

## A Complete Error Analysis of PINNs for Elliptic Equations Using Projected Stochastic Gradient Descent

Yuling Jiao<sup>2,3,5</sup>, Ruoxuan Li<sup>1</sup>, Defeng Sun<sup>6</sup>, Peiying Wu<sup>1</sup> and Jerry Zhijian Yang<sup>1,2,3,4,\*</sup>

<sup>1</sup> School of Mathematics and Statistics, Wuhan University, Wuhan 430072, P.R. China.

<sup>2</sup> National Center for Applied Mathematics in Hubei, Wuhan University, Wuhan 430072, P.R. China.

<sup>3</sup> Hubei Key Laboratory of Computational Science, Wuhan University, Wuhan 430072, P.R. China.

<sup>4</sup> Institute for Math & AI, Wuhan University, Wuhan 430072, P.R. China.

<sup>5</sup> School of Artificial Intelligence, Wuhan University, Wuhan 430072, P.R. China.

<sup>6</sup> Department of Applied Mathematics and Research Center for Intelligent Operations Research, The Hong Kong Polytechnic University, Hong Kong, P.R. China.

Received 19 April 2025; Accepted (in revised version) 29 August 2025

---

**Abstract.** Physics-informed neural networks (PINNs) have recently gained attention as a powerful and efficient tool for solving partial differential equations (PDEs). Despite their empirical success, the theoretical understanding of PINNs, especially in the context of over-parameterization, remains incomplete. This paper presents a complete error analysis of over-parameterized PINNs for elliptic equations using projected stochastic gradient descent (PSGD) optimization. Our analysis rigorously examines the interplay of approximation error, statistical error, and optimization error, offering a unified framework for understanding the convergence behavior of PINNs. By leveraging the properties of PSGD, we establish convergence rates and derive conditions on neural network architecture, training sample requirements, and optimization parameters to ensure specified accuracy.

**AMS subject classifications:** 65M15, 65N15, 65Y20

**Key words:** Physics-informed neural networks, projected stochastic gradient descent, complete error analysis, over-parameterization.

---

\*Corresponding author. *Email addresses:* yulingjiaomath@whu.edu.cn (Y. Jiao), ruoxuanli.math@whu.edu.cn (R. Li), defeng.sun@polyu.edu.hk (D. Sun), peiyingwu@whu.edu.cn (P. Wu), zjyang.math@whu.edu.cn (J. Z. Yang)

## 1 Introduction

Traditional numerical methods, such as the finite element method [9, 12], have proven to be very effective for solving low-dimensional PDEs. However, they encounter significant difficulties when applied to high-dimensional problems. The impressive success of deep learning in handling high-dimensional data has paved the way for employing deep neural networks in solving high-dimensional PDEs [1, 8, 16, 17, 19, 31, 35, 45, 47, 52]. Due to the excellent approximation power of deep neural networks, several numerical schemes have been proposed for solving PDEs, including the deep Ritz method [16], PDE-net [34], PINNs [45] and weak adversarial networks [52]. Among the various techniques developed, PINNs have emerged as a particularly powerful approach [45]. PINNs not only leverage the robust approximation abilities of deep learning but also seamlessly incorporate the underlying physical laws of the PDEs, making them highly effective for solving high-dimensional problems [25, 43, 44]. The success of PINNs has spurred deeper theoretical analysis, highlighting the need for comprehensive error analysis in deep learning, including approximation, generalization, and optimization errors.

While several studies have investigated the theoretical mechanisms of PINNs [13, 21–23, 26, 32, 36, 38, 39, 41, 46, 49, 50], these analyses exhibit two key limitations. First, they are typically conducted in scenarios where the number of neural network parameters is smaller than the number of training samples. Second, these analyses often do not address optimization errors, which are crucial for a comprehensive understanding of model performance. Specifically, over-parameterized deep neural networks, where the number of parameters significantly exceeds the sample size, are frequently employed in real-world applications due to their computational efficiency during training. Although extensive research has examined the role of over-parameterization in linear and kernel models, particularly in relation to the double descent phenomenon [2–7, 20, 33, 42], the underlying reasons for the effectiveness of over-parameterized deep neural networks remain unclear. Providing theoretical guarantees in such regimes continues to be a fundamental yet challenging problem. Recent studies have reported convergence results for over-parameterized norm-controlled neural networks in both regression and PDE settings [11, 29, 30, 51]. However, these analyses typically assume that optimization algorithms such as gradient descent (GD) or stochastic gradient descent (SGD) yield the empirical risk minimization (ERM) estimator, the theoretically optimal solution, thus neglecting the influence of optimization errors introduced during the training process.

In this work, we establish a complete error analysis of PINNs for elliptic equations in the over-parameterized setting using projected stochastic gradient descent (PSGD) optimization. Our analysis accounts for all three key error components: approximation error, statistical (generalization) error, and optimization error. By integrating these error components within a unified theoretical framework, we derive explicit convergence rates and precise conditions for achieving specified accuracy levels when solving elliptic boundary value problems. This represents a significant advancement in the theoretical understanding of PINNs.

