

1 **SELF-ADAPTIVE ReLU NEURAL NETWORK METHOD**
 2 **IN LEAST-SQUARES DATA FITTING***

3 ZHIQIANG CAI[†] AND MIN LIU[‡]

4 **Abstract.** This chapter provides a comprehensive introduction to a self-adaptive ReLU neural network method
 5 proposed recently in [11, 10, 5]. The purpose of the method is to design a nearly minimal neural network architecture
 6 to achieve the prescribed accuracy for a given task in scientific machine learning such as approximating a function
 7 or a solution of partial differential equation. Starting with a small one hidden-layer neural network, the method
 8 enhances the network adaptively by adding neurons in the current or new hidden-layer based on accuracy of the
 9 current approximation. In addition, the method provides a natural process for obtaining a good initialization in
 10 training the current network. Moreover, initialization of newly added neurons at each adaptive step is discussed in
 11 detail.

12 **Key words.** Self-adaptivity, Least-squares data fitting, Deep neural network, ReLU activation

13 **1. Introduction.** Given a data set $\{(\mathbf{x}_i, y_i)\}_{i=1}^M$ with $\mathbf{x}_i \in \Omega = [-1, 1]^d$ and positive weights
 14 $\{w_i\}_{i=1}^M$, consider the discrete least-squares problem: finding $f_{nn}(\mathbf{x}) \in \mathcal{M}(l)$ such that

15 (1.1)
$$f_{nn} = \arg \min_{v \in \mathcal{M}(l)} L(v),$$

16 where $\mathcal{M}(l)$ is a ReLU neuron network defined in section 2 with l hidden-layers and $L(\cdot)$ is a
 17 least-squares loss functional given by

18
$$L(v) = \sum_{i=1}^M w_i (v(\mathbf{x}_i) - y_i)^2.$$

19 For a prescribed tolerance $\epsilon > 0$, this chapter presents a self-adaptive algorithm, the adaptive
 20 neuron enhancement method (ANE), to adaptively construct a nearly optimal network \mathcal{M}^* such
 21 that the neural network approximation $f_{nn}(\mathbf{x})$ satisfies

22 (1.2)
$$L(f_{nn}) \leq \epsilon L(0),$$

23 where $L(0) = \sum_{i=1}^M w_i y_i^2$ is the square of the weighted l^2 norm of the output data $\{y_i\}_{i=0}^M$.

24 Multi-layer ReLU neural network is described in this chapter as a set of continuous *piece-wise*
 25 linear functions. Hence each network function is piece-wise linear with respect to a partition of
 26 the domain. This partition, referred as the (domain) physical partition (see section 3), provides
 27 geometric feature of the function and hence plays a critical role in the design of self-adaptive
 28 neural network method. Determination of this physical partition for a network function is in
 29 general computationally expensive, especially when the input dimension d is high. To circumvent
 30 this difficulty, we introduce a network indicator function that can easily determine such partition.

31 The idea of the ANE is similar to that of standard adaptive mesh-based numerical methods,
 32 and may be written as loops of the form

33 (1.3) **train → estimate → mark → enhance.**

*This work was supported in part by the National Science Foundation under grant DMS-2110571.

[†]Department of Mathematics, Purdue University, 150 N. University Street, West Lafayette, IN 47907-2067
 (caiz@purdue.edu).

[‡]School of Mechanical Engineering, Purdue University, 585 Purdue Mall, West Lafayette, IN 47907-
 2088(liu66@purdue.edu).

34 Starting with a small one hidden layer network, the step **train** is to iteratively solve the opti-
 35 mization problem of the current network; the step **estimate** is to compute error of the current
 36 approximation; the step **mark** is to identify local regions that need refinement; and the step
 37 **enhance** is to add new neurons to the current network with good initialization. This adaptive
 38 algorithm learns not only from given information (data, function, partial differential equation) but
 39 also from the current computer simulation.

40 When the current error does not satisfy (1.2), an efficient ANE method relies on strategies to
 41 address the following questions at each adaptive step:

- 42 (a) how many new neurons should be added at the last hidden layer?
- 43 (b) when should a new hidden layer be added?

44 By exploiting the geometric feature of the current approximation, the enhancement strategy (see
 45 section 4) determines the number of new neurons to be added at the last hidden layer. A new layer
 46 is added if a computable quantity measuring the improvement rate of two consecutive networks
 47 per the relative increase of parameters is small.

48 Problem (1.1) is a non-convex optimization that has many solutions, and the desired one is
 49 only attainable when one begins with an initial approximation that is sufficiently close. A common
 50 approach to obtaining a good initialization is through the method of continuation, as described in
 51 [1]. The ANE method offers a natural way to acquire a well-suited initialization. Essentially, the
 52 approximation provided by the previous network serves as a good starting point for the current
 53 network at each adaptive step. Additionally, we outline an approach for initializing the weights and
 54 biases of newly added neurons, leveraging the geometric properties of the current approximation,
 55 which is detailed in section 5.

56 **2. ReLU Neural Network.** A neural network defines a new class of approximating functions
 57 which is suitable for some computationally challenging problems. This section describes l -hidden-
 58 layer ReLU neural network as a set of continuous piece-wise linear functions and introduces related
 59 notations. This chapter is restricted to one dimensional output $n_{l+1} = 1$ for simplicity of presen-
 60 tation. Extension of materials covered by this chapter to multi-dimensional output $n_{l+1} > 1$ is
 61 straightforward.

62 ReLU refers to the rectified linear activation function defined by

$$63 \quad (2.1) \quad \sigma(t) = \max\{0, t\} = \begin{cases} t, & t > 0, \\ 0, & t \leq 0. \end{cases}$$

64 The $\sigma(t)$ is a continuous piece-wise linear function with one *breaking* point $t = 0$. For $k = 1, \dots, l$,
 65 let n_k denote the number of neurons at the k^{th} hidden-layer; denote by

$$66 \quad \mathbf{b}^{(k)} \in \mathbb{R}^{n_k} \quad \text{and} \quad \boldsymbol{\omega}^{(k)} \in \mathbb{R}^{n_k \times n_{k-1}}$$

67 the biases and weights of neurons at the k^{th} hidden-layer, respectively. Their i^{th} rows are denoted
 68 by $b_i^{(k)} \in \mathbb{R}$ and $\omega_i^{(k)} \in \mathbb{R}^{n_{k-1}}$, that are the bias and weights of the i^{th} neuron at the k^{th} hidden-
 69 layer, respectively. Introduce a vector-valued function $\mathbf{N}^{(k)} : \mathbb{R}^{n_{k-1}} \rightarrow \mathbb{R}^{n_k}$ as

$$70 \quad (2.2) \quad \mathbf{N}^{(k)}(\mathbf{x}^{(k-1)}) = \sigma\left(\boldsymbol{\omega}^{(k)}\mathbf{x}^{(k-1)} + \mathbf{b}^{(k)}\right) \quad \text{for } \mathbf{x}^{(k-1)} \in \mathbb{R}^{n_{k-1}},$$

71 where application of the activation function σ to a vector-valued function is defined component-
 72 wisesly and $n_0 = d$ is the input dimension.

73 A ReLU neural network with l hidden-layers and n_k neurons at the k^{th} hidden-layer may be
 74 defined as the collection of continuous piece-wise linear functions:

$$75 \quad (2.3) \quad \mathcal{M}(l) = \left\{ \mathbf{c}_1(\mathbf{N}^{(l)} \circ \dots \circ \mathbf{N}^{(1)}(\mathbf{x})) + c_0 : \begin{array}{l} (c_0, \mathbf{c}_1) \in \mathbb{R}^{n_l+1}, \boldsymbol{\omega}^{(k)} \in \mathbb{R}^{n_k \times n_{k-1}}, \\ \mathbf{b}^{(k)} \in \mathbb{R}^{n_k} \text{ for } k = 1, \dots, l \end{array} \right\},$$

76 where the symbol \circ denotes the composition of functions. The total number of parameters of $\mathcal{M}(l)$
77 is given by

$$78 \quad (2.4) \quad M(l) = (n_l + 1) + \sum_{k=1}^l n_k \times (n_{k-1} + 1).$$

79 As in [4], the biases and weights of all hidden-layers

$$80 \quad (2.5) \quad \Theta^{(l)} = \bigcup_{k=1}^l \left\{ \left(b_i^{(k)}, \omega_i^{(k)} \right) \right\}_{i=1}^{n_k} = \left\{ \left(\mathbf{b}^{(k)}, \boldsymbol{\omega}^{(k)} \right) \right\}_{k=1}^l$$

81 are referred as nonlinear parameters, and the output bias and weights

$$82 \quad \mathbf{c} = (c_0, \mathbf{c}_1) = (c_0, c_1, \dots, c_{n_l}) \in \mathbb{R}^{n_l+1}$$

83 are referred as linear parameters for a neural network function.

