x + y)^2 - x^2 - y^2$.
 1222 Thus, we need to reproduce x^2 locally by an EUAF network as shown in the following
 1223 lemma.

1224 **Lemma 17.** *There exists a function ϕ generated by an EUAF network with width 3 and*
 1225 *depth 2 such that*

$$1226 \quad \phi(x) = x^2 \quad \text{for any } x \in [-1, 1].$$

1227 *Proof.* Observe that

$$1228 \quad \sigma(y) + 1 = \frac{y}{|y| + 1} + 1 = \frac{y}{-y + 1} + 1 = \frac{1}{-y + 1} \quad \text{for any } y \leq 0.$$

1229 For any $x \in [-1, 1]$, we have $-x - 1 \leq 0$ and $-x - 2 \leq 0$, which implies

$$\begin{aligned} 1230 \quad \sigma(-x - 1) - \sigma(-x - 2) &= (\sigma(-x - 1) + 1) - (\sigma(-x - 2) + 1) \\ &= \frac{1}{-(-x - 1) + 1} - \frac{1}{-(-x - 2) + 1} \\ &= \frac{1}{x + 2} - \frac{1}{x + 3} = \frac{1}{(x + 2)(x + 3)}. \end{aligned}$$

1231 It follows from $1 - \frac{12}{(x+2)(x+3)} \leq 0$ for any $x \in [-1, 1]$ that

1232
$$\sigma\left(1 - \frac{12}{(x+2)(x+3)}\right) + 1 = \frac{1}{-\left(1 - \frac{12}{(x+2)(x+3)}\right) + 1} = \frac{x^2 + 5x + 6}{12},$$

1233 implying

1234
$$\begin{aligned} x^2 &= 12\sigma\left(1 - \frac{12}{(x+2)(x+3)}\right) + 12 - (5x + 6) \\ &= 12\sigma\left(1 - 12(\sigma(-x-1) - \sigma(-x-2))\right) + 11\frac{6-5x}{11} \\ &= 12\sigma\left(1 - 12\sigma(-x-1) + 12\sigma(-x-2)\right) + 11\sigma\left(\frac{6-5x}{11}\right) =: \phi(x), \end{aligned}$$

1235 where the equality $\frac{6-5x}{11} = \sigma\left(\frac{6-5x}{11}\right)$ comes from two facts: $\frac{6-5x}{11} \in [0, 1]$ since $x \in [-1, 1]$
 1236 and $\sigma(z) = z$ for any $z \in [0, 1]$.

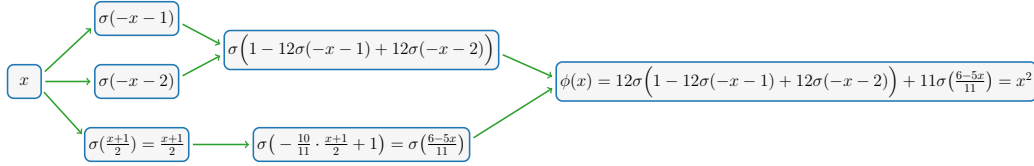


Figure 19: An illustration of the target EUAF network realizing $\phi(x) = x^2$ for $x \in [-1, 1]$.

1237 Then, x^2 can be generated by the network shown in Figure 19 for any $x \in [-1, 1]$. The
 1238 target network has width 3 and depth 2. So we finish the proof. ■

1239 With Lemma 17 at hand, we are ready to prove Lemma 16.

1240 *Proof of Lemma 16.* By Lemma 17, there exists a function $\tilde{\phi}$ generated by an EUAF net-
 1241 work such that $\tilde{\phi}(t) = t^2$ for any $t \in [-1, 1]$. Then, for any $x, y \in [-M, M]$, we have

1242
$$\begin{aligned} xy &= 2M^2\left(\left(\frac{x+y}{2M}\right)^2 - \left(\frac{x}{2M}\right)^2 - \left(\frac{y}{2M}\right)^2\right) \\ &= 2M^2\left(\tilde{\phi}\left(\frac{x+y}{2M}\right) - \tilde{\phi}\left(\frac{x}{2M}\right) - \tilde{\phi}\left(\frac{y}{2M}\right)\right) =: \phi(x, y). \end{aligned}$$

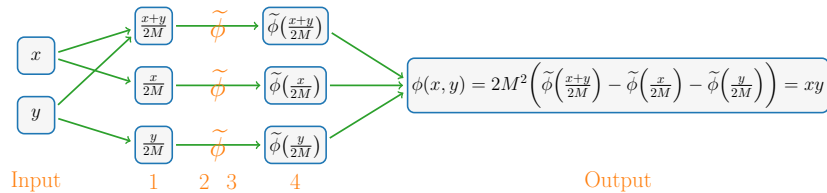


Figure 20: An illustration of the target network realizing $\phi(x) = xy$ for $x, y \in [-M, M]$.

1243 The target network realizing ϕ with width 9 and depth 4 is shown in Figure 20. Note
 1244 that we can reduce the depth by one if the activation function of each neuron in a hidden

1245 layer is the identity map. In fact, we can eliminate this hidden layer by composing two
1246 adjacent affine linear maps to generate a new one. The 1-st and 4-th hidden layers of the
1247 network in Figure 20 use the identity map as an activation function for each neuron. Thus,
1248 the network in Figure 20 can be interpreted as an EUAF network with width 9 and depth
1249 2. So we finish the proof. ■

1250 7. Proof of Proposition 7

1251 We will prove Proposition 7 in this section. The proof includes two main steps. First,
1252 we show how to simply generate a set of rationally independent numbers in Lemma 18
1253 below. Next, we prove that the target point set via a winding of the generated rationally
1254 independent numbers is dense in a hypercube. Such a proof relies on the fact that an
1255 irrational winding on the torus is dense (e.g., see Lemma 2 of (Yarotsky, 2021)) as shown
1256 in Lemma 19 below.

1257 **Lemma 18.** *Given any $K \in \mathbb{N}^+$, any transcendental number $\alpha \in \mathbb{R} \setminus \mathbb{A}$, and any pairwise
1258 distinct rational numbers $r_1, r_2, \dots, r_K \in \mathbb{Q}$, the set of numbers*

$$1259 \quad \left\{ \frac{1}{\alpha + r_k} : k = 1, 2, \dots, K \right\}$$

1260 *are rationally independent.*

1261 **Lemma 19.** *Given any rationally independent numbers a_1, a_2, \dots, a_K for any $K \in \mathbb{N}^+$ and
1262 an arbitrary periodic function $g : \mathbb{R} \rightarrow \mathbb{R}$ with period T , i.e., $g(x+T) = g(x)$ for any $x \in \mathbb{R}$,
1263 assume there exist $x_1, x_2 \in \mathbb{R}$ with $0 < x_2 - x_1 < T$ such that g is continuous on $[x_1, x_2]$.
1264 Then the following set*

$$1265 \quad \left\{ [g(wa_1), g(wa_2), \dots, g(wa_K)]^T : w \in \mathbb{R} \right\}$$

1266 *is dense in $[M_1, M_2]^K$, where $M_1 = \min_{x \in [x_1, x_2]} g(x)$ and $M_2 = \max_{x \in [x_1, x_2]} g(x)$.*

1267 The proofs of these two lemmas can be found in Sections 7.1 and 7.2, respectively.
1268 With these two lemmas at hand, the proof of Proposition 7 is straightforward. In fact,
1269 we can prove a more general result in Proposition 20 below, which implies Proposition 7
1270 immediately.

1271 **Proposition 20.** *Given an arbitrary periodic function $g : \mathbb{R} \rightarrow \mathbb{R}$ with period T , i.e.,
1272 $g(x+T) = g(x)$ for any $x \in \mathbb{R}$, assume there exist $x_1, x_2 \in \mathbb{R}$ with $0 < x_2 - x_1 < T$
1273 such that g is continuous on $[x_1, x_2]$. Then, for any $K \in \mathbb{N}^+$, any transcendental number
1274 $\alpha \in \mathbb{R} \setminus \mathbb{A}$, and any pairwise distinct rational numbers $r_1, r_2, \dots, r_K \in \mathbb{Q}$, the following set*

$$1275 \quad \left\{ \left[g\left(\frac{w}{\alpha + r_1}\right), g\left(\frac{w}{\alpha + r_2}\right), \dots, g\left(\frac{w}{\alpha + r_K}\right) \right]^T : w \in \mathbb{R} \right\}$$

1276 *is dense in $[M_1, M_2]^K$, where $M_1 = \min_{x \in [x_1, x_2]} g(x)$ and $M_2 = \max_{x \in [x_1, x_2]} g(x)$. In the case of
1277 $M_1 < M_2$, the following set*

$$1278 \quad \left\{ \left[u \cdot g\left(\frac{w}{\alpha + r_1}\right) + v, u \cdot g\left(\frac{w}{\alpha + r_2}\right) + v, \dots, u \cdot g\left(\frac{w}{\alpha + r_K}\right) + v \right]^T : u, v, w \in \mathbb{R} \right\}$$

1279 *is dense in \mathbb{R}^K .*

1280 Clearly, Proposition 7 is a special case of Proposition 20 with $g = \sigma_1$, $\alpha = \pi$, $r_k = k$ for
 1281 $k = 1, 2, \dots, K$. The transcendence of π is well known (e.g., see the Lindemann-Weierstrass
 1282 Theorem). By setting $x_1 = 0$ and $x_2 = 1$, we have $[M_1, M_2] = [0, 1]$ and σ_1 is continuous
 1283 on $[0, 1]$, which means that the following set

$$1284 \quad \left\{ \left[\sigma_1\left(\frac{w}{\pi+1}\right), \sigma_1\left(\frac{w}{\pi+2}\right), \dots, \sigma_1\left(\frac{w}{\alpha+K}\right) \right]^T : w \in \mathbb{R} \right\}$$

1285 is dense in $[0, 1]^K$ as desired.