This paper is structured as follows. Section 2 introduces the problem setup, describes the neural network structure, and outlines the PSGD algorithm for optimization. Section 3 presents a comprehensive end-to-end error analysis for solving second-order elliptic equations with PINNs. Our main result, presented in Theorem 3.5, provides a complete characterization of the convergence behavior under suitable choices of network architecture and optimization parameters.

### 1.1 Notation

We establish notation for this paper as follows. Bold-faced letters denote vectors, while capital letters represent matrices or fixed parameters. The symbol  $C$  denotes an absolute constant, whereas  $C(a,b)$  or  $C_i(a,b)$  represent constants that depend only on parameters  $a$  and  $b$ . For positive functions  $f(x)$  and  $g(x)$ , we write  $f(x) = \mathcal{O}(g(x))$  when  $f(x) \leq Cg(x)$  for some  $C > 0$ , and use  $\tilde{\mathcal{O}}(\cdot)$  to omit logarithmic factors.

Let  $\mathbb{N}$  denote the natural numbers and  $\mathbb{N}^+ := \{x \in \mathbb{N} \mid x > 0\}$  the positive integers. For  $x \in \mathbb{R}$ ,  $\lfloor x \rfloor := \max\{k \in \mathbb{N} : k \leq x\}$  and  $\lceil x \rceil := \min\{k \in \mathbb{N} : k \geq x\}$  denote the floor and ceiling functions. For  $N \in \mathbb{N}^+$ ,  $[N] := \{1, 2, \dots, N\}$  represents the set of integers from 1 to  $N$ . Given a vector  $\mathbf{x} = (x_1, \dots, x_d)^T \in \mathbb{R}^d$ , its  $\ell^2$ -norm and  $\ell^\infty$ -norm are defined as  $\|\mathbf{x}\|_2 := \sqrt{\sum_{i=1}^d x_i^2}$  and  $\|\mathbf{x}\|_\infty := \max_{1 \leq i \leq d} |x_i|$ . For a multi-index  $\boldsymbol{\alpha} \in \mathbb{N}^d$ , we define  $\|\boldsymbol{\alpha}\|_1 := \alpha_1 + \dots + \alpha_d$  and  $\boldsymbol{\alpha}! := \alpha_1! \cdots \alpha_d!$ .

For an open set  $\mathcal{D} \subset \mathbb{R}^d$  and a function  $f : \mathcal{D} \rightarrow \mathbb{R}$ , we define the  $L^p(\mathcal{D})$  norm as:

$$\|f\|_{L^p(\mathcal{D})} := \left( \int |f(\mathbf{x})|^p d\mathbf{x} \right)^{1/p}, \quad p \in (0, \infty); \quad \|f\|_{L^\infty(\mathcal{D})} := \sup_{\mathbf{x} \in \mathcal{D}} |f(\mathbf{x})|.$$

The derivative of order  $\boldsymbol{\alpha}$  of  $f$  is denoted by:

$$D^{\boldsymbol{\alpha}} f := \frac{\partial^{|\boldsymbol{\alpha}|_1} f}{\partial x_1^{\alpha_1} \partial x_2^{\alpha_2} \cdots \partial x_d^{\alpha_d}}.$$

We denote by  $C^s(\mathcal{D})$  the set of  $s$ -times continuously differentiable functions on  $\mathcal{D}$  for  $s \in \mathbb{N} \cup \{\infty\}$ . When  $\bar{\mathcal{D}}$  is compact and  $f \in C^s(\mathcal{D})$ , we define:

$$\|f\|_{C^s(\mathcal{D})} := \max_{0 \leq \|\boldsymbol{\alpha}\|_1 \leq s} \sup_{\mathbf{x} \in \bar{\mathcal{D}}} |D^{\boldsymbol{\alpha}} f(\mathbf{x})|.$$

For  $s \in \mathbb{N}$  and  $1 \leq p < \infty$ , the Sobolev space  $W^{s,p}(\mathcal{D})$  is defined as:

$$W^{s,p}(\mathcal{D}) := \{f \in L^p(\mathcal{D}) : D^{\boldsymbol{\alpha}} f \in L^p(\mathcal{D}), \forall \boldsymbol{\alpha} \in \mathbb{N}^d \text{ with } \|\boldsymbol{\alpha}\|_1 \leq s\}.$$

When  $p = 2$ , we denote  $H^s(\mathcal{D}) := W^{s,2}(\mathcal{D})$  for any  $s \in \mathbb{N}$ . For  $f \in W^{s,p}(\mathcal{D})$  with  $1 \leq p < \infty$ , the Sobolev norm is given by:

$$\|f\|_{W^{s,p}(\mathcal{D})} := \left( \sum_{0 \leq \|\boldsymbol{\alpha}\|_1 \leq s} \|D^{\boldsymbol{\alpha}} f\|_{L^p(\mathcal{D})}^p \right)^{1/p}.$$

For  $p = \infty$ , we define:

$$\|f\|_{W^{s,\infty}(\mathcal{D})} := \max_{0 \leq |\alpha|_1 \leq s} \|D^\alpha f\|_{L^\infty(\mathcal{D})}.$$

For  $f \in L^2(\mathcal{D})$ , we define:

$$\|f\|_{H^{1/2}(\mathcal{D})} := \left( \|f\|_{L^2(\mathcal{D})}^2 + \int_{\mathcal{D}} \int_{\mathcal{D}} \frac{|f(\mathbf{x}) - f(\mathbf{y})|^2}{\|\mathbf{x} - \mathbf{y}\|^{d+1}} d\mathbf{x} d\mathbf{y} \right)^{1/2}.$$

## 2 Preliminaries

In this section, we establish the mathematical foundation for our analysis. We begin with the problem setup by introducing the PINNs formulation for elliptic equations, followed by the neural network architecture, and conclude with the projected stochastic gradient descent algorithm used for optimization.