84 **REMARK 2.1.** *Domain of the nonlinear parameter $\Theta^{(l)}$ in (2.3) is too large in general and*
85 *hence admit infinite many global minimizers of (1.1). One may add some constraints to the domain*
86 *in order to reduce the number/dimension of the global minimizers. For example, the weights of*
87 *each neuron can be normalized (see, e.g., [4, 11, 8]).*

88 Linearity of the output parameter \mathbf{c} here means that \mathbf{c} is uniquely determined by a system of
89 linear algebraic equations with given nonlinear parameter $\Theta^{(l)}$. In the remainder of this section,
90 we introduce this linear system and show that the corresponding mass matrix is always symmetric;
91 moreover, it is positive definite under some condition. To this end, let

$$92 \quad (2.6) \quad \varphi_0^{(l)}(\mathbf{x}) = 1 \quad \text{and} \quad \varphi_i^{(l)}(\mathbf{x}) = \sigma \left(\boldsymbol{\omega}_i^{(l)} \left(\mathbf{N}^{(l-1)} \circ \dots \circ \mathbf{N}^{(1)}(\mathbf{x}) \right) + b_i^{(l)} \right),$$

93 then any function $v \in \mathcal{M}(l)$ has the form of

$$94 \quad (2.7) \quad v(\mathbf{x}) = \sum_{i=0}^{n_l} c_i \varphi_i^{(l)}(\mathbf{x}).$$

95 A solution f_{nn} of (1.1) satisfies the critical point equation

$$96 \quad (2.8) \quad \nabla_{\mathbf{c}} L(f_{nn}) = \mathbf{0}$$

97 for the linear parameter $\mathbf{c} = (c_0, \mathbf{c}_1)$. This implies that \mathbf{c} satisfies the following system of linear
98 algebraic equations

$$99 \quad (2.9) \quad \mathbf{M}^{(l)} \left(\Theta^{(l)} \right) \mathbf{c} = \mathbf{F}^{(l)} \left(\Theta^{(l)} \right),$$

100 where $\mathbf{M}^{(l)} \left(\Theta^{(l)} \right)$ and $\mathbf{F}^{(l)} \left(\Theta^{(l)} \right)$ are the discrete mass matrix and the right-hand side vector
101 given by

$$102 \quad (2.10) \quad \begin{cases} \mathbf{M}^{(l)} \left(\Theta^{(l)} \right) = \left(\sum_{e=1}^M w_e \varphi_j^{(l)}(\mathbf{x}_e) \varphi_i^{(l)}(\mathbf{x}_e) \right)_{(n_l+1) \times (n_l+1)} & \text{and} \\ \mathbf{F}^{(l)} \left(\Theta^{(l)} \right) = \left(\sum_{e=1}^M w_e y_e \varphi_i^{(l)}(\mathbf{x}_e) \right)_{(n_l+1) \times 1} \end{cases}.$$

103 LEMMA 2.2. The mass matrix $\mathbf{M}^{(l)}(\Theta^{(l)})$ defined in (2.10) is symmetric. Assume that func-
 104 tions $\{\varphi_i^{(l)}(\mathbf{x})\}_{i=0}^{n_l}$ are linearly independent, then $\mathbf{M}^{(l)}(\Theta^{(l)})$ is positive definite.

105 *Proof.* Obviously, $\mathbf{M}^{(l)}(\Theta^{(l)})$ is symmetric. For any $\mathbf{c} \in \mathbb{R}^{n_l+1}$, we have

$$106 \quad \mathbf{c}^T \mathbf{M}^{(l)}(\Theta^{(l)}) \mathbf{c} = \sum_{i,j=0}^{n_l+1} \sum_{e=1}^M c_i c_j w_e \varphi_j^{(l)}(\mathbf{x}_e) \varphi_i^{(l)}(\mathbf{x}_e) = \sum_{e=1}^M w_e \left(\sum_{i=0}^{n_l+1} c_i \varphi_i^{(l)}(\mathbf{x}_e) \right)^2,$$

107 which, together with the assumption, implies positive definiteness of $\mathbf{M}^{(l)}(\Theta^{(l)})$. \square

108 Even though $\mathbf{M}^{(l)}(\Theta^{(l)})$ is symmetric, positive definite, it could be highly ill-conditioned.
 109 This fact, in turn, implies inefficiency of the optimization methods of gradient descent type.

110 **3. Physical Partition.** A neural network function in $\mathcal{M}(l)$ has the form of

$$111 \quad (3.1) \quad v(\mathbf{x}) = \mathbf{c}_1 \left(\mathbf{N}^{(l)} \circ \dots \circ \mathbf{N}^{(1)}(\mathbf{x}) \right) + c_0 = \sum_{i=0}^{n_l} c_i \varphi_i^{(l)}(\mathbf{x}),$$

112 where $\varphi_i^{(l)}(\mathbf{x})$ is defined in (2.6). Obviously, $v(\mathbf{x})$ is a continuous *piece-wise* linear (CPWL) function
 113 defined in \mathbb{R}^d . This means that there exists a partition of \mathbb{R}^d such that $v(\mathbf{x})$ is linear on all
 114 subdomains of this partition. This section studies such a partition for a given neural network
 115 function $v(\mathbf{x})$ of the form in (3.1).

116 DEFINITION 3.1. For a given network function $v(\mathbf{x})$ of the form in (3.1) defined in $\Omega =$
 117 $[-1, 1]^d$, a partition $\mathcal{K}^{(l)}(v)$ of Ω is said to be the physical partition of $v(\mathbf{x})$ with respect to Ω
 118 if

119 (i) $\mathcal{K}^{(l)}(v)$ is a partition of Ω , i.e.,

$$120 \quad \Omega = \bigcup_{K \in \mathcal{K}^{(l)}(v)} \bar{K} \quad \text{and} \quad K \cap T = \emptyset \quad \text{if} \quad K \neq T \quad \text{for all} \quad K, T \in \mathcal{K}^{(l)}(v).$$

121 (ii) for each subdomain $K \in \mathcal{K}^{(l)}(v)$, the restriction of $v(\mathbf{x})$ on K is a linear function.

122 REMARK 3.2. The physical partition $\mathcal{K}^{(l)}(v)$ defined in Definition 3.1 depends on the nonlinear
 123 parameter $\Theta^{(l)}$ but not on the linear parameter \mathbf{c} .

124 For a shallow neural network $\mathcal{M}(1)$, each function $v \in \mathcal{M}(1)$ has the form of

$$125 \quad v(\mathbf{x}) = \mathbf{c}_1 \left(\mathbf{N}^{(1)}(\mathbf{x}) \right) + c_0 = \sum_{i=1}^{n_1} c_i \sigma \left(\omega_i^{(1)} \mathbf{x} + b_i^{(1)} \right) + c_0 = \sum_{i=0}^{n_1} c_i \varphi_i^{(1)}(\mathbf{x}).$$

126 where $\Theta^{(1)} = \left\{ \theta_i^{(1)} \right\}_{i=1}^{n_1} := \left\{ \left(b_i^{(1)}, \omega_i^{(1)} \right) \right\}_{i=1}^{n_1}$ is nonlinear parameter. For $i = 1, \dots, n_1$, denote
 127 the pre-activation function of the i^{th} neuron by

$$128 \quad (3.2) \quad g_i^{(1)}(\mathbf{x}) = \omega_i^{(1)} \mathbf{x} + b_i^{(1)}$$

129 and its zero level set, called the *breaking hyper-plane*, by

$$130 \quad (3.3) \quad \mathcal{P}_i^{(1)} \left(\theta_i^{(1)} \right) = \left\{ \mathbf{x} \in \Omega : g_i^{(1)}(\mathbf{x}) = 0 \right\} = \left\{ \mathbf{x} \in \Omega : \omega_i^{(1)} \mathbf{x} + b_i^{(1)} = 0 \right\}.$$

131 For fixed $\Theta^{(1)}$, the physical partition $\mathcal{K}^{(1)}(v)$ is formed by the set of the breaking hyper-planes
 132 $\left\{\mathcal{P}_i^{(1)}\right\}_{i=1}^{n_1}$ and the boundary of the domain Ω .