1286 Finally, let us prove Proposition 20 by assuming Lemmas 18 and 19 are true.

1287 *Proof of Proposition 20.* By Lemma 18, the set of numbers

$$1288 \quad \left\{ \frac{1}{\alpha+r_k} : k = 1, 2, \dots, K \right\}$$

1289 are rationally independent. Denote $a_k = \frac{1}{\alpha+r_k}$ for $k = 1, 2, \dots, K$. Then, by Lemma 19,

$$1290 \quad \begin{aligned} & \left\{ \left[g(wa_1), g(wa_2), \dots, g(wa_K) \right]^T : w \in \mathbb{R} \right\} \\ &= \left\{ \left[g\left(\frac{w}{\alpha+r_1}\right), g\left(\frac{w}{\alpha+r_2}\right), \dots, g\left(\frac{w}{\alpha+r_K}\right) \right]^T : w \in \mathbb{R} \right\} \end{aligned}$$

1291 is dense in $[M_1, M_2]^K$.

1292 Next, let us consider the case $M_1 < M_2$ for the latter result. For any $\varepsilon > 0$ and any
 1293 $\mathbf{x} \in \mathbb{R}^K$, by setting $J = \|\mathbf{x}\|_\infty + 1 > 0$, we have $\frac{\mathbf{x}+J}{2J} \in [0, 1]^K$, and hence

$$1294 \quad \mathbf{y} := \frac{\mathbf{x}+J}{2J}(M_2 - M_1) + M_1 \in [M_1, M_2]^K.$$

1295 By the former result, there exists $w_0 \in \mathbb{R}$ such that

$$1296 \quad \left\| \mathbf{y} - \left[g\left(\frac{w_0}{\alpha+r_1}\right), g\left(\frac{w_0}{\alpha+r_2}\right), \dots, g\left(\frac{w_0}{\alpha+r_K}\right) \right]^T \right\|_\infty < \frac{M_2-M_1}{2J}\varepsilon$$

1297 It follows from $\mathbf{y} = \frac{\mathbf{x}+J}{2J}(M_2 - M_1) + M_1$ that

$$1298 \quad \mathbf{x} = \frac{2J}{M_2-M_1}\mathbf{y} + \frac{J(M_1+M_2)}{M_1-M_2} =: u_0\mathbf{y} + v_0,$$

1299 where $u_0 = \frac{2J}{M_2-M_1}$ and $v_0 = \frac{J(M_1+M_2)}{M_1-M_2}$. Therefore,

$$1300 \quad \begin{aligned} & \left\| \mathbf{x} - \left[u_0g\left(\frac{w_0}{\alpha+r_1}\right) + v_0, u_0g\left(\frac{w_0}{\alpha+r_2}\right) + v_0, \dots, u_0g\left(\frac{w_0}{\alpha+r_K}\right) + v_0 \right]^T \right\|_\infty \\ &= \left\| u_0\mathbf{y} + v_0 - \left[u_0g\left(\frac{w_0}{\alpha+r_1}\right) + v_0, u_0g\left(\frac{w_0}{\alpha+r_2}\right) + v_0, \dots, u_0g\left(\frac{w_0}{\alpha+r_K}\right) + v_0 \right]^T \right\|_\infty \\ &< u_0 \frac{M_2-M_1}{2J}\varepsilon = \frac{2J}{M_2-M_1} \frac{M_2-M_1}{2J}\varepsilon = \varepsilon. \end{aligned}$$

1301 Since $\varepsilon > 0$ and $\mathbf{x} \in \mathbb{R}^K$ are arbitrary, the following set

$$1302 \quad \left\{ \left[u \cdot g\left(\frac{w}{\alpha+r_1}\right) + v, u \cdot g\left(\frac{w}{\alpha+r_2}\right) + v, \dots, u \cdot g\left(\frac{w}{\alpha+r_K}\right) + v \right]^T : u, v, w \in \mathbb{R} \right\}$$

1303 is dense in \mathbb{R}^K . So we finish the proof. ■

1304 **7.1 Proof of Lemma 18**

1305 Before proving Lemma 18, let us first briefly discuss related concepts. Recall that a complex
 1306 number α is an algebraic number if and only if there exist $\lambda_0, \lambda_1, \dots, \lambda_J \in \mathbb{Q}$ with
 1307 $\sum_{j=0}^J \lambda_j \alpha^j = 0$. The set of all algebraic numbers is denoted by \mathbb{A} . We say a complex number
 1308 is **transcendental** if it is not in \mathbb{A} . Almost all complex numbers are transcendental
 1309 since the set \mathbb{A} is countable. The best known transcendental numbers are π (the ratio of a
 1310 circle's circumference to its diameter) and e (the natural logarithmic base).

1311 In order to prove Lemma 18, we need an auxiliary lemma below, characterizing some
 1312 properties of coefficients of Lagrange basis polynomials. Recall that, for any given pairwise
 1313 distinct numbers $x_1, x_2, \dots, x_K \in \mathbb{R}$, the Lagrange basis polynomials are

$$1314 \quad p_k(x) := \prod_{\substack{j \in \{1, 2, \dots, K\} \\ j \neq k}} \frac{x - x_j}{x_k - x_j} = \frac{x - x_1}{x_k - x_1} \dots \frac{x - x_{k-1}}{x_k - x_{k-1}} \frac{x - x_{k+1}}{x_k - x_{k+1}} \dots \frac{x - x_K}{x_k - x_K} \quad (23)$$

1315 for $k = 1, 2, \dots, K$. They are polynomials of degree $\leq K - 1$, which means we can represent
 1316 each p_k by

$$1317 \quad p_k(x) = \sum_{j=1}^K a_{k,j} x^{j-1} = a_{k,1} + a_{k,2}x + \dots + a_{k,K}x^{K-1}$$

1318 for $k = 1, 2, \dots, K$ and any $x \in \mathbb{R}$. Thus, the coefficients of these K Lagrange basis
 1319 polynomials p_1, p_2, \dots, p_K form a matrix

$$1320 \quad \mathbf{A} = (a_{i,j}) = \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,K} \\ a_{2,1} & a_{2,2} & \dots & a_{2,K} \\ \vdots & \vdots & \ddots & \vdots \\ a_{K,1} & a_{K,2} & \dots & a_{K,K} \end{bmatrix} \in \mathbb{R}^{K \times K}. \quad (24)$$

1321 The lemma below essentially characterizes the linear independence of Lagrange basis
 1322 polynomials.

1323 **Lemma 21.** *With the same setting just above, the matrix \mathbf{A} given in Equation (24) is*
 1324 *invertible.*

1325 *Proof.* For any $\mathbf{y} = [y_1, y_2, \dots, y_K] \in \mathbb{R}^K$, by the definition of Lagrange basis polyno-
 1326 mials $p_k(x)$ for $k = 1, 2, \dots, K$ in Equation (23), $p(x) = \sum_{k=1}^K y_k p_k(x)$ is the target in-
 1327 terpolation polynomial for sample points $(x_1, y_1), (x_2, y_2), \dots, (x_K, y_K)$. That is, for any
 1328 $\ell \in \{1, 2, \dots, K\}$, we have

$$1329 \quad \begin{aligned} y_\ell = p(x_\ell) &= \sum_{k=1}^K y_k p_k(x_\ell) = \sum_{k=1}^K y_k \sum_{j=1}^K a_{k,j} x_\ell^{j-1} \\ &= [y_1, y_2, \dots, y_K] \cdot \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,K} \\ a_{2,1} & a_{2,2} & \dots & a_{2,K} \\ \vdots & \vdots & \ddots & \vdots \\ a_{K,1} & a_{K,2} & \dots & a_{K,K} \end{bmatrix} \cdot \begin{bmatrix} x_\ell^0 \\ x_\ell^1 \\ \vdots \\ x_\ell^{K-1} \end{bmatrix} = \mathbf{y}^T \mathbf{A} \begin{bmatrix} x_\ell^0 \\ x_\ell^1 \\ \vdots \\ x_\ell^{K-1} \end{bmatrix}. \end{aligned}$$

1330 It follows that

$$1331 \quad \mathbf{y}^T = [y_1, y_2, \dots, y_K] = \mathbf{y}^T \mathbf{A} \begin{bmatrix} x_1^0 & x_2^0 & \cdots & x_K^0 \\ x_1^1 & x_2^1 & \cdots & x_K^1 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{K-1} & x_2^{K-1} & \cdots & x_K^{K-1} \end{bmatrix}.$$