### 2.1 Problem set-up

We first recall the PINNs proposed in [45]. Within the  $d$ -dimensional unit cube  $[0,1]^d$ , let  $\mathcal{D}$  be a bounded convex domain with boundary  $\partial\mathcal{D}$ . Consider the following linear second-order elliptic equation with Dirichlet boundary conditions:

$$\begin{cases} -\sum_{i,j=1}^d \kappa_{ij} \partial_{s_i} \partial_{s_j} v + \sum_{i=1}^d \nu_i \partial_{s_i} v + \mu v = f & \text{in } \mathcal{D}, \\ v = \psi & \text{on } \partial\mathcal{D}, \end{cases} \quad (2.1)$$

where  $\kappa_{ij} \in C(\bar{\mathcal{D}})$ ,  $\nu_i, \mu \in L^\infty(\mathcal{D})$ ,  $f \in L^2(\mathcal{D})$ ,  $\psi \in L^2(\partial\mathcal{D})$  with the strictly elliptic condition, i.e., there exists a constant  $\tau_0 > 0$  such that  $\sum_{i,j} \kappa_{ij} \xi_i \xi_j \geq \tau_0 |\xi|^2$ ,  $\forall \mathbf{s} \in \mathcal{D}, \xi \in \mathbb{R}^d$ . Assume that (2.1) has a unique strong solution  $v_* \in C^m(\bar{\mathcal{D}})$  with  $m \geq 3$ . Define the coefficient norms  $R_\kappa = \max_{i,j} \{\|\kappa_{ij}\|_{C(\bar{\mathcal{D}})}\}$  and  $R_\nu = \max_i \{\|\nu_i\|_{L^\infty(\mathcal{D})}\}$ . Further, let  $R_\mu = \|\mu\|_{L^\infty(\mathcal{D})}$ ,  $R_f = \|f\|_{L^2(\mathcal{D})}$ , and  $R_\psi = \|\psi\|_{L^2(\partial\mathcal{D})}$ . We set  $R_0 = \max\{R_\kappa, R_\nu, R_\mu, R_f, R_\psi\}$ . The residual of (2.1) is defined as

$$\mathcal{J}(v) := \int_{\mathcal{D}} \left( -\sum_{i,j=1}^d \kappa_{ij} \partial_{s_i} \partial_{s_j} v + \sum_{i=1}^d \nu_i \partial_{s_i} v + \mu v - f \right)^2 ds + \int_{\partial\mathcal{D}} (v - \psi)^2 ds. \quad (2.2)$$

The main idea of PINNs is to use deep learning techniques to minimize the residual (2.2). First, it is rewritten as

$$\begin{aligned} \mathcal{J}(v) = |\mathcal{D}| \mathbb{E}_{S \sim \mathcal{U}(\mathcal{D})} & \left[ -\sum_{i,j=1}^d \kappa_{ij}(S) \partial_{s_i} \partial_{s_j} v(S) + \sum_{i=1}^d \nu_i(S) \partial_{s_i} v(S) \right. \\ & \left. + \mu(S)v(S) - f(S) \right]^2 + |\partial\mathcal{D}| \mathbb{E}_{T \sim \mathcal{U}(\partial\mathcal{D})} [v(T) - \psi(T)]^2, \end{aligned}$$

where  $\mathcal{U}(\mathcal{D})$  and  $\mathcal{U}(\partial\mathcal{D})$  are uniform distribution on  $\mathcal{D}$  and  $\partial\mathcal{D}$ , respectively.

To facilitate numerical computation, a discrete version of  $\mathcal{J}$  is given by

$$\begin{aligned} \widehat{\mathcal{J}}(v) := & \frac{|\mathcal{D}|}{N} \sum_{k=1}^N \left[ - \sum_{i,j=1}^d \kappa_{ij}(S_k) \partial_{s_i} \partial_{s_j} v(S_k) + \sum_{i=1}^d v_i(S_k) \partial_{s_i} v(S_k) \right. \\ & \left. + \mu(S_k) v(S_k) - f(S_k) \right]^2 + \frac{|\partial\mathcal{D}|}{N} \sum_{k=1}^N [v(T_k) - \psi(T_k)]^2, \end{aligned}$$

where  $\{S_k\}_{k=1}^N$  and  $\{T_k\}_{k=1}^N$  are i.i.d. Monte Carlo sample points following  $\mathcal{U}(\mathcal{D})$  on  $\mathcal{D}$  and  $\mathcal{U}(\partial\mathcal{D})$  on  $\partial\mathcal{D}$ , respectively. Then, we select a deep neural network class  $\mathcal{F}_\omega$ , within which we will minimize  $\widehat{\mathcal{J}}(v_\omega)$  for  $v_\omega \in \mathcal{F}_\omega$ .