133 The breaking hyper-planes $\left\{\mathcal{P}_i^{(1)}\right\}_{i=1}^{n_1}$ in one dimension ($d = 1$) degenerate to the breaking
 134 points $\left\{-\frac{b_i^{(1)}}{\omega_i^{(1)}}\right\}_{i=1}^{n_1}$ (blue dots in Fig. 1(a)), which partitions the interval $\Omega = [-1, 1]$ into sub-
 135 intervals. The breaking hyper-planes in two dimensions ($d = 2$) degenerate to the breaking lines
 136 (blue lines in Fig. 1(b))

$$137 \quad \mathcal{P}_i^{(1)}\left(\theta_i^{(1)}\right) = \left\{\mathbf{x} = (x_1, x_2) \in \Omega = [-1, 1]^2 : \omega_{i1}^{(1)}x_1 + \omega_{i2}^{(1)}x_2 + b_i^{(1)} = 0\right\},$$

138 which partition the domain $\Omega \in \mathbb{R}^2$ into irregular, polygonal sub-domains (see also Fig. 2(c)).

139 For a two-hidden-layer neural network $\mathcal{M}(2)$, each function $v \in \mathcal{M}(2)$ has the form of

$$140 \quad v(\mathbf{x}) = \mathbf{c}_1 \left(\mathbf{N}^{(2)} \circ \mathbf{N}^{(1)}(\mathbf{x})\right) + c_0 = \sum_{i=0}^{n_2} c_i \varphi_i^{(2)}(\mathbf{x}),$$

141 where $\left\{\varphi_i^{(2)}(\mathbf{x})\right\}_{i=1}^{n_2}$ are similarly defined as in (2.6) by

$$142 \quad \varphi_0^{(2)}(\mathbf{x}) = 1 \quad \text{and} \quad \varphi_i^{(2)}(\mathbf{x}) = \sigma\left(\omega_i^{(2)}\left(\mathbf{N}^{(1)}(\mathbf{x})\right) + b_i^{(2)}\right) = \sigma\left(\omega_i^{(2)}\sigma\left(\omega^{(1)}\mathbf{x} + \mathbf{b}^{(1)}\right) + b_i^{(2)}\right)$$

143 for $i = 1, \dots, n_2$ and nonlinear parameters are given by

$$144 \quad (3.4) \quad \Theta^{(2)} = \Theta^{(1)} \cup \left\{\theta_i^{(2)}\right\}_{i=1}^{n_2} := \Theta^{(1)} \cup \left\{\left(b_i^{(2)}, \omega_i^{(2)}\right)\right\}_{i=1}^{n_2}.$$

145 Similar to the shallow network $\mathcal{M}(1)$, denote pre-activation functions of neurons at the 2^{nd} hidden-
 146 layer by

$$147 \quad (3.5) \quad g_i^{(2)}(\mathbf{x}) = \omega_i^{(2)}\sigma\left(\omega^{(1)}\mathbf{x} + \mathbf{b}^{(1)}\right) + b_i^{(2)} \quad \text{for } i = 1, \dots, n_2$$

148 and their zero level sets, called the *breaking poly-hyper-planes*, by

$$149 \quad (3.6) \quad \mathcal{P}_i^{(2)}\left(\Theta^{(1)}, \theta_i^{(2)}\right) = \left\{\mathbf{x} \in \Omega : \omega_i^{(2)}\sigma\left(\omega^{(1)}\mathbf{x} + \mathbf{b}^{(1)}\right) + b_i^{(2)} = 0\right\}.$$

150 **REMARK 3.3.** Note that $g_i^{(2)}(\mathbf{x})$ is a single-valued, continuous piece-wise linear function. This
 151 fact implies that $\mathcal{P}_i^{(2)}\left(\Theta^{(1)}, \theta_i^{(2)}\right)$ as a zero level set is either empty or consists of poly-hyper-planes
 152 that do not intersect. Here, the poly-hyper-plane means a continuous hyper-plane that is composed
 153 of one or more connected hyper-plane segments. Moreover, each poly-hyper-plane is either closed
 154 or from part of the boundary to another part of the boundary.

155 **REMARK 3.4.** The physical partition $\mathcal{K}^{(2)}(v)$ is the refinement of the partition $\mathcal{K}^{(1)}(v)$ by using
 156 the breaking poly-hyper-planes $\left\{\mathcal{P}_i^{(2)}\left(\Theta^{(1)}, \theta_i^{(2)}\right)\right\}_{i=1}^{n_2}$.

157 In one dimension, $\mathcal{K}^{(2)}(v)$ is the refinement of $\mathcal{K}^{(1)}(v)$ by adding the 2^{nd} layer breaking points
 158 (red crosses in Fig. 1(a)) satisfying

$$159 \quad \sum_{j=1}^{n_2} \omega_{ij}^{(2)}\sigma\left(\omega_j^{(1)}x + b_j^{(1)}\right) + b_i^{(2)} = 0 \quad \text{for } i = 1, \dots, n_2.$$

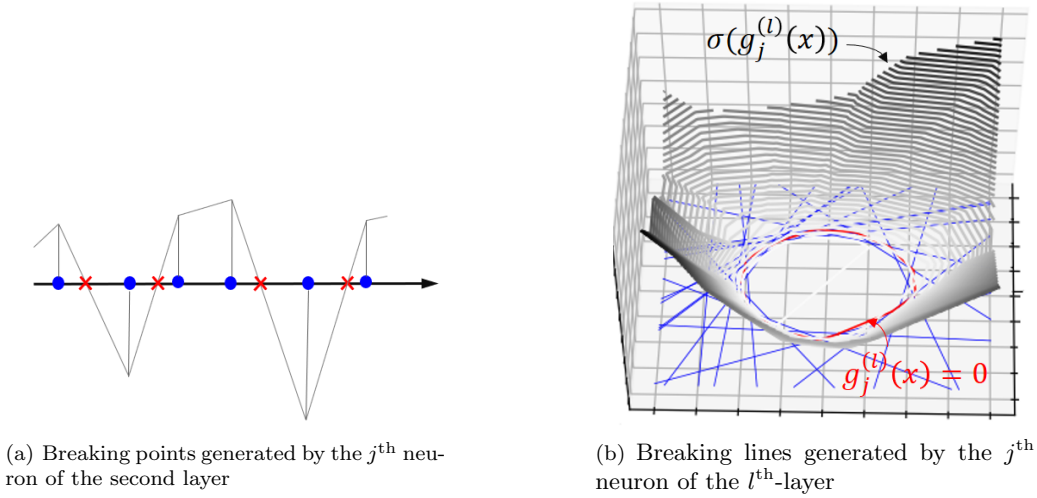


FIG. 1. Breaking points/lines in the first two hidden layers.