1332 Since $\mathbf{y} \in \mathbb{R}^K$ is arbitrary, we have

$$1333 \quad \mathbf{A} \begin{bmatrix} x_1^0 & x_2^0 & \cdots & x_K^0 \\ x_1^1 & x_2^1 & \cdots & x_K^1 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{K-1} & x_2^{K-1} & \cdots & x_K^{K-1} \end{bmatrix} = \mathbf{I}_K,$$

1334 where $\mathbf{I}_K \in \mathbb{R}^{K \times K}$ is the identity matrix. Recall that x_1, x_2, \dots, x_K are pairwise distinct,
1335 which implies the Vandermonde matrix

$$1336 \quad \begin{bmatrix} x_1^0 & x_2^0 & \cdots & x_K^0 \\ x_1^1 & x_2^1 & \cdots & x_K^1 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{K-1} & x_2^{K-1} & \cdots & x_K^{K-1} \end{bmatrix}$$

1337 is invertible. Thus, \mathbf{A} is also invertible. So we complete the proof. ■

1338 With Lemma 21 at hand, we are ready to prove Lemma 18.

1339 *Proof of Lemma 18.* Let $x_k = -r_k \in \mathbb{Q}$ for $k = 1, 2, \dots, K$ and define the Lagrange basis
1340 polynomials as

$$1341 \quad p_k(x) := \prod_{\substack{j \in \{1, 2, \dots, K\} \\ j \neq k}} \frac{x - x_j}{x_k - x_j} = w_k \prod_{\substack{j \in \{1, 2, \dots, K\} \\ j \neq k}} (x - x_j),$$

1342 where

$$1343 \quad w_k = \prod_{\substack{j \in \{1, 2, \dots, K\} \\ j \neq k}} \frac{1}{x_k - x_j} \neq 0 \quad \text{for } k = 1, 2, \dots, K.$$

1344 It follows from $x_k \in \mathbb{Q}$ that w_k is rational and nonzero, i.e., $w_k \in \mathbb{Q}/\{0\}$ for any k . Clearly,
1345 each p_k is a polynomial of degree $\leq K - 1$. That means we can represent p_k by

$$1346 \quad p_k(x) = \sum_{j=1}^K a_{k,j} x^{j-1} = a_{k,1} + a_{k,2}x + \cdots + a_{k,K}x^{K-1}$$

1347 for $k = 1, 2, \dots, K$ and any $x \in \mathbb{R}$, where each coefficient $a_{k,j}$ is rational. Therefore, the
1348 coefficients of p_1, p_2, \dots, p_K form a matrix

$$1349 \quad \mathbf{A} = (a_{i,j}) = \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,K} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,K} \\ \vdots & \vdots & \ddots & \vdots \\ a_{K,1} & a_{K,2} & \cdots & a_{K,K} \end{bmatrix} \in \mathbb{Q}^{K \times K}.$$

1350 Now assume there exist rational numbers $\lambda_1, \lambda_2, \dots, \lambda_K \in \mathbb{Q}$ such that $\sum_{k=1}^K \lambda_k \cdot \frac{1}{\alpha + r_k} =$
 1351 0. Our goal is to prove $\lambda_1 = \lambda_2 = \dots = \lambda_K = 0$. Clearly, we have

$$\begin{aligned}
 0 &= \sum_{k=1}^K \frac{\lambda_k}{\alpha + r_k} = \sum_{k=1}^K \frac{\lambda_k}{\alpha - x_k} = \prod_{j=1}^K (\alpha - x_j) \cdot \underbrace{\sum_{k=1}^K \frac{\lambda_k}{\alpha - x_k}}_{=0} = \sum_{k=1}^K \frac{\lambda_k}{w_k} \cdot w_k \prod_{\substack{j \in \{1, 2, \dots, K\} \\ j \neq k}} (\alpha - x_j) \\
 &= \sum_{k=1}^K \frac{\lambda_k}{w_k} \cdot p_k(\alpha) = \sum_{k=1}^K \frac{\lambda_k}{w_k} \sum_{j=1}^K a_{k,j} \alpha^{j-1} = \sum_{j=1}^K \left(\underbrace{\sum_{k=1}^K \frac{\lambda_k}{w_k} a_{k,j}}_{=0 \text{ since } \alpha \in \mathbb{R} \setminus \mathbb{A}} \right) \cdot \alpha^{j-1}.
 \end{aligned}$$

1353 For any $k, j \in \{1, 2, \dots, K\}$, we have $\lambda_k, w_k, a_{k,j} \in \mathbb{Q}$, implying $\sum_{k=1}^K \frac{\lambda_k}{w_k} a_{k,j} \in \mathbb{Q}$. Since
 1354 $\alpha \in \mathbb{R} \setminus \mathbb{A}$ is a transcendental number, the coefficients must be 0, i.e.,

$$\sum_{k=1}^K \frac{\lambda_k}{w_k} a_{k,j} = 0 \quad \text{for } j = 1, 2, \dots, K.$$

1356 It follows that

$$\mathbf{0} = \begin{bmatrix} \frac{\lambda_1}{w_1} & \frac{\lambda_2}{w_2} & \dots & \frac{\lambda_K}{w_K} \end{bmatrix} \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,K} \\ a_{2,1} & a_{2,2} & \dots & a_{2,K} \\ \vdots & \vdots & \ddots & \vdots \\ a_{K,1} & a_{K,2} & \dots & a_{K,K} \end{bmatrix} = \begin{bmatrix} \frac{\lambda_1}{w_1} & \frac{\lambda_2}{w_2} & \dots & \frac{\lambda_K}{w_K} \end{bmatrix} \mathbf{A}.$$

1358 By Lemma 21, \mathbf{A} is invertible. Thus, $\begin{bmatrix} \frac{\lambda_1}{w_1} & \frac{\lambda_2}{w_2} & \dots & \frac{\lambda_K}{w_K} \end{bmatrix} = \mathbf{0}$, which implies $\lambda_1 = \lambda_2 = \dots =$
 1359 $\lambda_K = 0$. Hence, the set of numbers $\left\{ \frac{1}{\alpha + r_k} : k = 1, 2, \dots, K \right\}$ are rationally independent,
 1360 which means we finish the proof. ■

1361 7.2 Proof of Lemma 19

1362 The proof of Lemma 19 is mainly based on the fact that an irrational winding is dense on
 1363 the torus (e.g., see Lemma 2 of (Yarotsky, 2021)). For completeness, we establish a lemma
 1364 below and give its detailed proof.

1365 **Lemma 22.** *Given any $K \in \mathbb{N}^+$ and an arbitrary set of rationally independent numbers*
 1366 *$\{a_k : k = 1, 2, \dots, K\} \subseteq \mathbb{R}$, the following set*

$$\left\{ \left[\tau(wa_1), \tau(wa_2), \dots, \tau(wa_K) \right]^T : w \in \mathbb{R} \right\} \subseteq [0, 1)^K$$

1368 *is dense in $[0, 1]^K$, where $\tau(x) := x - \lfloor x \rfloor$ for any $x \in \mathbb{R}$.*

1369 The proof of Lemma 22 can be found later in this section. Now let us first prove
 1370 Lemma 19 by assuming Lemma 22 is true.

1371 *Proof of Lemma 19.* Define $\tilde{g}(x) := g(Tx)$ for any $x \in \mathbb{R}$. Clearly, \tilde{g} is periodic with period
 1372 1 since g is periodic with period T . The continuity of g on $[x_1, x_2]$ implies \tilde{g} is continuous

1373 on $[\frac{x_1}{T}, \frac{x_2}{T}]$ and therefore uniformly continuous on $[\frac{x_1}{T}, \frac{x_2}{T}]$. For any $\varepsilon > 0$, there exists
 1374 $\delta \in (0, \frac{x_2 - x_1}{T})$ such that

$$1375 \quad |\tilde{g}(u) - \tilde{g}(v)| < \varepsilon \quad \text{for any } u, v \in [\frac{x_1}{T}, \frac{x_2}{T}] \text{ with } |u - v| < \delta. \quad (25)$$

1376 Given any $\boldsymbol{\xi} = [\xi_1, \xi_2, \dots, \xi_K] \in [M_1, M_2]^K$, by the extreme value theorem and the
 1377 intermediate value theorem, there exists $z_1, z_2, \dots, z_K \in [x_1, x_2]$ such that

$$1378 \quad g(z_k) = \xi_k \quad \text{for any } k = 1, 2, \dots, K. \quad (26)$$

1379 For $k = 1, 2, \dots, K$, set $y_k = z_k/T \in [\frac{x_1}{T}, \frac{x_2}{T}]$ and

$$1380 \quad \tilde{y}_k = y_k + \frac{\delta}{2} \cdot \mathbf{1}_{\{y_k \leq \frac{x_1}{T} + \frac{\delta}{2}\}} - \frac{\delta}{2} \cdot \mathbf{1}_{\{y_k \geq \frac{x_2}{T} - \frac{\delta}{2}\}}.$$