## 2.2 Neural network class

Let  $D, d \in \mathbb{N}$ ,  $H_0 = d$  and  $H_D = 1$ . Consider a neural network function  $\Phi_\omega : \mathbb{R}^d \rightarrow \mathbb{R}$  with following structure:

$$\begin{aligned} \Phi^{[0]}(\mathbf{s}) &= \mathbf{s}, \\ \Phi^{[\ell]}(\mathbf{s}) &= \sigma(\mathbf{W}_{\ell-1} \Phi^{[\ell-1]}(\mathbf{s}) + \mathbf{c}_{\ell-1}), \quad \ell = 1, \dots, D-1, \\ \Phi_\omega(\mathbf{s}) &= \Phi^{[D]}(\mathbf{s}) = \mathbf{W}_{D-1} \Phi^{[D-1]}(\mathbf{s}) + \mathbf{c}_{D-1}, \end{aligned}$$

where  $\mathbf{W}_\ell = (w_{i,j}^{[\ell]}) \in \mathbb{R}^{H_{\ell+1} \times H_\ell}$ ,  $\mathbf{c}_\ell = (c_i^{[\ell]}) \in \mathbb{R}^{H_{\ell+1}}$  and

$$\boldsymbol{\omega} = (w_{1,1}^{[0]}, \dots, w_{H_D, H_{D-1}}^{[D-1]}, c_1^{[0]}, \dots, c_{H_D}^{[D-1]}) \in \Omega.$$

We say  $\Phi_\omega$  belongs to  $\mathcal{F}_\sigma(H, D, R_\omega)$  if it satisfies:

$$\max\{H_1, \dots, H_D\} = H, \quad \|\boldsymbol{\omega}\|_\infty \leq R_\omega.$$

When solving PDEs with deep learning, the neural network must be differentiable to satisfy the PDE constraints. In this paper, we let  $\sigma = \tanh$ , and abbreviate the network class by  $\mathcal{F}(H, D, R_\omega)$ .

For  $\Phi_\omega \in \mathcal{F}(H, D, R_\omega)$ , let  $n_\ell$  be the number of its nonzero weights in the first  $\ell$  layers. It directly follows that  $n_D \leq \mathfrak{G}(H, D, d)$ , where

$$\mathfrak{G}(H, D, d) := (H+1)[(D-2)H+d+1]. \quad (2.3)$$

Since  $H$  typically exceeds  $d$ , we will also use the bound  $n_D \leq H(H+1)D$ . Note that any weight vector  $\boldsymbol{\omega}$  can be extended to a  $\mathfrak{G}(H, D, d)$ -dimensional vector  $\boldsymbol{\omega}'$  by zero-padding, with  $\Phi_\omega = \Phi_{\boldsymbol{\omega}'} \in \mathcal{F}(H, D, R_\omega)$ . Therefore, the domain  $\Omega$  of neural network parameters can be formalized as:

$$\Omega = [-R_\omega, R_\omega]^{\mathfrak{G}(H, D, d)}.$$

Further, we introduce a parallel neural network class  $\mathcal{PF}_{q,K}(H,D,R_\omega)$ , defined as a linear combination of several sub-network classes  $\mathcal{F}(H,D,R_\omega)$ . Specifically, any function  $v_{q,\omega} \in \mathcal{PF}_{q,K}(H,D,R_\omega)$  can be expressed as

$$v_{q,\omega}(\mathbf{s}) = \sum_{k=1}^q \delta_k \Phi_\omega^k(\mathbf{s}), \quad \delta_k \in \mathbb{R},$$

where  $\Phi_\omega^k(\mathbf{s}) \in \mathcal{F}(H,D,R_\omega)$ ,  $\sum_{k=1}^q |\delta_k| \leq K$ . Let  $\boldsymbol{\omega}_{\text{int}}^q := (\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_q)$  with

$$\boldsymbol{\omega}_k = (w_{k,1,1}^{[0]}, \dots, w_{k,H_D,H_{D-1}}^{[D-1]}, c_{k,1}^{[0]}, \dots, c_{k,H_D}^{[D-1]}) \in \Omega$$

parameterizing  $\Phi_\omega^k$ . Let  $\boldsymbol{\omega}_{\text{ext}}^q := (\delta_1, \dots, \delta_q)$ ,  $\boldsymbol{\omega}_{\text{all}}^q := (\boldsymbol{\omega}_{\text{int}}^q, \boldsymbol{\omega}_{\text{ext}}^q) \in \Omega^q$ . The symbol  $\mathcal{PF}$  serves both as an abbreviation for parallel network class, and we may write  $\mathcal{PF}$  instead of  $\mathcal{PF}_{q,K}(H,D,R_\omega)$  when the context is clear.