160 In two dimensions, the $\mathcal{K}^{(2)}(v)$ is the refinement of $\mathcal{K}^{(1)}(v)$ by adding the 2^{nd} layer breaking
 161 poly-lines (red poly-lines in Figs. 1(b) and 2(d)) satisfying

$$162 \quad \sum_{j=1}^{n_2} \omega_{ij}^{(2)} \sigma(\omega_j^{(1)} \mathbf{x} + b_j^{(1)}) + b_i^{(2)} = 0 \quad \text{for } i = 1, \dots, n_2.$$

163 For $k = 1, \dots, l-1$, denote the nonlinear parameter of the first k hidden-layers of $v(\mathbf{x})$ by

$$164 \quad (3.7) \quad \Theta^{(k)} = \Theta^{(k-1)} \cup \left\{ \theta_i^{(k)} \right\}_{i=1}^{n_k} := \Theta^{(k-1)} \cup \left\{ \left(b_i^{(k)}, \omega_i^{(k)} \right) \right\}_{i=1}^{n_k}.$$

165 Let $\mathcal{K}^{(k)}(v)$ denote the physical partition determined by the nonlinear parameters $\Theta^{(k)}$. Then
 166 the physical partition $\mathcal{K}^{(l)}(v)$ may be described through a refinement process starting from the
 167 physical partition $\mathcal{K}^{(1)}(v)$. For $k = 2, \dots, l$, the physical partition $\mathcal{K}^{(k)}(v)$ is the refinement of the
 168 previous physical partition $\mathcal{K}^{(k-1)}(v)$ by adding the following poly-hyper-planes

$$169 \quad (3.8) \quad \mathcal{P}_i^{(k)} \left(\Theta^{(k-1)}, \theta_i^{(k)} \right) = \left\{ \mathbf{x} \in \mathbb{R}^d : g_i^{(k)}(\mathbf{x}) = 0 \right\} \quad \text{for } i = 1, \dots, n_k,$$

170 where $g_i^{(k)}(\mathbf{x})$ is the pre-activation function of the i^{th} neuron at the k^{th} hidden-layer given by

$$171 \quad (3.9) \quad g_i^{(k)}(\mathbf{x}) = \omega_i^{(k)} \left(\mathbf{N}^{(k-1)} \circ \dots \circ \mathbf{N}^{(1)}(\mathbf{x}) \right) + b_i^{(k)}.$$

172 The procedure for determining the physical partition $\mathcal{K}^{(l)}(v)$ of the domain Ω involves calculating
 173 the arrangement of a domain formed by a set of hyper-planes and poly-hyper-planes. This may be
 174 computationally expensive, especially when the input dimension d is high.

175 In practice, computation is usually done over a set of points in Ω , e.g., the input data set
 176 $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^M$ for problem (1.1) and integration point set as in [6, 10]. This motivates introduction
 177 of the *data physical partition*, i.e., the physical partition of $v(\mathbf{x}) \in \mathcal{M}(l)$ with respect to a given
 178 data set $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^M$. In a similar fashion as Definition 3.1, we define the data physical partition
 179 $\mathcal{D}^{(k)}$ as follows

$$180 \quad (3.10) \quad \mathcal{D}^{(k)} = \mathcal{D} \cap \mathcal{K}^{(k)} = \left\{ \mathcal{D} \cap K : K \in \mathcal{K}^{(k)} \right\}$$

181 for $k = 1, \dots, l$.

182 Next, we describe how to form $\mathcal{D}^{(l)}$ for a given nonlinear parameters $\Theta^{(l)}$. To this end, let
 183 $H(t)$ be the Heaviside step function given by

$$184 \quad H(t) = \begin{cases} 1, & t > 0, \\ 0, & t < 0. \end{cases}$$

185 For $k = 1, \dots, l$, introduce vector-valued *layer indicator function* $\mathbf{I}^{(k)} : \mathbb{R}^d \rightarrow \mathbb{R}^{n_k}$ as

$$186 \quad (3.11) \quad \mathbf{I}^{(k)}(\mathbf{x}) = H\left(\mathbf{g}^{(k)}(\mathbf{x})\right),$$

187 where $\mathbf{g}^{(k)}(\mathbf{x}) = \left(g_i^{(k)}(\mathbf{x})\right)_{n_k \times 1}$ is the pre-activation function defined in (3.9) and application of H
 188 to a vector-valued function is defined component-wisely. For a given nonlinear parameter $\Theta^{(l)}$, we
 189 define the *network indicator function* by

$$190 \quad \mathcal{I}^{(l)}(\mathbf{x}) = \left(\mathbf{I}^{(1)}(\mathbf{x}), \dots, \mathbf{I}^{(l)}(\mathbf{x})\right).$$

191 Let the data physical partition $\mathcal{D}^{(l)}$ be of the form

$$192 \quad (3.12) \quad \mathcal{D}^{(l)} = \left\{ D_j^{(l)} \right\}_{j=1}^{m_l},$$

193 where m_l is the number of disjoint elements of the data physical partition $\mathcal{D}^{(l)}$ and each element
 194 $D_j^{(l)}$ is a subset of the input data set \mathcal{D} such that the value of the network indicator function $\mathcal{I}^{(l)}$
 195 is same for all points in $D_j^{(l)}$. Denote this value by $\mathcal{I}_{D_j^{(l)}}^{(l)}$, then we have

$$196 \quad (3.13) \quad \mathcal{I}^{(l)}(\mathbf{x}) = \mathcal{I}_{D_j^{(l)}}^{(l)} \quad \text{for } \mathbf{x} \in D_j^{(l)}.$$

197 For each element $D \in \mathcal{D}^{(l)}$, denote the centroid of D by

$$198 \quad (3.14) \quad \mathbf{x}_D = \frac{1}{|D|} \sum_{\mathbf{x}_i \in D} w_i \mathbf{x}_i,$$

199 and the covariance matrix of D formed by vectors $\mathbf{x}_i - \mathbf{x}_D$ for all $\mathbf{x}_i \in D$ by

$$200 \quad (3.15) \quad \text{CoV}_D = \sum_{\mathbf{x}_i \in D} [\mathbf{x}_i - \mathbf{x}_D]^T [\mathbf{x}_i - \mathbf{x}_D],$$

201 where $\mathbf{x}_i - \mathbf{x}_D$ is a d -dimensional row vector. Then each element $D \in \mathcal{D}^{(l)}$ has d principal directions
 202 that correspond to the eigenvectors of CoV_D .

203 **4. Adaptive Network Enhancement Method.** This section describes the adaptive net-
 204 work enhancement method (ANE) for problem (1.1).

205 To this end, denote the current neural network, approximation, and error at the k^{th} adaptive
 206 step by

$$207 \quad \mathcal{M}^{(k)}(l_k), \quad f^{(k)}(\mathbf{x}), \quad \text{and} \quad \xi^{(k)} = L\left(f^{(k)}\right),$$

208 respectively, where l_k is the number of hidden-layers of the network $\mathcal{M}^{(k)}(l_k)$. When accuracy of
 209 the current approximation $f^{(k)}(\mathbf{x})$ is not within the prescribed tolerance, i.e., $\xi^{(k)} > \epsilon L(0)$, the

210 network $\mathcal{M}^{(k)}(l_k)$ is enhanced by adding neurons at the either l_k -th or $(l_k + 1)$ -th hidden-layer.
 211 The latter means that we starts a new hidden-layer.

212 To determine the number of neurons to be added at the l_k -th hidden-layer, we use the *local*
 213 *network enhancement strategy* based on the data physical partition of $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^M$:

$$214 \quad \mathcal{D}^{(l_k)}(f^{(k)}) = \mathcal{D} \cap \mathcal{K}^{(l_k)}(f^{(k)}) = \left\{ \mathcal{D} \cap K : K \in \mathcal{K}^{(l_k)}(f^{(k)}) \right\}$$

215 by the current approximation $f^{(k)}(\mathbf{x})$. Specifically, we divide $\mathcal{D}^{(l_k)}(f^{(k)})$ into two disjoint subsets,

$$216 \quad \mathcal{D}^{(l_k)}(f^{(k)}) = \hat{\mathcal{D}}^{(l_k)}(f^{(k)}) \cup \left(\mathcal{D}^{(l_k)}(f^{(k)}) \setminus \hat{\mathcal{D}}^{(l_k)}(f^{(k)}) \right)$$

217 where $\hat{\mathcal{D}}^{(l_k)}(f^{(k)}) = \mathcal{D} \cap \hat{\mathcal{K}}^{(l_k)}(f^{(k)})$ is a subset of $\mathcal{D}^{(l_k)}(f^{(k)})$ consisting of elements in $\mathcal{D}^{(l_k)}(f^{(k)})$
 218 such that $f^{(k)}(\mathbf{x})$ is not yet a good approximation. Then the *enhancement strategy* is to add
 219 $|\hat{\mathcal{D}}^{(l_k)}(f^{(k)})|$ new neurons to the l_k -th hidden-layer, where

$$220 \quad (4.1) \quad |\hat{\mathcal{D}}^{(l_k)}(f^{(k)})| = \text{the number of elements of } \hat{\mathcal{D}}^{(l_k)}(f^{(k)}).$$