1381 Then, for $k = 1, 2, \dots, K$, we have

$$1382 \quad \tilde{y}_k = y_k + \frac{\delta}{2} \cdot \mathbf{1}_{\{y_k \leq \frac{x_1}{T} + \frac{\delta}{2}\}} - \frac{\delta}{2} \cdot \mathbf{1}_{\{y_k \geq \frac{x_2}{T} - \frac{\delta}{2}\}} \in [\frac{x_1}{T} + \frac{\delta}{2}, \frac{x_2}{T} - \frac{\delta}{2}]$$

1383 and

$$1384 \quad |\tilde{y}_k - y_k| \leq \left| \frac{\delta}{2} \cdot \mathbf{1}_{\{y_k \leq \frac{x_1}{T} + \frac{\delta}{2}\}} - \frac{\delta}{2} \cdot \mathbf{1}_{\{y_k \geq \frac{x_2}{T} - \frac{\delta}{2}\}} \right| \leq \delta/2.$$

1385 Define $\tau(x) := x - \lfloor x \rfloor$ for any $x \in \mathbb{R}$. Clearly, $[\tau(\tilde{y}_1), \tau(\tilde{y}_2), \dots, \tau(\tilde{y}_K)]^T \in [0, 1]^K$. Then,
 1386 by Lemma 22, there exists $w_0 \in \mathbb{R}$ such that

$$1387 \quad |\tau(w_0 a_k) - \tau(\tilde{y}_k)| < \delta/2 \quad \text{for } k = 1, 2, \dots, K.$$

1388 It follows that

$$1389 \quad \left| \tau(w_0 a_k) + \lfloor \tilde{y}_k \rfloor - \tilde{y}_k \right| = \left| \tau(w_0 a_k) - (\tilde{y}_k - \lfloor \tilde{y}_k \rfloor) \right| = |\tau(w_0 a_k) - \tau(\tilde{y}_k)| < \delta/2$$

1390 for $k = 1, 2, \dots, K$. Since $\tilde{y}_k \in [\frac{x_1}{T} + \frac{\delta}{2}, \frac{x_2}{T} - \frac{\delta}{2}]$, we have $\tau(w_0 a_k) + \lfloor \tilde{y}_k \rfloor \in [\frac{x_1}{T}, \frac{x_2}{T}]$. Besides,

$$1391 \quad \left| \tau(w_0 a_k) + \lfloor \tilde{y}_k \rfloor - y_k \right| \leq \left| \tau(w_0 a_k) + \lfloor \tilde{y}_k \rfloor - \tilde{y}_k \right| + |\tilde{y}_k - y_k| < \delta/2 + \delta/2 = \delta$$

1392 for $k = 1, 2, \dots, K$. Then, by Equation (25), we have

$$1393 \quad \left| \tilde{g}(\tau(w_0 a_k) + \lfloor \tilde{y}_k \rfloor) - \tilde{g}(y_k) \right| < \varepsilon \quad \text{for } k = 1, 2, \dots, K.$$

1394 Recall that \tilde{g} is periodic with period 1, from which we deduce

$$1395 \quad \tilde{g}(\tau(w_0 a_k) + \lfloor \tilde{y}_k \rfloor) = \tilde{g}(w_0 a_k - \lfloor w_0 a_k \rfloor + \lfloor \tilde{y}_k \rfloor) = \tilde{g}(w_0 a_k) = g(T \cdot w_0 a_k)$$

1396 for $k = 1, 2, \dots, K$. Also, we have

$$1397 \quad \tilde{g}(y_k) = g(T y_k) = g(z_k) = \xi_k \quad \text{for } k = 1, 2, \dots, K,$$

1398 where the last equality comes from Equation (26). It follows that

$$1399 \quad |g(T \cdot w_0 a_k) - \xi_k| = \left| \tilde{g}(\tau(w_0 a_k) + \lfloor \tilde{y}_k \rfloor) - \tilde{g}(y_k) \right| < \varepsilon \quad \text{for } k = 1, 2, \dots, K.$$

1400 That is

$$1401 \quad \left\| [g(w_1 a_1), g(w_1 a_2), \dots, g(w_1 a_K)]^T - \boldsymbol{\xi} \right\|_\infty < \varepsilon,$$

1402 where $w_1 = T \cdot w_0 \in \mathbb{R}$. Since $\boldsymbol{\xi} \in [M_1, M_2]^K$ and $\varepsilon > 0$ are arbitrary, the following set

$$1403 \quad \left\{ [g(w a_1), g(w a_2), \dots, g(w a_K)]^T : w \in \mathbb{R} \right\}$$

1404 is dense in $[M_1, M_2]^K$ as desired. So we finish the proof. ■

1405 Finally, let us present the detailed proof of Lemma 22.

1406 *Proof of Lemma 22.* We prove this lemma by mathematical induction. First, we consider
 1407 the case $K = 1$. Note that $a_1 \neq 0$ since it is rationally independent. Thus, we have
 1408 $\{\tau(w a_1) : w \in \mathbb{R}\} = [0, 1)$, which implies $\{\tau(w a_1) : w \in \mathbb{R}\}$ is dense in $[0, 1]$.

1409 Now assume this lemma holds for $K = J - 1 \in \mathbb{N}^+$. Our goal is to prove the case $K = J$.
 1410 Given any $\varepsilon \in (0, 1/100)$ and an arbitrary $\boldsymbol{\xi} = [\xi_1, \xi_2, \dots, \xi_J]^T \in [0, 1]^J$, our goal is to find
 1411 a proper $w \in \mathbb{R}$ such that

$$1412 \quad |\tau(w a_j) - \xi_j| < C\varepsilon \quad \text{for } j = 1, 2, \dots, J, \quad \text{where } C \text{ is an absolute constant.} \quad (27)$$

1413 We remark that the constant C in the above equation is actually equal to 11 in our proof.
 1414 As we shall see later, we need an assumption that the given point is in $[6\varepsilon, 1 - 6\varepsilon]^J$. Thus,
 1415 we slightly modify $\boldsymbol{\xi}$ by setting

$$1416 \quad \tilde{\xi}_j = \xi_j + 6\varepsilon \cdot \mathbf{1}_{\{\xi_j \leq 6\varepsilon\}} - 6\varepsilon \cdot \mathbf{1}_{\{\xi_j \geq 1 - 6\varepsilon\}} \quad \text{for } j = 1, 2, \dots, J.$$

1417 Then, we have

$$1418 \quad \tilde{\xi}_j \in [6\varepsilon, 1 - 6\varepsilon] \quad \text{for } j = 1, 2, \dots, J \quad (28)$$

1419 and

$$1420 \quad |\xi_j - \tilde{\xi}_j| = |6\varepsilon \cdot \mathbf{1}_{\{\xi_j \leq 6\varepsilon\}} - 6\varepsilon \cdot \mathbf{1}_{\{\xi_j \geq 1 - 6\varepsilon\}}| \leq 6\varepsilon \quad \text{for } j = 1, 2, \dots, J. \quad (29)$$

1421 For any $n \in \mathbb{N}^+$, we define

$$1422 \quad \hat{\xi}_j := \tau(\tilde{\xi}_j - \frac{\tilde{\xi}_j}{a_j} a_j) \quad \text{for } j = 1, 2, \dots, J.$$

1423 Then $\hat{\xi}_J = 0$ and $\hat{\xi}_j \in [0, 1)$ for $j = 1, 2, \dots, J - 1$. To approximate $[\hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_{J-1}]^T \in$
 1424 $[0, 1)^{J-1}$, we only need to consider $J - 1$ indices, and, therefore, we can use the induction
 1425 hypothesis to continue our proof.

1426 Clearly, the rational independence of a_1, a_2, \dots, a_J implies none of them is equal to zero.
 1427 Define

$$1428 \quad \mathbf{b}_n := \left[\tau\left(\frac{n}{a_J} a_1\right), \tau\left(\frac{n}{a_J} a_2\right), \dots, \tau\left(\frac{n}{a_J} a_{J-1}\right) \right]^T \in [0, 1)^{J-1}.$$

1429 Then, the bounded sequence $(\mathbf{b}_n)_{n=1}^\infty$ has a convergent subsequence by the Bolzano-Weierstrass
 1430 Theorem. Thus, there exist $n_1, n_2 \in \mathbb{N}^+$ with $n_1 < n_2$ such that $\|\mathbf{b}_{n_2} - \mathbf{b}_{n_1}\|_\infty < \varepsilon$, i.e.,

$$1431 \quad \left| \tau\left(\frac{n_2}{a_J} a_j\right) - \tau\left(\frac{n_1}{a_J} a_j\right) \right| < \varepsilon \quad \text{for } j = 1, 2, \dots, J - 1.$$

1432 Set $\hat{n} = n_2 - n_1 \in \mathbb{N}^+$ and

$$1433 \quad k_j = \lfloor \frac{n_1}{a_j} a_j \rfloor - \lfloor \frac{n_2}{a_j} a_j \rfloor \in \mathbb{Z} \quad \text{for } j = 1, 2, \dots, J-1.$$

1434 Then, by defining

$$1435 \quad \hat{a}_j := \frac{\hat{n}}{a_j} a_j + k_j \quad \text{for } j = 1, 2, \dots, J-1,$$