As implied above, we will select a  $\mathcal{PF}_{q,K}$  class with  $q$  sub-networks for implementation in this paper. Consequently,  $\widehat{\mathcal{J}}(v_{q,\omega})$  is given by

$$\begin{aligned} \widehat{\mathcal{J}}(v_{q,\omega}) = & \frac{|\mathcal{D}|}{N} \sum_{k=1}^N \left[ - \sum_{i,j=1}^d \kappa_{ij}(S_k) \partial_{s_i s_j}^2 v_{q,\omega}(S_k) + \sum_{i=1}^d v_i(S_k) \partial_{s_i} v_{q,\omega}(S_k) \right. \\ & \left. + \mu(S_k) v_{q,\omega}(S_k) - f(S_k) \right]^2 + \frac{|\partial\mathcal{D}|}{N} \sum_{k=1}^N [v_{q,\omega}(T_k) - \psi(T_k)]^2. \end{aligned} \quad (2.4)$$

We now introduce some auxiliary function classes, which will play a crucial role in the subsequent analysis of optimization and statistical errors. First, let us define the squared residual classes:

$$\begin{aligned} \mathcal{E}'_1 := & \left\{ \pm e: \mathcal{D} \rightarrow \mathbb{R} \mid \exists v_{q,\omega} \in \mathcal{PF}_{q,K}(H,D,R_\omega), s.t. \right. \\ & \left. e(\cdot; \boldsymbol{\omega}) = \left[ - \sum_{i,j=1}^d \kappa_{ij} \partial_{s_i s_j}^2 v_{q,\omega} + \sum_{i=1}^d v_i \partial_{s_i} v_{q,\omega} + \mu v_{q,\omega} - f \right]^2 \right\}, \\ \mathcal{E}'_2 := & \left\{ \pm e: \partial\mathcal{D} \rightarrow \mathbb{R} \mid \exists v_{q,\omega} \in \mathcal{PF} s.t. e(\cdot; \boldsymbol{\omega}) = (v_{q,\omega} - \psi)^2 \right\}. \end{aligned}$$

The corresponding residual classes without the square operation are defined as:

$$\begin{aligned} \mathcal{E}_1 := & \left\{ \pm e: \mathcal{D} \rightarrow \mathbb{R} \mid \exists v_{q,\omega} \in \mathcal{PF}_{q,K}(H,D,R_\omega), s.t. \right. \\ & \left. e(\cdot; \boldsymbol{\omega}) = - \sum_{i,j=1}^d \kappa_{ij} \partial_{s_i s_j}^2 v_{q,\omega} + \sum_{i=1}^d v_i \partial_{s_i} v_{q,\omega} + \mu v_{q,\omega} - f \right\}, \\ \mathcal{E}_2 := & \left\{ \pm e: \partial\mathcal{D} \rightarrow \mathbb{R} \mid \exists v_{q,\omega} \in \mathcal{PF} s.t. e(\cdot; \boldsymbol{\omega}) = v_{q,\omega} - \psi \right\}. \end{aligned}$$

Finally, we define the sub-network function classes:

$$\begin{aligned}\mathcal{E}_1^{sub} &:= \left\{ \pm e: \mathcal{D} \rightarrow \mathbb{R} \mid \exists \Phi_\omega \in \mathcal{F}(H, D, R_\omega) \text{ s.t.} \right. \\ &\quad \left. e(\cdot; \omega) = - \sum_{i,j=1}^d \kappa_{ij} \partial_{s_i, s_j}^2 \Phi_\omega + \sum_{i=1}^d \nu_i \partial_{s_i} \Phi_\omega + \mu \Phi_\omega \right\}, \\ \mathcal{E}_2^{sub} &:= \left\{ \pm e: \partial \mathcal{D} \rightarrow \mathbb{R} \mid \exists \Phi_\omega \in \mathcal{F}(H, D, R_\omega) \text{ s.t. } e(\cdot; \omega) = \Phi_\omega \right\}.\end{aligned}$$

### 2.3 Projected stochastic gradient descent

To minimize the residual (2.4), we apply the PSGD algorithm, which is well-suited for large scale optimization problems with constraints.

We start by initializing the network parameters. The linear coefficients  $\omega_{\text{ext}}^q$  are initialized to zero, i.e.,  $(\omega_{\text{ext}}^q)^{(0)} = (\delta_1, \dots, \delta_q)^{(0)} = \mathbf{0}$ ; Each sub-network parameter in  $(\omega_{\text{int}}^q)^{(0)}$  is initialized independently following a uniform distribution  $\mathcal{U}[-U, U]$ . Overall,

$$(w_{k,i,j}^{[q]})^{(0)} \sim_{\text{i.i.d.}} \mathcal{U}[-U, U], \quad (c_{k,i}^{[q]})^{(0)} \sim_{\text{i.i.d.}} \mathcal{U}[-U, U], \quad (\delta_k)^{(0)} = 0. \quad (2.5)$$

Then, for positive parameters  $\xi$  and  $\rho$ , we define the constraint sets:

- (i)  $X_\xi$ : the (random) collection of weight vectors  $\omega_{\text{int}}^q$  satisfying

$$\|\omega_{\text{int}}^q - (\omega_{\text{int}}^q)^{(0)}\|_2 \leq \xi. \quad (2.6)$$