221 To generate $\hat{\mathcal{D}}^{(l_k)}(f^{(k)})$, we employ the so-called marking strategy. There are two commonly
 222 used marking strategies for adaptive mesh refinement. One is the average marking strategy and
 223 the other is the bulk marking strategy. To describe these marking strategies, let us first introduce
 224 the following local error indicator

$$225 \quad (4.2) \quad \xi_D^{(k)} = \left(\sum_{\mathbf{x}_i \in D} w_i \left(f^{(k)}(\mathbf{x}_i) - y_i \right)^2 \right)^{1/2}$$

226 for each element $D \in \mathcal{D}^{(l_k)}(f^{(k)})$. Clearly, we have

$$227 \quad (4.3) \quad \xi^{(k)} = \left(\sum_{i=1}^M w_i \left(f^{(k)}(\mathbf{x}_i) - y_i \right)^2 \right)^{1/2} = \left(\sum_{D \in \mathcal{D}^{(l_k)}(f^{(k)})} \left(\xi_D^{(k)} \right)^2 \right)^{1/2}.$$

228 The average marking strategy is given by

$$229 \quad (4.4) \quad \hat{\mathcal{D}}^{(l_k)}(f^{(k)}) = \left\{ D \in \mathcal{D}^{(l_k)}(f^{(k)}) : \xi_D^{(k)} \geq \frac{1}{|\mathcal{D}^{(l_k)}(f^{(k)})|} \sum_{D \in \mathcal{D}^{(l_k)}(f^{(k)})} \xi_D^{(k)} \right\}.$$

230 The bulk marking strategy is to find a minimal subset $\hat{\mathcal{D}}^{(l_k)}(f^{(k)})$ such that

$$231 \quad (4.5) \quad \sum_{D \in \hat{\mathcal{D}}^{(l_k)}(f^{(k)})} \left(\xi_D^{(k)} \right)^2 \geq \gamma_1 \sum_{D \in \mathcal{D}^{(l_k)}(f^{(k)})} \left(\xi_D^{(k)} \right)^2 \quad \text{for } \gamma_1 \in (0, 1).$$

232 The enhancement strategy adding $|\hat{\mathcal{D}}^{(l_k)}(f^{(k)})|$ new neurons is suitable for all hidden-layers.
 233 Nevertheless, it may not be efficient for hidden-layers beyond the first hidden-layer. Notice that a
 234 multi-layer network is capable of generating piece-wise breaking hyper-planes in connected subdo-
 235 mains by one neuron. This observation motivates the notion of the *reduced number of elements* in
 236 $\hat{\mathcal{D}}^{(l_k)}(f^{(k)})$. To this end, let

$$237 \quad (4.6) \quad \hat{\mathcal{D}}^{(l)}(f^{(k)}) = \left\{ \hat{\mathcal{D}}_j^{(l)} \right\}_{j=1}^{\hat{m}_l}.$$

238 That is, there are \hat{m}_l marked elements in $\mathcal{D}^{(l)}(f^{(k)})$. Any two elements in $\hat{\mathcal{D}}^{(l_k)}(f^{(k)})$ are said
 239 to be disconnected if there is no pass connecting these two elements by elements of $\hat{\mathcal{D}}^{(l_k)}(f^{(k)})$.
 240 Let us group connected elements of $\hat{\mathcal{D}}^{(l)}(f^{(k)})$ to form a set, whose elements are disconnected,
 241 denoted by

$$242 \quad (4.7) \quad \tilde{\mathcal{D}}^{(l)}(f^{(k)}) = \left\{ \tilde{D}_j^{(l)} \right\}_{j=1}^{\tilde{m}_l},$$

243 where each element $\tilde{D}_j^{(l)} \in \tilde{\mathcal{D}}^{(l)}(f^{(k)})$ is either an element of $\hat{\mathcal{D}}^{(l)}(f^{(k)})$ or a union of connected
 244 elements in $\hat{\mathcal{D}}^{(l)}(f^{(k)})$. Obviously, $\tilde{m}_l \leq \hat{m}_l$. Now, we define the reduced number of elements in
 245 $\hat{\mathcal{D}}^{(l_k)}(f^{(k)})$ by

$$246 \quad (4.8) \quad \left| \hat{\mathcal{D}}^{(l_k)}(f^{(k)}) \right|_r = \begin{cases} \hat{m}_l, & l_k = 1, \\ \tilde{m}_l, & l_k \geq 2. \end{cases}$$

247 where any two elements in $\hat{\mathcal{D}}^{(l_k)}(f^{(k)})$ are disjoint if there is no pass connecting these two elements
 248 by elements of $\hat{\mathcal{D}}^{(l_k)}(f^{(k)})$.

249 **REMARK 4.1.** For any two elements in $\hat{\mathcal{D}}^{(l_k)}(f^{(k)})$, if values of their network indicator func-
 250 tion differ only for one neuron, e.g., the i^{th} neuron at the k^{th} hidden-layer, then these two elements
 251 are neighbor and share part of the poly-hyper-plane $\mathcal{P}_i^{(k)}(\Theta^{(k-1)}, \theta_i^{(k)})$ defined in (3.8).

252 To address question (b) in section 1, i.e., when to add a new hidden-layer, we introduce a
 253 computable quantity, referred as the *improvement rate*, defined by

$$254 \quad (4.9) \quad \eta_r^{(k)} = \left(\frac{\xi^{(k-1)} - \xi^{(k)}}{\xi^{(k-1)}} \right) / \left(\frac{(M^{(k)}(l_k))^r - (M^{(k-1)}(l_{k-1}))^r}{(M^{(k)}(l_k))^r} \right),$$

255 where $M^{(k-1)}(l_{k-1})$ and $M^{(k)}(l_k)$ denote the numbers of parameters of the networks $\mathcal{M}^{(k-1)}(l_{k-1})$
 256 and $\mathcal{M}^{(k)}(l_k)$, respectively; and r is the order of the approximation with respect to the number
 257 of parameters and may depend on the activation function and the layer. The improvement rate
 258 measure a rate of improvement of two consecutive networks per the relative increase of parameters.
 259 If the improvement rate $\eta_r^{(k)}$ is less than or equal to a prescribed expectation rate $\delta \in (0, 2)$, i.e.,

$$260 \quad (4.10) \quad \eta_r^{(k)} \leq \delta,$$

261 for two consecutive adaptive steps, then the ANE adds a new hidden-layer. Otherwise, the ANE
 262 adds neurons to the l_k -th hidden-layer of the current network $\mathcal{M}^{(k)}(l_k)$.

263 The ANE method for generating a nearly minimal multi-layer neural network is described in
 264 Algorithm 3.1.

265 **5. Initialization of training.** This section discusses initialization strategies of parameters
 266 of neural network in two dimensions. Extensions to three dimensions are straightforward.

267 The optimization problem in Step (7) of Algorithm 3.1 is non-convex and, hence, computa-
 268 tionally intensive and complicated. Currently, this problem is often solved by either the first- or
 269 second-order iterative optimization methods such as gradient-based methods or Newton-like meth-
 270 ods (see survey papers [2, 3] and references therein). Since non-convex optimizations usually have
 271 many solutions and/or many local minimums, it is then critical to start with a good initial guess
 272 in order to obtain the desired solution.

273 The ANE method itself is a natural continuation process for generating good initialization.
 274 That is, the approximation $f^{(k)}(\mathbf{x})$ of the previous network $\mathcal{M}^{(k)}(l_k)$ is in general a good approxi-
 275 mation to $f^{(k+1)}(\mathbf{x})$ defined in Step (7) of Algorithm 3.1 for the enhanced network $\mathcal{M}^{(k+1)}(l_{k+1})$.

Algorithm 3.1 Adaptive Network Enhancement.