1436 we have

$$1437 \quad \begin{aligned} |\hat{a}_j| &= \left| \frac{\hat{n}}{a_j} a_j + k_j \right| = \left| \frac{n_2}{a_j} a_j - \frac{n_1}{a_j} a_j + \lfloor \frac{n_1}{a_j} a_j \rfloor - \lfloor \frac{n_2}{a_j} a_j \rfloor \right| \\ &= \left| \left(\frac{n_2}{a_j} a_j - \lfloor \frac{n_2}{a_j} a_j \rfloor \right) - \left(\frac{n_1}{a_j} a_j - \lfloor \frac{n_1}{a_j} a_j \rfloor \right) \right| \\ &= \left| \tau\left(\frac{n_2}{a_j} a_j\right) - \tau\left(\frac{n_1}{a_j} a_j\right) \right| < \varepsilon. \end{aligned} \quad (30)$$

1438 It is easy to verify that $\hat{a}_1, \hat{a}_2, \dots, \hat{a}_{J-1}$ are rationally independent. To see this, assume
1439 there exist $\lambda_1, \lambda_2, \dots, \lambda_{J-1} \in \mathbb{Q}$ such that

$$1440 \quad 0 = \sum_{j=1}^{J-1} \lambda_j \hat{a}_j = \sum_{j=1}^{J-1} \lambda_j \left(\frac{\hat{n}}{a_j} a_j + k_j \right) = \sum_{j=1}^{J-1} \lambda_j \frac{\hat{n}}{a_j} a_j + \sum_{j=1}^{J-1} \lambda_j k_j.$$

1441 It follows that

$$1442 \quad 0 = \sum_{j=1}^{J-1} \lambda_j \hat{n} a_j + \left(\sum_{j=1}^{J-1} \lambda_j k_j \right) a_J.$$

1443 Recall that $\hat{n} \in \mathbb{N}^+$, $k_j \in \mathbb{Z}$, and $\lambda_j \in \mathbb{Q}$ for any j . That means the coefficients $\lambda_j \hat{n}$ and
1444 $\sum_{j=1}^{J-1} \lambda_j k_j$ are rational for any j . Since a_1, a_2, \dots, a_J are rationally independent, we have

$$1445 \quad \lambda_j \hat{n} = 0 \quad \text{and} \quad \sum_{j=1}^{J-1} \lambda_j k_j = 0 \quad \text{for } j = 1, 2, \dots, J-1.$$

1446 It follows from $\hat{n} = n_2 - n_1 > 0$ that $\lambda_1 = \lambda_2 = \dots = \lambda_{J-1} = 0$. Therefore, $\hat{a}_1, \hat{a}_2, \dots, \hat{a}_{J-1}$
1447 are rationally independent as desired.

1448 By the induction hypothesis, the following set

$$1449 \quad \left\{ \left[\tau(s \cdot \hat{a}_1), \tau(s \cdot \hat{a}_2), \dots, \tau(s \cdot \hat{a}_{J-1}) \right]^T : s \in \mathbb{R} \right\} \subseteq [0, 1]^{J-1}$$

1450 is dense in $[0, 1]^{J-1}$. Recall that $\hat{\xi}_j = \tau(\tilde{\xi}_j - \frac{\tilde{\xi}_j}{a_j} a_j) \in [0, 1]$ for $j = 1, 2, \dots, J-1$, implying

$$1451 \quad \hat{\xi}_j + 3\varepsilon \cdot \mathbb{1}_{\{\hat{\xi}_j \leq 3\varepsilon\}} - 3\varepsilon \cdot \mathbb{1}_{\{\hat{\xi}_j \geq 1-3\varepsilon\}} \in [3\varepsilon, 1-3\varepsilon].$$

1452 Hence, there exists $s_0 \in \mathbb{R}$ such that

$$1453 \quad \left| \tau(s_0 \hat{a}_j) - \left(\hat{\xi}_j + 3\varepsilon \cdot \mathbb{1}_{\{\hat{\xi}_j \leq 3\varepsilon\}} - 3\varepsilon \cdot \mathbb{1}_{\{\hat{\xi}_j \geq 1-3\varepsilon\}} \right) \right| < \varepsilon$$

1454 for $j = 1, 2, \dots, J-1$. It follows that

$$1455 \quad \tau(s_0 \hat{a}_j) \in [2\varepsilon, 1-2\varepsilon] \quad \text{for } j = 1, 2, \dots, J-1$$

1456 and

$$1457 \quad \left| \tau(s_0 \widehat{a}_j) - \widehat{\xi}_j \right| < \varepsilon + \left| 3\varepsilon \cdot \mathbf{1}_{\{\widehat{\xi}_j \leq 3\varepsilon\}} - 3\varepsilon \cdot \mathbf{1}_{\{\widehat{\xi}_j \geq 1-3\varepsilon\}} \right| \leq 4\varepsilon \quad (31)$$

1458 for $j = 1, 2, \dots, J-1$.

1459 To estimate $\tau(\lfloor s_0 \rfloor \widehat{a}_j) - \widehat{\xi}_j$, we need to bound $\tau(s_0 \widehat{a}_j) - \tau(\lfloor s_0 \rfloor \widehat{a}_j)$. To this end, we need
1460 an observation for any $x, y \in \mathbb{R}$ as follows.

$$1461 \quad |x - y| < \varepsilon \quad \text{and} \quad \tau(x) \in [2\varepsilon, 1 - 2\varepsilon] \quad \implies \quad |\tau(x) - \tau(y)| < \varepsilon. \quad (32)$$

1462 In fact, $\tau(x) \in [2\varepsilon, 1 - 2\varepsilon]$ implies $\varepsilon \leq \tau(x) - \varepsilon \leq \tau(x) + \varepsilon \leq 1 - \varepsilon$, from which we deduce

$$1463 \quad \begin{aligned} y \in [x - \varepsilon, x + \varepsilon] &= \left[[x] + \underbrace{\tau(x) - \varepsilon}_{\geq \varepsilon}, [x] + \underbrace{\tau(x) + \varepsilon}_{\leq 1 - \varepsilon} \right] \\ &\subseteq \left[[x] + \varepsilon, [x] + 1 - \varepsilon \right] \subseteq \left[[x], [x] + 1 \right). \end{aligned}$$

1464 Then, we have $\lfloor y \rfloor = \lfloor x \rfloor$, which implies

$$1465 \quad \begin{aligned} |\tau(x) - \tau(y)| &= |\tau(x) - \tau(y) + \lfloor x \rfloor - \lfloor y \rfloor| \\ &= \left| (\tau(x) + \lfloor x \rfloor) - (\tau(y) + \lfloor y \rfloor) \right| = |x - y| < \varepsilon. \end{aligned}$$

1466 Thus, Equation (32) is proved.

1467 By Equation (30), we have

$$1468 \quad \left| s_0 \widehat{a}_j - \lfloor s_0 \rfloor \widehat{a}_j \right| \leq \left| s_0 - \lfloor s_0 \rfloor \right| \cdot |\widehat{a}_j| \leq |\widehat{a}_j| < \varepsilon \quad \text{for } j = 1, 2, \dots, J-1.$$

1469 Recall that

$$1470 \quad \tau(s_0 \widehat{a}_j) \in [2\varepsilon, 1 - 2\varepsilon] \quad \text{for } j = 1, \dots, J-1.$$

1471 Then, for each $j \in \{1, 2, \dots, J-1\}$, by the observation above in Equation (32) (set $x = s_0 \widehat{a}_j$
1472 and $y = \lfloor s_0 \rfloor \widehat{a}_j$ therein), we have $|\tau(s_0 \widehat{a}_j) - \tau(\lfloor s_0 \rfloor \widehat{a}_j)| < \varepsilon$.

1473 Recall that $\widehat{\xi}_j = \tau(\widetilde{\xi}_j - \frac{\widetilde{\xi}_j}{a_j} a_j)$ for $j = 1, 2, \dots, J$. Therefore, by Equation (31), we have

$$1474 \quad \begin{aligned} \left| \tau(\lfloor s_0 \rfloor \widehat{a}_j) - \tau(\widetilde{\xi}_j - \frac{\widetilde{\xi}_j}{a_j} a_j) \right| &= \left| \tau(\lfloor s_0 \rfloor \widehat{a}_j) - \widehat{\xi}_j \right| \\ &\leq \left| \tau(\lfloor s_0 \rfloor \widehat{a}_j) - \tau(s_0 \widehat{a}_j) \right| + \left| \tau(s_0 \widehat{a}_j) - \widehat{\xi}_j \right| < \varepsilon + 4\varepsilon = 5\varepsilon, \end{aligned}$$

1475 for $j = 1, 2, \dots, J-1$.

1476 Observe that, for any $x, y \in \mathbb{R}$, there exist $z \in \mathbb{Z}$ such that $\tau(x) - \tau(y) = x - y - z$. To
1477 see this, we set $z = \lfloor x \rfloor - \lfloor y \rfloor \in \mathbb{Z}$ and then $\tau(x) - \tau(y) = x - \lfloor x \rfloor - (y - \lfloor y \rfloor) = x - y - z$.