- (ii)  $Y_\rho$ : all weight vectors  $\omega_{\text{ext}}^q$  such that

$$\|\omega_{\text{ext}}^q\|_1 = \sum_{k=1}^q |\delta_k| \leq \rho. \quad (2.7)$$

Since the Monte Carlo samples  $\{S_k, T_k\}_{k=1}^N$  remain unchanged during optimization,  $\widehat{\mathcal{J}}(v_{q,\omega})$  depends only on the network parameters  $\omega_{\text{all}}^q$ , and we denote it as  $\mathcal{L}(\omega_{\text{all}}^q)$ :

$$\widehat{\mathcal{J}}(v_{q,\omega}) =: \mathcal{L}(\omega_{\text{all}}^q) = \frac{1}{N} \sum_{k=1}^N \mathcal{L}_k(\omega_{\text{all}}^q) = \frac{1}{N} \sum_{k=1}^N \mathcal{L}_k(\omega_{\text{int}}^q, \omega_{\text{ext}}^q).$$

Here,  $\mathcal{L}_k(\omega_{\text{all}}^q)$  represents the contribution from the  $k$ -th Monte Carlo pair  $\{S_k, T_k\}$ .

At each iteration  $t \in [T_{\text{iter}}]$ , a random index  $k_t \in [N]$  is uniformly selected, and the gradient of  $\mathcal{L}_{k_t}$  with respect to  $\omega_{\text{all}}^q$  is computed, followed by a gradient descent step with step size  $\gamma$ . The updated parameters are then projected onto  $X_\xi \times Y_\rho$ .

Formally, the iteration is given by

$$\begin{aligned}\mathbf{y}^{(t)} &= \nabla_{\omega_{\text{all}}^q} \mathcal{L}_{k_t}(\omega_{\text{int}}^{q,(t-1)}, \omega_{\text{ext}}^{q,(t-1)}), \\ (\omega_{\text{int}}^q, \omega_{\text{ext}}^q)^{(t)} &= \text{Proj}_{X_\xi \times Y_\rho} \left\{ (\omega_{\text{int}}^q, \omega_{\text{ext}}^q)^{(t-1)} - \gamma \mathbf{y}^{(t)} \right\}.\end{aligned} \quad (2.8)$$

By [14], the projection onto the  $\ell_1$  ball  $Y_\rho$  can be implemented efficiently with linear time complexity relative to the dimensionality, while the projection onto the  $\ell_2$  ball  $X_\xi$  can be computed in closed form. It is worth emphasizing that the projection step plays a vital role in our error analysis: by constraining the neural network parameters, it enables us to bound the Rademacher complexity of the overparameterized  $\mathcal{PF}$  class and derive size-independent statistical error.

Finally, we define the numerical PDE solution  $v_{\mathcal{A}}$  of the PSGD algorithm as  $v_{q,\omega}$  parameterized with  $(\omega_{\text{int}}^q, \omega_{\text{ext}}^q)^{(t_*)}$ , where the optimal stopping time  $t_*$  is defined as

$$t_* := \arg \min_{t=0, \dots, T} \mathbb{E}_{\text{SGD}} [\mathcal{L}(\omega_{\text{int}}^{q,(t)}, \omega_{\text{ext}}^{q,(t)})].$$

Here,  $\mathbb{E}_{\text{SGD}}$  denotes the conditional expectation of the loss, taken over the random indices  $\{k_t\}_{t=1}^{T_{\text{iter}}}$ , given a fixed training set and network initialization  $(\omega_{\text{int}}^q)^{(0)}$ . Consequently,  $t_*$  remains a random variable depending on these conditions.

### 3 Complete error analysis

Building on the preliminaries above, we conduct a comprehensive error analysis that includes approximation, optimization, and statistical errors. This analysis provides guidelines for choosing the neural network architecture and optimization parameters to achieve an  $\epsilon$ -accuracy between  $v_{\mathcal{A}}$  and  $v_*$ . Section 3.1 introduces the error decomposition, followed by detailed analyses in Sections 3.2 to 3.4, and the main result, Theorem 3.5, in Section 3.5.

#### 3.1 Error decomposition

We first characterize a specific ‘optimization error’ for the PSGD algorithm as:

$$\varepsilon_{\text{opt}} := \mathbb{E}_{\text{SGD}} [\widehat{\mathcal{J}}(v_{\mathcal{A}})] - \widehat{\mathcal{J}}(\bar{v}).$$

The term  $\bar{v}$  denotes the best approximation in some parallel network class  $\mathcal{PF}'$  that may not coincide with the algorithm’s  $\mathcal{PF}$  class, defined by:

$$\bar{v} \in \operatorname{argmin}_{v \in \mathcal{PF}'} \|v - v_*\|_{C^2(\mathcal{D})}^2. \quad (3.1)$$

This leads to our error decomposition theorem:

**Theorem 3.1.** *Consider the PSGD algorithm solution  $v_{\mathcal{A}}$  from Section 2.3 for solving Eq. (2.1) via PINNs, along with  $\bar{v} \in \mathcal{PF}'$  specified in (3.1). The  $H^{1/2}$  error between  $v_{\mathcal{A}}$  and the true*

solution  $v_*$  satisfies:

$$\begin{aligned} & \mathbb{E}_{\text{SGD}} [\|v_{\mathcal{A}} - v_*\|_{H^{1/2}(\mathcal{D})}^2] \\ & \leq C(\mathcal{D}, d, R_0) \left\{ \underbrace{\sup_{v \in \mathcal{PF}} |\mathcal{J}(v) - \widehat{\mathcal{J}}(v)|}_{\varepsilon_{sta}} + \underbrace{[\mathbb{E}_{\text{SGD}} [\widehat{\mathcal{J}}(v_{\mathcal{A}})] - \widehat{\mathcal{J}}(\bar{v})]}_{\varepsilon_{opt}} + \underbrace{\|\bar{v} - v_*\|_{C^2(\bar{\mathcal{D}})}^2}_{\varepsilon_{app}} \right\}. \end{aligned}$$