Given a data set $\{(\mathbf{x}_i, y_i)\}_{i=1}^M$ with $\mathbf{x}_i \in \Omega = [-1, 1]^d$, positive weights $\{w_i\}_{i=1}^M$, and a tolerance $\epsilon > 0$ for accuracy, starting with a one hidden-layer network $\mathcal{M}^{(0)}(l_0)$ with a small number of neurons, compute $f^{(0)} = \arg \min_{v \in \mathcal{M}^{(0)}(l_0)} L(v)$ by an iterative solver, then for $k = 0, 1, 2, \dots$,

- (1) use the network indicator function to determine the data physical partition $\mathcal{D}^{(l_k)}(f^{(k)})$;
- (2) for each $D \in \mathcal{D}^{(l_k)}(f^{(k)})$, compute the local indicator $\xi_D^{(k)}$ in (4.2) and the estimator $\xi^{(k)}$ in (4.3);
- (3) if $\xi < \epsilon$, then stop; otherwise, go to Step (4);
- (4) use a marking strategy to form the subset $\hat{\mathcal{D}}^{(l_k)}(f^{(k)})$ and calculate $\left| \hat{\mathcal{D}}^{(l_k)}(f^{(k)}) \right|_r$;
- (5) for a prescribed expectation rate $\delta \in (0, 2)$, if (4.10) holds for two consecutive steps, then set $l_{k+1} = l_k + 1$; otherwise, set $l_{k+1} = l_k$;
- (6) form network $\mathcal{M}^{(k+1)}(l_{k+1})$ by adding $\left| \hat{\mathcal{D}}^{(l_k)}(f^{(k)}) \right|_r$ new neurons to the l_{k+1} -th hidden-layer;
- (7) compute $f^{(k+1)} = \arg \min_{v \in \mathcal{M}^{(k+1)}(l_{k+1})} L(v)$ by an iterative solver.

276 Therefore, the trained nonlinear parameters of $\mathcal{M}^{(k)}(l_k)$ for $f^{(k)}(\mathbf{x})$ are good initials for the corre-
 277 sponding nonlinear parameters of the enhanced network $\mathcal{M}^{(k+1)}(l_{k+1})$. Based on this observation,
 278 below we discuss our initialization strategies for (1) parameters of the network $\mathcal{M}^{(0)}(l_0)$, (2) pa-
 279 rameters of newly added neurons, and (3) linear (output) parameters of $f^{(k+1)}$.

280 Starting with a one hidden-layer network $\mathcal{M}^{(0)}(l_0)$ with relatively small number n_{l_0} of neurons,
 281 the approximation $f^{(0)}(\mathbf{x})$ has of the form

$$282 \quad f^{(0)}(\mathbf{x}) = \mathbf{c}_1 \sigma(\boldsymbol{\omega}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}) + c_0,$$

283 where $\boldsymbol{\Theta}^{(l_0)} = \left\{ \boldsymbol{\theta}_i^{(1)} \right\}_{i=1}^{n_{l_0}} := \left\{ \left(b_i^{(1)}, \boldsymbol{\omega}_i^{(1)} \right) \right\}_{i=1}^{n_1}$ are nonlinear parameters and $\mathbf{c} = (c_0, \mathbf{c}_1) \in \mathbb{R}^{n_{l_0}+1}$
 284 are linear parameters. Initial of $\boldsymbol{\Theta}^{(l_0)}$ is chosen such that the hyper-lines

$$285 \quad \mathcal{P}_i^{(1)}(\boldsymbol{\theta}_i^{(1)}) : \boldsymbol{\omega}_i^{(1)} \mathbf{x} + b_i^{(1)} = 0 \quad \text{for } i = 1, \dots, n_1$$

286 partition the domain $\Omega = (0, 1)^2$ uniformly. Initial of \mathbf{c} is set to be the solution of the system of
 287 linear algebraic equations

$$288 \quad (5.1) \quad \mathbf{M}^{(l_0)}(\boldsymbol{\Theta}^{(l_0)}) \mathbf{c} = \mathbf{F}^{(l_0)}(\boldsymbol{\Theta}^{(l_0)})$$

289 defined in a similar fashion as (2.9).

290 Next, we discuss how to initialize the biases and weights of newly added neurons of the network
 291 $\mathcal{M}^{(k+1)}(l_{k+1})$. There are three cases:

$$292 \quad (1) l_{k+1} = 1, \quad (2) l_{k+1} = l_k + 1, \quad \text{and} \quad (3) l_{k+1} = l_k \geq 2.$$

293 Case (1) means that the new neurons are added at the first hidden-layer. By associating each new
 294 neuron with an element $D \in \mathcal{D}^{(l_k)}(f^{(k)})$, we initialize this neuron by setting its corresponding
 295 breaking line to pass through the centroid \mathbf{x}_D and orthogonal to the principal direction that
 296 corresponds to the smallest eigenvalue of the covariance matrix CoV_D .

307 Consider Case (2). When $l_{k+1} = l_k + 1$, we start a new hidden-layer with $\left| \hat{\mathcal{D}}^{(l_k)}(f^{(k)}) \right|_r$
 308 neurons. By the definition in (4.8), we associate each neuron at the new hidden-layer $l_{k+1} = l_k + 1$
 309 with an isolated element or an element consisting of several connected elements in $\hat{\mathcal{D}}^{(l_k)}(f^{(k)})$ and
 310 denote its bias and weights by

$$301 \quad (5.2) \quad \boldsymbol{\theta}^{(l_{k+1})} = \left(b^{(l_{k+1})}, \boldsymbol{\omega}^{(l_{k+1})} \right) \in \mathbb{R}^{n_{l_k}+1} = \left(b^{(l_{k+1})}, \omega_1^{(l_{k+1})}, \dots, \omega_{n_{l_k}}^{(l_{k+1})} \right) \in \mathbb{R}^{n_{l_k}+1}.$$

302 As section 2, denote the corresponding pre-activation function of the neuron by

$$303 \quad g^{(l_{k+1})}(\mathbf{x}) = \boldsymbol{\omega}^{(l_{k+1})} \left(\mathbf{N}^{(l_k)} \circ \dots \circ \mathbf{N}^{(1)}(\mathbf{x}) \right) + b^{(l_{k+1})}$$

304 If the corresponding element D is an isolated element in $\hat{\mathcal{D}}^{(l_k)}(f^{(k)})$, let $l_D(\mathbf{x}) = 0$ be the line
 305 that passes through the centroid \mathbf{x}_D of D and is orthogonal to the direction vector with the lowest
 306 variance of D (see section 3). Denote by \mathbf{x}_d the projection of a point in D onto the line $l_D(\mathbf{x}) = 0$
 307 and whose distance to \mathbf{x}_D is the largest among projections of all points in D onto the line. Then
 308 initial $\boldsymbol{\theta}_D^{(l_{k+1})}$ of the parameter $\boldsymbol{\theta}^{(l_{k+1})}$ is set to be

$$309 \quad (5.3) \quad \boldsymbol{\theta}_D^{(l_{k+1})} = \arg \min_{\boldsymbol{\theta}^{(l_{k+1})} \in \mathbb{R}^{n_{l_k}+1}} \left\{ \left(g^{(l_{k+1})}(\mathbf{x}_D) \right)^2 + \left(g^{(l_{k+1})}(\mathbf{x}_d) \right)^2 \right\}.$$

310 When the corresponding element D consists of several connected elements in $\hat{\mathcal{D}}^{(l_k)}(f^{(k)})$, denote
 311 the collection of these connected elements by \mathcal{C} . For each element $C \in \mathcal{C}$, denote by \mathbf{x}_C the centroid
 312 of C . Then initial $\boldsymbol{\theta}_C^{(l_{k+1})}$ of the parameter $\boldsymbol{\theta}^{(l_{k+1})}$ is set to be

$$313 \quad (5.4) \quad \boldsymbol{\theta}_C^{(l_{k+1})} = \arg \min_{\boldsymbol{\theta}^{(l_{k+1})} \in \mathbb{R}^{n_{l_k}+1}} \sum_{C \in \mathcal{C}} \left(g^{(l_{k+1})}(\mathbf{x}_C) \right)^2.$$