1478 Therefore, for $j = 1, 2, \dots, J-1$, there exists $z_j \in \mathbb{Z}$ such that

$$1479 \quad \tau(\lfloor s_0 \rfloor \widehat{a}_j) - \tau(\widetilde{\xi}_j - \frac{\widetilde{\xi}_j}{a_j} a_j) = \lfloor s_0 \rfloor \widehat{a}_j - (\widetilde{\xi}_j - \frac{\widetilde{\xi}_j}{a_j} a_j) - z_j = \lfloor s_0 \rfloor \widehat{a}_j + \frac{\widetilde{\xi}_j}{a_j} a_j - (z_j + \widetilde{\xi}_j),$$

1480 which implies

$$1481 \quad \left| \lfloor s_0 \rfloor \widehat{a}_j + \frac{\widetilde{\xi}_j}{a_j} a_j - (z_j + \widetilde{\xi}_j) \right| = \left| \tau(\lfloor s_0 \rfloor \widehat{a}_j) - \tau(\widetilde{\xi}_j - \frac{\widetilde{\xi}_j}{a_j} a_j) \right| < 5\varepsilon.$$

1482 It follows that, for $j = 1, 2, \dots, J - 1$,

$$1483 \quad \lfloor s_0 \rfloor \widehat{a}_j + \frac{\widetilde{\xi}_J}{a_J} a_j \in [z_j + \underbrace{\widetilde{\xi}_j - 5\varepsilon}_{\geq \varepsilon}, z_j + \underbrace{\widetilde{\xi}_j + 5\varepsilon}_{\leq 1 - \varepsilon}] \subseteq [z_j + \varepsilon, z_j + 1 - \varepsilon],$$

1484 where the fact $\varepsilon \leq \widetilde{\xi}_j - 5\varepsilon \leq \widetilde{\xi}_j + 5\varepsilon \leq 1 - \varepsilon$ comes from Equation (28). Therefore, we have

$$1485 \quad \left\lfloor \lfloor s_0 \rfloor \widehat{a}_j + \frac{\widetilde{\xi}_J}{a_J} a_j \right\rfloor = z_j \quad \text{for } j = 1, 2, \dots, J - 1,$$

1486 implying

$$1487 \quad \tau(\lfloor s_0 \rfloor \widehat{a}_j + \frac{\widetilde{\xi}_J}{a_J} a_j) = (\lfloor s_0 \rfloor \widehat{a}_j + \frac{\widetilde{\xi}_J}{a_J} a_j) - z_j \in [\widetilde{\xi}_j - 5\varepsilon, \widetilde{\xi}_j + 5\varepsilon].$$

1488 Clearly, we have

$$1489 \quad \lfloor s_0 \rfloor \widehat{a}_j + \frac{\widetilde{\xi}_J}{a_J} a_j = \lfloor s_0 \rfloor \left(\frac{\widehat{n}}{a_J} a_j + k_j \right) + \frac{\widetilde{\xi}_J}{a_J} a_j = \frac{\lfloor s_0 \rfloor \widehat{n} + \widetilde{\xi}_J}{a_J} a_j + \underbrace{k_j \lfloor s_0 \rfloor}_{\in \mathbb{Z}}$$

1490 for $j = 1, 2, \dots, J - 1$, which implies

$$1491 \quad \tau\left(\frac{\lfloor s_0 \rfloor \widehat{n} + \widetilde{\xi}_J}{a_J} a_j\right) = \tau(\lfloor s_0 \rfloor \widehat{a}_j + \frac{\widetilde{\xi}_J}{a_J} a_j) \in [\widetilde{\xi}_j - 5\varepsilon, \widetilde{\xi}_j + 5\varepsilon].$$

1492 We also need to consider the case $j = J$. By Equation (28), we have $\widetilde{\xi}_J \in [6\varepsilon, 1 - 6\varepsilon]$, from
1493 which we deduce

$$1494 \quad \tau\left(\frac{\lfloor s_0 \rfloor \widehat{n} + \widetilde{\xi}_J}{a_J} a_J\right) = \tau(\underbrace{\lfloor s_0 \rfloor \widehat{n}}_{\in \mathbb{Z}} + \widetilde{\xi}_J) = \widetilde{\xi}_J.$$

1495 Thus, for $j = 1, 2, \dots, J$, we have

$$1496 \quad \left| \tau\left(\frac{\lfloor s_0 \rfloor \widehat{n} + \widetilde{\xi}_J}{a_J} a_j\right) - \widetilde{\xi}_j \right| \leq 5\varepsilon.$$

1497 By Equation (29), we have $|\widetilde{\xi}_j - \xi_j| < 6\varepsilon$ for $j = 1, 2, \dots, J$, which implies

$$1498 \quad \left| \tau\left(\frac{\lfloor s_0 \rfloor \widehat{n} + \widetilde{\xi}_J}{a_J} a_j\right) - \xi_j \right| \leq \left| \tau\left(\frac{\lfloor s_0 \rfloor \widehat{n} + \widetilde{\xi}_J}{a_J} a_j\right) - \widetilde{\xi}_j \right| + |\widetilde{\xi}_j - \xi_j| \leq 5\varepsilon + 6\varepsilon = 11\varepsilon.$$

1499 That means $w_0 = \frac{\lfloor s_0 \rfloor \widehat{n} + \widetilde{\xi}_J}{a_J}$ is the desired w in Equation (27) and the constant $C > 0$ therein
1500 is 11. Therefore,

$$1501 \quad \left| \tau(w_0 a_j) - \xi_j \right| \leq 11\varepsilon \quad \text{for } j = 1, 2, \dots, J.$$

1502 Since $\boldsymbol{\xi} = [\xi_1, \xi_2, \dots, \xi_J]^T \in [0, 1]^J$ and $\varepsilon > 0$ are arbitrary, the following set

$$1503 \quad \left\{ [\tau(w a_1), \tau(w a_2), \dots, \tau(w a_J)]^T : w \in \mathbb{R} \right\} \subseteq [0, 1]^J$$

1504 is dense in $[0, 1]^J$ as desired. We finish the process of mathematical induction and therefore
1505 finish the proof by the principle of mathematical induction. \blacksquare

1506 We remark that the target parameter $w_0 = \frac{\lfloor s_0 \rfloor \widehat{n} + \widetilde{\xi}_J}{a_J}$ designed in the above proof may
1507 not be bounded uniformly for any approximation error ε since \widehat{n} can be arbitrarily large as
1508 ε goes to 0. Therefore, the network in Theorem 1 may require sufficiently large parameters
1509 to achieve an arbitrarily small error ε .

1510 8. Conclusion

1511 This paper studies the super approximation power of deep feed-forward neural networks
1512 activated by EUAF with a fixed size. It is proved by construction that there exists an EUAF
1513 network architecture with d input neurons, a maximum width $36d(2d+1)$, 11 hidden layers,
1514 and at most $5437(d+1)(2d+1)$ nonzero parameters, achieving the universal approximation
1515 property by only adjusting its finitely many parameters. That is, without changing the
1516 network size, our EUAF network can approximate any continuous function $f : [a, b]^d \rightarrow \mathbb{R}$
1517 within an arbitrarily small error $\varepsilon > 0$ with appropriate parameters depending on f , ε , d ,
1518 a , and b . Moreover, augmenting this EUAF network using one more layer with 2 neurons
1519 can exactly realize a classification function $\sum_{j=1}^J r_j \cdot \mathbb{1}_{E_j}$ in $\bigcup_{j=1}^J E_j$ for any $J \in \mathbb{N}^+$, where
1520 r_1, r_2, \dots, r_J are distinct rational numbers and E_1, E_2, \dots, E_J are arbitrary pairwise disjoint
1521 bounded closed subsets of \mathbb{R}^d .

1522 While we are interested in the analysis of the approximation error here, it would be very
1523 interesting to investigate the generalization and optimization errors of EUAF networks.
1524 Acting as a proof of concept, our experimentation shows the numerical advantages of EUAF
1525 compared to ReLU. We believe our EUAF activation function could be further developed
1526 and applied to real-world applications.

1527 Acknowledgments

1528 Z. Shen was supported by Distinguished Professorship of National University of Singapore.
1529 H. Yang was partially supported by the US National Science Foundation under award DMS-
1530 2244988, DMS-2206333, and the Office of Naval Research Young Investigator Award.

1531 References

1532 Andrew R. Barron. Universal approximation bounds for superpositions of a sigmoidal
1533 function. *IEEE Transactions on Information Theory*, 39(3):930–945, May 1993. ISSN
1534 0018-9448. doi: [10.1109/18.256500](https://doi.org/10.1109/18.256500).

1535 Andrea D. Beck, Jonas Zeifang, Anna Schwarz, and David G. Flad. A neural network
1536 based shock detection and localization approach for discontinuous Galerkin methods.
1537 *Journal of Computational Physics*, 423:article 109824, 2020. ISSN 0021-9991. doi:
1538 [10.1016/j.jcp.2020.109824](https://doi.org/10.1016/j.jcp.2020.109824).

1539 Christian Beck, Martin Hutzenthaler, Arnulf Jentzen, and Benno Kuckuck. An overview
1540 on deep learning-based approximation methods for partial differential equations. *arXiv*
1541 *e-prints*, page arXiv:2012.12348, December 2020. URL [https://arxiv.org/abs/2012.](https://arxiv.org/abs/2012.12348)
1542 [12348](https://arxiv.org/abs/2012.12348).