*Proof.* By Proposition 2.2 in [26], we have

$$\begin{aligned} & \mathbb{E}_{\text{SGD}} [\|v_{\mathcal{A}} - v_*\|_{H^{1/2}(\mathcal{D})}^2] \leq C(\mathcal{D}, d, R_0) \{ \mathbb{E}_{\text{SGD}} [\mathcal{J}(v_{\mathcal{A}})] - \mathcal{J}(v_*) \} \\ & = C(\mathcal{D}, d, R_0) \left\{ \mathbb{E}_{\text{SGD}} [\mathcal{J}(v_{\mathcal{A}}) - \widehat{\mathcal{J}}(v_{\mathcal{A}})] + \mathbb{E}_{\text{SGD}} [\widehat{\mathcal{J}}(v_{\mathcal{A}})] - \widehat{\mathcal{J}}(\bar{v}) \right. \\ & \quad \left. + \widehat{\mathcal{J}}(\bar{v}) - \mathcal{J}(\bar{v}) + \mathcal{J}(\bar{v}) - \mathcal{J}(v_*) \right\} \\ & \leq C(\mathcal{D}, d, R_0) \left\{ 2 \sup_{v \in \mathcal{PF}} |\mathcal{J}(v) - \widehat{\mathcal{J}}(v)| + \mathbb{E}_{\text{SGD}} [\widehat{\mathcal{J}}(v_{\mathcal{A}})] - \widehat{\mathcal{J}}(\bar{v}) + \mathcal{J}(\bar{v}) - \mathcal{J}(v_*) \right\} \\ & \leq C(\mathcal{D}, d, R_0) \left\{ 2 \sup_{v \in \mathcal{PF}} |\mathcal{J}(v) - \widehat{\mathcal{J}}(v)| + [\mathbb{E}_{\text{SGD}} [\widehat{\mathcal{J}}(v_{\mathcal{A}})] - \widehat{\mathcal{J}}(\bar{v})] \right. \\ & \quad \left. + 6|\mathcal{D}|d^2R_0^2\|\bar{v} - v_*\|_{H^2(\mathcal{D})}^2 + |\partial\mathcal{D}|\|\bar{v} - v_*\|_{C(\partial\mathcal{D})}^2 \right\} \\ & \leq C(\mathcal{D}, d, R_0) \left\{ \underbrace{\sup_{v \in \mathcal{PF}} |\mathcal{J}(v) - \widehat{\mathcal{J}}(v)|}_{\varepsilon_{sta}} + \underbrace{[\mathbb{E}_{\text{SGD}} [\widehat{\mathcal{J}}(v_{\mathcal{A}})] - \widehat{\mathcal{J}}(\bar{v})]}_{\varepsilon_{opt}} + \underbrace{\|\bar{v} - v_*\|_{C^2(\bar{\mathcal{D}})}^2}_{\varepsilon_{app}} \right\}. \end{aligned}$$

This completes the proof.  $\square$

### 3.2 Approximation error

The approximation error  $\varepsilon_{app}$  is defined as

$$\varepsilon_{app} := \|\bar{v} - v_*\|_{C^2(\bar{\mathcal{D}})}^2,$$

where  $\bar{v}$  is given in (3.1). The following approximation result in  $C^2(\bar{\mathcal{D}})$  follows directly from a similar argument as in Theorem 3.2 of [28]. The proof is provided in Appendix D.

**Theorem 3.2.** *Given any  $v_* \in C^m(\bar{\mathcal{D}})$  with  $m \geq 3$ , for some small  $\varepsilon^* > 0$  and any  $0 < \varepsilon < \varepsilon^*$ , there exists  $v_{\bar{q}, \bar{\omega}} \in \mathcal{PF}_{\bar{q}, \bar{R}}(\bar{H}, \bar{\mathcal{D}}, R_{\bar{\omega}})$  with*

$$\begin{aligned} \bar{q} &= \lceil C_1 \cdot \varepsilon^{-\frac{d}{m-2l-2}} \rceil, \quad \bar{K} = C_2 \cdot \varepsilon^{-\frac{d}{m-2l-2}}, \quad \bar{H} = 2^{\lceil \log_2(d+m-1) \rceil + 1}, \\ \bar{D} &= \lceil \log_2(d+m-1) \rceil + 2, \quad R_{\bar{\omega}} = C_3 \cdot \varepsilon^{-\frac{2d+2m}{m-2l-2}}, \end{aligned}$$

such that

$$\|v_{\bar{q},\bar{\omega}} - v_*\|_{C^2(\bar{D})} \leq \epsilon,$$

where  $0 < \iota < 1$ ; Constants  $C_1, C_2$ , and  $C_3$  depend exclusively on  $d$  and  $m$ .