314 Now, let us consider Case (3) where new neurons are added at the current layer $l_{k+1} = l_k$. Let
 315 $s \in \{1, \dots, k-1\}$ be the largest integer such that $l_{k-s} = l_k - 1$. Then $\mathcal{M}^{(k-s)}(l_{k-s})$ is the final
 316 network with $l_{k-s} = l_k - 1$ hidden-layers. Hence the weights and bias of each neuron associated
 317 with an element in $\hat{\mathcal{D}}^{(l_k)}(f^{(k)})$ has the form of

$$318 \quad (5.5) \quad \boldsymbol{\theta}^{(l_{k+1})} = \left(b^{(l_{k+1})}, \boldsymbol{\omega}^{(l_{k+1})} \right) = \left(b^{(l_{k+1})}, \omega_1^{(l_{k+1})}, \dots, \omega_{n_{l_{k-s}}}^{(l_{k+1})} \right) \in \mathbb{R}^{n_{l_{k-s}}+1}.$$

319 Initial of $\boldsymbol{\theta}^{(l_{k+1})}$ in (5.5) can then be defined in a similar fashion as Case (2). Specifically, we have

$$320 \quad (5.6) \quad \begin{cases} \boldsymbol{\theta}_D^{(l_{k+1})} = \arg \min_{\boldsymbol{\theta}^{(l_{k+1})} \in \mathbb{R}^{n_{l_{k-s}}+1}} \left\{ \left(g^{(l_{k+1})}(\mathbf{x}_D) \right)^2 + \left(g^{(l_{k+1})}(\mathbf{x}_d) \right)^2 \right\} & \text{and} \\ \boldsymbol{\theta}_C^{(l_{k+1})} = \arg \min_{\boldsymbol{\theta}^{(l_{k+1})} \in \mathbb{R}^{n_{l_{k-s}}+1}} \sum_{C \in \mathcal{C}} \left(g^{(l_{k+1})}(\mathbf{x}_C) \right)^2, \end{cases}$$

321 where \mathbf{x}_D , \mathbf{x}_d , and \mathbf{x}_C are defined in a similar way as in Case (2).

322 Finally, initial of the linear parameter $\mathbf{c} = (c_0, \mathbf{c}_1) \in \mathbb{R}^{n_{l_{k+1}}+1}$ of $f^{k+1}(\mathbf{x})$ is set to be the
 323 solution of the system of algebraic linear equations

$$324 \quad (5.7) \quad \mathbf{M}^{(l_{k+1})} \left(\boldsymbol{\Theta}^{(l_{k+1})} \right) \mathbf{c} = \mathbf{F}^{(l_{k+1})} \left(\boldsymbol{\Theta}^{(l_{k+1})} \right)$$

325 defined in a similar fashion as (2.9).

326 **6. Numerical Experiment.** In this section, we report the numerical experiment on using
 327 the ANE method to approximate a function using the least-squares loss. The target function is
 328 defined on the domain $\Omega = [-1, 1]^2$, and is given by

$$329 \quad (6.1) \quad f(x, y) = \tanh\left(\frac{1}{\alpha}(x^2 + y^2 - \frac{1}{4})\right) - \tanh\left(\frac{3}{4\alpha}\right).$$

330 For small constant α , this function exhibits a sharp transitional layer across a circular interface.

331 For this experiment, we set a small $\alpha = 0.01$ to test approximation accuracy using ANE. and
 332 the corresponding target function f is depicted in Fig. 2(a). A data set \mathcal{D} for training network is
 333 generated using a fixed set 200×200 of quadrature points that are uniformly distributed in the
 334 domain Ω .

335 During the ANE process, we adopt the bulk marking strategy defined in (4.5) with $\gamma_1 = 0.5$
 336 and choose the expectation rate $\delta = 0.6$ with $r = 1$ in (4.9); and the expected precision $\epsilon =$
 337 0.05 . The ANE method started with an initial network of 12 neurons in one hidden layer. The
 338 corresponding breaking lines $\{\mathcal{P}_i\}_{i=1}^{12}$ of these 12 neurons were uniformly initialized within the
 339 domain. Specifically, half of breaking lines are parallel to the x -axis

$$340 \quad \omega_i^{(1)} = (0, 1) \quad \text{and} \quad b_i^{(1)} = -1 + \frac{1}{3}i \quad \text{for} \quad i = 0, \dots, 5$$

341 and the other half are parallel to the y -axis

$$342 \quad \omega_i^{(1)} = (1, 0) \quad \text{and} \quad b_i^{(1)} = -1 + \frac{1}{3}(i - 6) \quad \text{for} \quad i = 6, \dots, 12.$$

343 In addition, the output weights and bias are initialized by solving the linear system in (5.1).

344 For each iteration of the ANE process, the corresponding minimization problem in (1.1) is
 345 solved iteratively using the Adam version of gradient descent [9] with a fixed learning rate 0.005.
 346 Adam’s iterative solver is terminated when the relative change of the loss function $\|f - \hat{f}\|_\tau$ is less
 347 than 10^{-3} per 2000 iterations.

TABLE 1
 Numerical results for using ANE to approximate function with a circular transitional layer

NN structure	# parameters	Approximation accuracy $\ f - \hat{f}\ _\tau / \ f\ $	Improvement rate η
2-12-1	37	0.357414	–
2-18-1	55	0.323118	0.293198
2-26-1	93	0.272614	0.382528
2-18-5-1	137	0.025483	1.538967

348 The ANE process is automatically terminated after four loops (see Table 1 for the interme-
 349 diate and final result), and the final network model generated by the ANE is 2-18-5-1¹ with 137
 350 parameters. The final network approximation model and the corresponding physical partition are
 351 shown in Figs. 2 (e) and (d). Using a relatively small set of parameters, ANE is able to accurately
 352 approximate a function with a thin transition layer without any oscillations. This remarkable
 353 approximation property can be explained by the fact that the circular interface of the underlying
 354 function is captured very effectively by a few breaking poly-lines generated in *the second hidden*
 355 *layer*, see the closed breaking lines formed by the 5 neurons in the second hidden layer in Fig. 2
 356 (d).

¹The structure of a two- or three-hidden-layer network is expressed as 2- n_1 -1 or 2- n_1 - n_2 -1, respectively, where n_1 and n_2 are the number of neurons at the first and second hidden-layer.

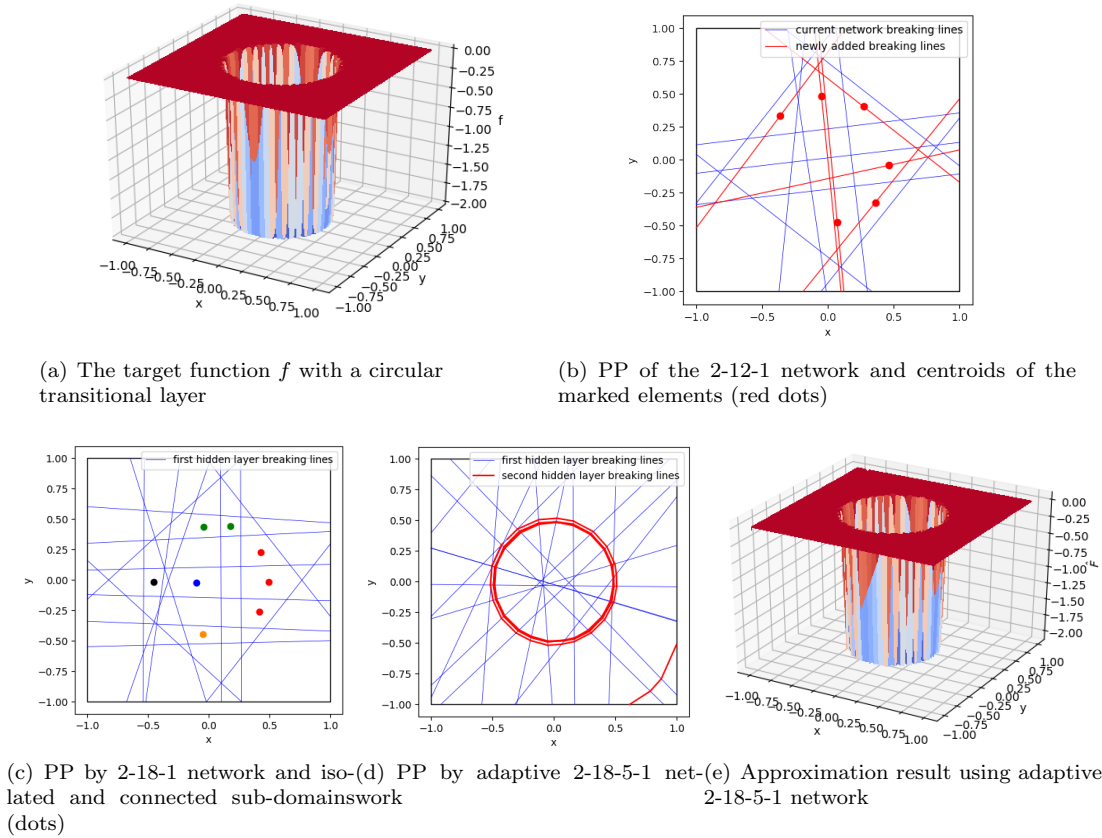


FIG. 2. Adaptive approximation results for function with a transitional layer

357 Figs. 2 (b)-(c) plot the physical partitions of the NN models at the intermediate adaptive
 358 process. In Fig. 2 (b), the centroids of the marked elements are illustrated by red dots; the
 359 breaking lines corresponding to the current and newly added neurons are shown by blue and red
 360 lines, respectively. Notice that the newly added neurons are initialized with break lines that pass
 361 through the centroids and align with the principal directions of the marked elements. Fig. 2(c)
 362 shows that there are 8 marked elements and 5 disjoint elements, which explains that 5 neurons are
 363 added to the second hidden layer during the neuron enhancement step.