1543 David E. Bernholdt, Mark R. Cianciosa, David L. Green, Jin M. Park, Kody J. H. Law, and
1544 Clement Etienam. Cluster, classify, regress: A general method for learning discontinuous
1545 functions. *Foundations of Data Science*, 1(4):491–506, 2019. doi: [10.3934/fods.2019020](https://doi.org/10.3934/fods.2019020).

- 1546 Helmut Bölcskei, Philipp Grohs, Gitta Kutyniok, and Philipp Petersen. Optimal approxi-
 1547 mation with sparsely connected deep neural networks. *SIAM Journal on Mathematics of*
 1548 *Data Science*, 1(1):8–45, Jan 2019. ISSN 2577-0187. doi: [10.1137/18m118709x](https://doi.org/10.1137/18m118709x).
- 1549 Andrea Bonito, Ronald DeVore, Peter Jantsch Diane Guignard, and Guergana Petrova.
 1550 Polynomial approximation of anisotropic analytic functions of several variables. *Con-*
 1551 *structive Approximation*, 53:319–348, 2021. doi: [10.1007/s00365-020-09511-4](https://doi.org/10.1007/s00365-020-09511-4).
- 1552 Liang Chen and Congwei Wu. A note on the expressive power of deep rectified linear unit
 1553 networks in high-dimensional spaces. *Mathematical Methods in the Applied Sciences*, 42
 1554 (9):3400–3404, 2019. doi: [10.1002/mma.5575](https://doi.org/10.1002/mma.5575).
- 1555 Albert Cohen, Ronald DeVore, Guergana Petrova, and Przemyslaw Wojtaszczyk. Optimal
 1556 stable nonlinear approximation. *Foundations of Computational Mathematics*, 22:607–648,
 1557 2022. doi: [10.1007/s10208-021-09494-z](https://doi.org/10.1007/s10208-021-09494-z).
- 1558 George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics*
 1559 *of Control, Signals, and Systems*, 2:303–314, 1989. doi: [10.1007/BF02551274](https://doi.org/10.1007/BF02551274).
- 1560 Ingrid Daubechies, Ronald DeVore, Simon Foucart, Boris Hanin, and Guergana Petrova.
 1561 Nonlinear approximation and (deep) ReLU networks. *Constructive Approximation*, 55:
 1562 127–172, 2022. doi: [10.1007/s00365-021-09548-z](https://doi.org/10.1007/s00365-021-09548-z).
- 1563 Ronald A. DeVore. Nonlinear approximation. *Acta Numerica*, 7:51–150, 1998. doi:
 1564 [10.1017/S0962492900002816](https://doi.org/10.1017/S0962492900002816).
- 1565 Weinan E and Qingcan Wang. Exponential convergence of the deep neural network ap-
 1566 proximation for analytic functions. *Science China Mathematics*, 61:1733–1740, 2018. doi:
 1567 [10.1007/s11425-018-9387-x](https://doi.org/10.1007/s11425-018-9387-x).
- 1568 Weinan E and Stephan Wojtowytsch. A priori estimates for classification problems using
 1569 neural networks. *CoRR*, abs/2009.13500, 2020. URL [https://arxiv.org/abs/2009.](https://arxiv.org/abs/2009.13500)
 1570 [13500](https://arxiv.org/abs/2009.13500).
- 1571 Weinan E and Stephan Wojtowytsch. Representation formulas and pointwise properties for
 1572 Barron functions. *Calculus of Variations and Partial Differential Equations*, 61:article
 1573 46, 2022. ISSN 0944-2669. doi: [10.1007/s00526-021-02156-6](https://doi.org/10.1007/s00526-021-02156-6).
- 1574 Weinan E, Chao Ma, and Qingcan Wang. A priori estimates of the population risk for
 1575 residual networks. *CoRR*, abs/1903.02154, 2019a. URL [http://arxiv.org/abs/1903.](http://arxiv.org/abs/1903.02154)
 1576 [02154](http://arxiv.org/abs/1903.02154).
- 1577 Weinan E, Chao Ma, and Lei Wu. A priori estimates of the population risk for two-layer
 1578 neural networks. *Communications in Mathematical Sciences*, 17(5):1407–1425, 2019b.
 1579 doi: [10.4310/CMS.2019.v17.n5.a11](https://doi.org/10.4310/CMS.2019.v17.n5.a11).
- 1580 Dennis Elbrächter, Philipp Grohs, Arnulf Jentzen, and Christoph Schwab. DNN expression
 1581 rate analysis of high-dimensional PDEs: Application to option pricing. *Constructive*
 1582 *Approximation*, 55:3–71, 2022. doi: [10.1007/s00365-021-09541-6](https://doi.org/10.1007/s00365-021-09541-6).

- 1583 Johannes Gedeon, Jonathan Schmidt, Matthew J.P. Hodgson, Jack Wetherell, Car-
1584 los L. Benavides-Riveros, and Miguel A. L. Marques. Machine learning the deriva-
1585 tive discontinuity of density-functional theory. *Machine Learning: Science and Tech-*
1586 *nology*, 3:article 015011, 2021. URL [http://iopscience.iop.org/article/10.1088/](http://iopscience.iop.org/article/10.1088/2632-2153/ac3149)
1587 [2632-2153/ac3149](http://iopscience.iop.org/article/10.1088/2632-2153/ac3149).
- 1588 Vanshika Gupta, Sharad Kumar Gupta, and Jungrack Kim. Automated discontinuity de-
1589 tection and reconstruction in subsurface environment of Mars using deep learning: A case
1590 study of sharad observation. *Applied Sciences*, 10(7):article 2279, 2020. ISSN 2076-3417.
1591 URL <https://www.mdpi.com/2076-3417/10/7/2279>.
- 1592 Jun Han and Claudio Moraga. The influence of the sigmoid function parameters on the speed
1593 of backpropagation learning. In José Mira and Francisco Sandoval, editors, *From Natural*
1594 *to Artificial Neural Computation*, pages 195–201, Berlin, Heidelberg, 1995. Springer Berlin
1595 Heidelberg. ISBN 978-3-540-49288-7. doi: [10.1007/3-540-59497-3_175](https://doi.org/10.1007/3-540-59497-3_175).
- 1596 Juncai He, Xiaodong Jia, Jinchao Xu, Lian Zhang, and Liang Zhao. Make ℓ_1 regularization
1597 effective in training sparse CNN. *Computational Optimization and Applications*, 77(1):
1598 163–182, 2020. doi: [10.1007/s10589-020-00202-1](https://doi.org/10.1007/s10589-020-00202-1).
- 1599 Geoffrey E. Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan
1600 Salakhutdinov. Improving neural networks by preventing co-adaptation of feature de-
1601 tectors. *CoRR*, abs/1207.0580, 2012. URL <http://arxiv.org/abs/1207.0580>.
- 1602 Sean Hon and Haizhao Yang. Simultaneous neural network approximations in Sobolev
1603 spaces. *arXiv e-prints*, page arXiv:2109.00161, August 2021. URL [https://arxiv.org/](https://arxiv.org/abs/2109.00161)
1604 [abs/2109.00161](https://arxiv.org/abs/2109.00161).
- 1605 Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Net-*
1606 *works*, 4(2):251–257, 1991. ISSN 0893-6080. doi: [10.1016/0893-6080\(91\)90009-T](https://doi.org/10.1016/0893-6080(91)90009-T).
- 1607 Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks
1608 are universal approximators. *Neural Networks*, 2(5):359–366, 1989. ISSN 0893-6080. doi:
1609 [10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8).
- 1610 Wei-Fan Hu, Te-Sheng Lin, and Ming-Chih Lai. A discontinuity capturing shallow neural
1611 network for elliptic interface problems. *arXiv e-prints*, page arXiv:2106.05587, June 2021.
1612 URL <https://arxiv.org/abs/2106.05587>.
- 1613 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network train-
1614 ing by reducing internal covariate shift. In *Proceedings of the 32nd International Con-*
1615 *ference on International Conference on Machine Learning - Volume 37, ICML’15*, pages
1616 448–456. JMLR.org, 2015. URL <https://proceedings.mlr.press/v37/ioffe15.html>.
- 1617 Yuling Jiao, Yanming Lai, Xiliang Lu, and Jerry Zhijian Yang. Deep neural networks with
1618 ReLU-sine-exponential activations break curse of dimensionality on Hölder class. *CoRR*,
1619 abs/2103.00542, 2021. URL <https://arxiv.org/abs/2103.00542>.