TABLE 2
 Numerical results of adaptive and fixed networks for function with a transitional layer

Network structure	# parameters	Approximation accuracy $\ f - \hat{f}\ _{\tau} / \ f\ $
2-18-5-1 (Adaptive)	137	2.5483 %
2-18-5-1 (Fixed)	137	4.6199%
2-174-1 (Fixed)	523	11.1223%

364 For the purpose of a comparative study, we conducted function approximation experiments
 365 using two fixed network structures. As outlined in Table 2, when utilizing the same network
 366 structure (2-18-5-1), the resulting approximation accuracy is inferior to that achieved by the ANE

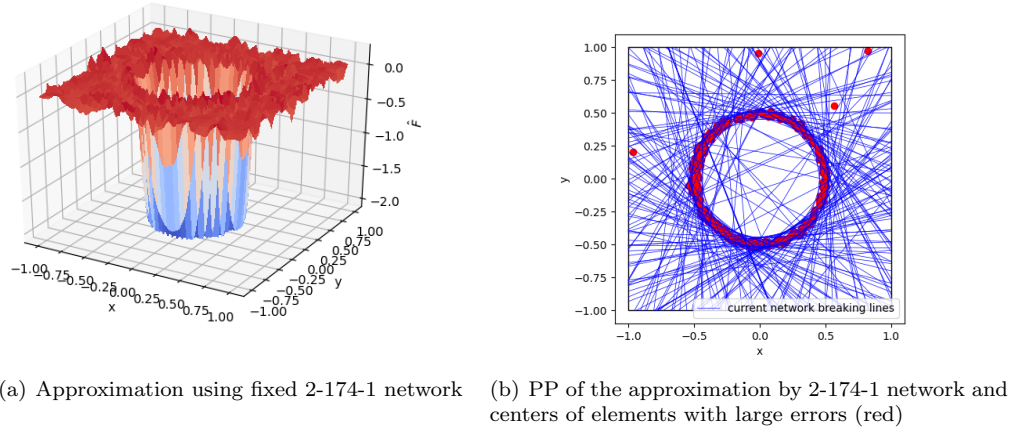


FIG. 3. Approximation results generated by a fixed 2-174-1 network for function with a transitional layer

367 method. The first two rows of Table 2 suggest that the ANE method provides a good initialization,
 368 which may simplify the non-convex optimization problem.

369 In the second experiment, we employed a fixed one-hidden-layer network (2-174-1) with nearly
 370 four times the number of parameters compared to the adaptive network. Despite the increased
 371 degrees of freedom and complexity, its approximation is less accurate (refer to the third row of
 372 Table 2). Furthermore, the approximated NN model exhibits a certain degree of oscillation (see
 373 Fig. 3 (a)), although the corresponding physical partition (Fig. 3(b)) still captures the narrow
 374 transition layer.

375 In general, a one-hidden-layer network necessitates dense breaking lines to approximate a cir-
 376 cular interface, and oscillations along the interface can be attributed to the global basis functions
 377 generated from the first hidden layer. This experiment highlights that a deeper network, as il-
 378 lustrated by the two-hidden-layer network in this example, is more efficient in approximating a
 379 function with a thin nonlinear transition layer or interface. This experimental observation aligns
 380 with the theoretical findings presented in [7].

381 **7. Conclusion.** Designing an optimal deep neural network for a given task is important
 382 and challenging in many machine learning applications. This chapter provides a comprehensive
 383 introduction to the adaptive network enhancement (ANE) method, proposed recently in [11, 10, 5],
 384 which generates a nearly optimal multi-layer neural network for a given task within some prescribed
 385 accuracy. This self-adaptive algorithm is based on the novel network enhancement strategies
 386 that determine when a new hidden-layer and how many new neurons should be added when
 387 the current network is not sufficient for the task. This adaptive algorithm learns not only from
 388 given information (data, function, partial differential equation) but also from the current computer
 389 simulation, and it is therefore a learning algorithm at a level which is more advanced than common
 390 machine learning algorithms.

391 The resulting non-convex optimization at each adaptive step is computationally intensive and
 392 complicated with possible many global/local minimums. The ANE method provides a natural
 393 process for obtaining a good initialization that assists training significantly. Moreover, to provide
 394 a better initial guess, this chapter discusses an advanced procedure for initializing newly added
 395 neurons at the current or next hidden-layer.

396 Functions and partial differential equations with sharp transitions or discontinuities at un-
 397 known location have been computationally challenging, when approximated using other functional
 398 classes such as polynomials or piece-wise polynomials with fixed meshes. It was demonstrated

399 numerically in [11, 10, 5] that the ANE method can automatically design a nearly minimal two-
400 or multi-hidden-layer network to learn functions exhibiting sharp transitional layers as well as
401 continuous/discontinuous solutions of partial differential equations.

402

REFERENCES

- 403 [1] E. L. Allgower and K. Georg. *Numerical Continuation Methods: An Introduction*. Springer-Verlag, Berlin,
404 1990.
- 405 [2] L. Bottou, F. E. Curtis, and J. Nocedal. Optimization methods for large-scale machine learning. *SIAM Review*,
406 60:223–311, 2018.
- 407 [3] L. Bottou, F. E. Curtis, and J. Nocedal. Optimization methods for large-scale machine learning. *SIAM Review*,
408 60(2):223–311, 2018.
- 409 [4] Z. Cai, J. Chen, and M. Liu. Least-squares ReLU neural network (LSNN) method for linear advection-reaction
410 equation. *Journal of Computational Physics*, 443 (2021) 110514.
- 411 [5] Z. Cai, J. Chen, and M. Liu. Self-adaptive deep neural network: Numerical approximation to functions and
412 PDEs. *J. Comput. Phys.*, 455 (2022) 111021.
- 413 [6] Z. Cai, J. Chen, M. Liu, and X. Liu. Deep least-squares methods: An unsupervised learning-based numerical
414 method for solving elliptic pdes. *J. Comput. Phys.*, 420 (2020) 109707.
- 415 [7] Z. Cai, J. Choi, and M. Liu. Least-squares neural network (LSNN) method for linear advection-reaction
416 equation: general discontinuous interface. *arXiv:2301.06156v3[math.NA]*, 2023.
- 417 [8] Z. Cai, T. Ding, M. Liu, X. Liu, and J. Xia. A damped block gauss-newton method for shallow relu neural
418 network. *manuscript*, 2023.
- 419 [9] D. P. Kingma and J. Ba. ADAM: A method for stochastic optimization. In *International Conference on*
420 *Representation Learning, San Diego*, 2015; arXiv preprint arXiv:1412.6980.
- 421 [10] M. Liu and Z. Cai. Adaptive two-layer relu neural network: II. Ritz approximation to elliptic PDEs. *Comput.*
422 *Math. Appl.*, 113:103–116, 2022.
- 423 [11] M. Liu, Z. Cai, and J. Chen. Adaptive two-layer ReLU neural network: I. best least-squares approximation.
424 *Comput. Math. Appl.*, 113:34–44, 2022.