- 1620 Kenji Kawaguchi. Deep learning without poor local minima. In D. D. Lee, M. Sugiyama,
1621 U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information*
1622 *Processing Systems 29*, pages 586–594. Curran Associates, Inc., 2016. URL [http://](http://papers.nips.cc/paper/6112-deep-learning-without-poor-local-minima.pdf)
1623 papers.nips.cc/paper/6112-deep-learning-without-poor-local-minima.pdf.
- 1624 Kenji Kawaguchi and Yoshua Bengio. Depth with nonlinearity creates no bad local
1625 minima in resnets. *Neural Networks*, 118:167–174, 2019. ISSN 0893-6080. doi:
1626 [10.1016/j.neunet.2019.06.009](https://doi.org/10.1016/j.neunet.2019.06.009).
- 1627 Andrei Nikolaevich Kolmogorov. On the representation of continuous functions of many
1628 variables by superposition of continuous functions of one variable and addition. *Doklady*
1629 *Akademii Nauk SSSR*, 114(5):953–956, 1957. URL <http://mi.mathnet.ru/dan22050>.
- 1630 Phong Le and Willem Zuidema. Compositional distributional semantics with long short
1631 term memory. In *Proceedings of the Fourth Joint Conference on Lexical and Computa-*
1632 *tional Semantics*, pages 10–19, Denver, Colorado, June 2015. Association for Computa-
1633 tional Linguistics. doi: [10.18653/v1/S15-1002](https://doi.org/10.18653/v1/S15-1002).
- 1634 Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and dynamics of
1635 stochastic gradient algorithms I: Mathematical foundations. *Journal of Machine Learning*
1636 *Research*, 20(40):1–47, 2019. URL <http://jmlr.org/papers/v20/17-526.html>.
- 1637 Qianxiao Li, Ting Lin, and Zuowei Shen. Deep learning via dynamical systems: An ap-
1638 proximation perspective. *Journal of European Mathematical Society*, to appear. doi:
1639 [10.4171/JEMS/1221](https://doi.org/10.4171/JEMS/1221).
- 1640 Hongzhou Lin and Stefanie Jegelka. Resnet with one-neuron hidden layers is a universal
1641 approximator. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi,
1642 and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31.
1643 Curran Associates, Inc., 2018. URL [https://proceedings.neurips.cc/paper/2018/](https://proceedings.neurips.cc/paper/2018/file/03bfc1d4783966c69cc6aef8247e0103-Paper.pdf)
1644 [file/03bfc1d4783966c69cc6aef8247e0103-Paper.pdf](https://proceedings.neurips.cc/paper/2018/file/03bfc1d4783966c69cc6aef8247e0103-Paper.pdf).
- 1645 Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and
1646 Jiawei Han. On the variance of the adaptive learning rate and beyond. In *International*
1647 *Conference on Learning Representations*, 2020. URL [https://openreview.net/forum?](https://openreview.net/forum?id=rkgz2aEKDr)
1648 [id=rkgz2aEKDr](https://openreview.net/forum?id=rkgz2aEKDr).
- 1649 Jianfeng Lu, Zuowei Shen, Haizhao Yang, and Shijun Zhang. Deep network approximation
1650 for smooth functions. *SIAM Journal on Mathematical Analysis*, 53(5):5465–5506, 2021.
1651 doi: [10.1137/20M134695X](https://doi.org/10.1137/20M134695X).
- 1652 Vitaly Maiorov and Allan Pinkus. Lower bounds for approximation by MLP neural net-
1653 works. *Neurocomputing*, 25(1):81–91, 1999. ISSN 0925-2312. doi: [10.1016/S0925-](https://doi.org/10.1016/S0925-2312(98)00111-8)
1654 [2312\(98\)00111-8](https://doi.org/10.1016/S0925-2312(98)00111-8).
- 1655 Hadrien Montanelli, Haizhao Yang, and Qiang Du. Deep ReLU networks overcome the
1656 curse of dimensionality for generalized bandlimited functions. *Journal of Computational*
1657 *Mathematics*, 39(6):801–815, 2021. ISSN 1991-7139. doi: [10.4208/jcm.2007-m2019-0239](https://doi.org/10.4208/jcm.2007-m2019-0239).

- 1658 Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro.
1659 The role of over-parametrization in generalization of neural networks. In *International*
1660 *Conference on Learning Representations*, 2019. URL [https://openreview.net/forum?](https://openreview.net/forum?id=BygfgHAcYX)
1661 [id=BygfgHAcYX](https://openreview.net/forum?id=BygfgHAcYX).
- 1662 Quynh Nguyen and Matthias Hein. The loss surface of deep and wide neural networks.
1663 In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Con-*
1664 *ference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*,
1665 pages 2603–2612. PMLR, 06–11 Aug 2017. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v70/nguyen17a.html)
1666 [v70/nguyen17a.html](https://proceedings.mlr.press/v70/nguyen17a.html).
- 1667 Philipp Petersen and Felix Voigtlaender. Optimal approximation of piecewise smooth func-
1668 tions using deep ReLU neural networks. *Neural Networks*, 108:296–330, 2018. ISSN
1669 0893-6080. doi: [10.1016/j.neunet.2018.08.019](https://doi.org/10.1016/j.neunet.2018.08.019).
- 1670 Zuowei Shen, Haizhao Yang, and Shijun Zhang. Deep network approximation characterized
1671 by number of neurons. *Communications in Computational Physics*, 28(5):1768–1811,
1672 2020. ISSN 1991-7120. doi: [10.4208/cicp.OA-2020-0149](https://doi.org/10.4208/cicp.OA-2020-0149).
- 1673 Zuowei Shen, Haizhao Yang, and Shijun Zhang. Deep network with approximation error
1674 being reciprocal of width to power of square root of depth. *Neural Computation*, 33(4):
1675 1005–1036, 2021a. ISSN 0899-7667. doi: [10.1162/neco.a.01364](https://doi.org/10.1162/neco.a.01364).
- 1676 Zuowei Shen, Haizhao Yang, and Shijun Zhang. Neural network approximation: Three
1677 hidden layers are enough. *Neural Networks*, 141:160–173, 2021b. ISSN 0893-6080. doi:
1678 [10.1016/j.neunet.2021.04.011](https://doi.org/10.1016/j.neunet.2021.04.011).
- 1679 Zuowei Shen, Haizhao Yang, and Shijun Zhang. Optimal approximation rate of ReLU
1680 networks in terms of width and depth. *Journal de Mathématiques Pures et Appliquées*,
1681 157:101–135, 2022. ISSN 0021-7824. doi: [10.1016/j.matpur.2021.07.009](https://doi.org/10.1016/j.matpur.2021.07.009).
- 1682 Jonathan W. Siegel and Jinchao Xu. Optimal approximation rates and metric entropy of
1683 ReLU^k and cosine networks. *arXiv e-prints*, page arXiv:2101.12365, January 2021. URL
1684 <https://www.doi.org/abs/2101.12365>.
- 1685 Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhut-
1686 dinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of*
1687 *Machine Learning Research*, 15(56):1929–1958, 2014. URL [http://jmlr.org/papers/](http://jmlr.org/papers/v15/srivastava14a.html)
1688 [v15/srivastava14a.html](http://jmlr.org/papers/v15/srivastava14a.html).
- 1689 Joseph Turian, James Bergstra, and Yoshua Bengio. Quadratic features and deep architec-
1690 tures for chunking. In *Proceedings of Human Language Technologies: The 2009 Annual*
1691 *Conference of the North American Chapter of the Association for Computational Lin-*
1692 *guistics, Companion Volume: Short Papers*, NAACL-Short '09, pages 245–248, USA,
1693 2009. Association for Computational Linguistics. URL [https://aclanthology.org/](https://aclanthology.org/N09-2062/)
1694 [N09-2062/](https://aclanthology.org/N09-2062/).
- 1695 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset
1696 for benchmarking machine learning algorithms. *arXiv e-prints*, page arXiv:1708.07747,
1697 August 2017. URL <https://arxiv.org/abs/1708.07747>.

- 1698 Yunfei Yang, Zhen Li, and Yang Wang. Approximation in shift-invariant spaces with deep
1699 ReLU neural networks. *Neural Networks*, 153:269–281, 2022. ISSN 0893-6080. doi:
1700 [10.1016/j.neunet.2022.06.013](https://doi.org/10.1016/j.neunet.2022.06.013).
- 1701 Dmitry Yarotsky. Optimal approximation of continuous functions by very deep ReLU net-
1702 works. In Sébastien Bubeck, Vianney Perchet, and Philippe Rigollet, editors, *Proceedings*
1703 *of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learn-*
1704 *ing Research*, pages 639–649. PMLR, 06–09 Jul 2018. URL [http://proceedings.mlr.](http://proceedings.mlr.press/v75/yarotsky18a.html)
1705 [press/v75/yarotsky18a.html](http://proceedings.mlr.press/v75/yarotsky18a.html).
- 1706 Dmitry Yarotsky. Elementary superexpressive activations. In Marina Meila and Tong
1707 Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*,
1708 volume 139 of *Proceedings of Machine Learning Research*, pages 11932–11940. PMLR,
1709 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/yarotsky21a.html>.
- 1710 Dmitry Yarotsky and Anton Zhevnerchuk. The phase diagram of approximation rates for
1711 deep neural networks. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and
1712 H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages
1713 13005–13015. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/](https://proceedings.neurips.cc/paper/2020/file/979a3f14bae523dc5101c52120c535e9-Paper.pdf)
1714 [paper/2020/file/979a3f14bae523dc5101c52120c535e9-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/979a3f14bae523dc5101c52120c535e9-Paper.pdf).
- 1715 Shijun Zhang. Deep neural network approximation via function compositions. *PhD The-*
1716 *sis, National University of Singapore*, 2020. URL [https://scholarbank.nus.edu.sg/](https://scholarbank.nus.edu.sg/handle/10635/186064)
1717 [handle/10635/186064](https://scholarbank.nus.edu.sg/handle/10635/186064).