### 3.3 Optimization error

In this section, we provide a complete analysis of the optimization error  $\varepsilon_{opt}$ . Let  $\bar{v} = v_{\bar{q},\bar{\omega}} \in \mathcal{PF}_{\bar{q},\bar{K}}(\bar{H},\bar{D},R_{\bar{\omega}})$  in Theorem 3.2, then  $\varepsilon_{opt}$  is defined as

$$\varepsilon_{opt} = \mathbb{E}_{\text{SGD}} [\hat{\mathcal{J}}(v_{\mathcal{A}})] - \hat{\mathcal{J}}(v_{\bar{q},\bar{\omega}}). \quad (3.2)$$

Let  $\tau > 0$ , and let  $G, J \in \mathbb{N}$  with  $J$  sufficiently large. Set the number of sub-networks  $q = \bar{q} \cdot G \cdot J$ . As indicated in (3.2), the weights of  $v_{\bar{q},\bar{\omega}}$  are treated as 'target parameters'. Thus, we set the sub-network width  $H$  in our implemented  $v_{q,\omega} \in \mathcal{PF}$  to  $\bar{H}$ , the sub-network depth  $D$  to  $\bar{D}$ , and the uniform distribution range  $U$  to  $R_{\bar{\omega}}$ . For random initialization  $(\omega_{\text{int}}^q)^{(0)} = (\omega_1^{(0)}, \dots, \omega_q^{(0)})$ , we aim to define an event which contains all the 'sufficiently good' initialization with respect to the target  $(\bar{\omega}_1, \dots, \bar{\omega}_{\bar{q}})$ .

Let us define the event  $E_{q,\bar{q},G,\tau}$  as follows: For each target parameter  $\bar{\omega}_k$ , there exist at least  $G$  sub-networks among the total  $q$  sub-networks, such that their initial parameters  $(\omega_{\cdot})^{(0)}$  are within a  $\tau$ -neighborhood of  $\bar{\omega}_k$  in the infinity norm, i.e.,  $\|(\omega_{\cdot})^{(0)} - \bar{\omega}_k\|_{\infty} \leq \tau$ . In other words, event  $E_{q,\bar{q},G,\tau}$  ensures after initialization, each target parameter  $\bar{\omega}_k$  is sufficiently approximated by a minimum of  $G$  sub-networks. We now derive the probability of  $E_{q,\bar{q},G,\tau}$  with the help of its complement  $E_{q,\bar{q},G,\tau}^c$ .

(i) Consider the case where  $\bar{q} = G = 1$  and  $q = J$ . The event  $E_{q,1,1,\tau}^c$  occurs when all initial weight vectors  $(\omega_i)^{(0)}$ , for  $i = 1, \dots, J$ , deviate from  $\bar{\omega}_1$  by more than  $\tau$  in the infinity norm, i.e.,  $\|(\omega_i)^{(0)} - \bar{\omega}_1\|_{\infty} > \tau$ . Given that  $\bar{\omega}_1$  is a fixed vector in  $[-R_{\bar{\omega}}, R_{\bar{\omega}}]^{\mathfrak{G}(\bar{H},\bar{D},d)}$ , where  $\mathfrak{G}(H,D,d)$  is defined in (2.3), and the sub-network parameters are initialized independently according to  $\mathcal{U}[-R_{\bar{\omega}}, R_{\bar{\omega}}]$ , we have

$$\mathbb{P} \left[ \|(\omega_i)^{(0)} - \bar{\omega}_1\|_{\infty} \leq \tau \right] \geq \left( \frac{\tau}{2R_{\bar{\omega}}} \right)^{\mathfrak{G}(H,D,d)} \geq \left( \frac{\tau}{2R_{\bar{\omega}}} \right)^{H(\bar{H}+1)\bar{D}}$$

for any  $i \in \{1, \dots, J\}$ . Thus, we have

$$\mathbb{P}(E_{q,1,1,\tau}^c) = \mathbb{P} \left[ \forall i \in \{1, \dots, J\} : \|(\omega_i)^{(0)} - \bar{\omega}_1\|_{\infty} > \tau \right] \leq \left[ 1 - \tau^{H(\bar{H}+1)\bar{D}} (2R_{\bar{\omega}})^{-H(\bar{H}+1)\bar{D}} \right]^J.$$

(ii) For general  $\bar{q}, G \in \mathbb{N}$ ,  $E_{q,\bar{q},G,\tau}$  occurs when each target parameter  $\bar{\omega}_k$  in the sequence  $(\bar{\omega}_1, \dots, \bar{\omega}_{\bar{q}})$  is approximated by at least  $G$  initial weight vectors from the set  $\{(\omega_i)^{(0)}\}_{i=1}^{\bar{q} \cdot G \cdot J}$  within a  $\tau$ -neighborhood in the infinity norm. We observe that

$$E_{q,\bar{q},G,\tau} \supseteq \bigcap_{j=1}^G \bigcap_{k=1}^{\bar{q}} \left\{ \exists i \in \{[(j-1)\bar{q} + k - 1]J, \dots, [(j-1)\bar{q} + k]J\} : \|(\omega_i)^{(0)} - \bar{\omega}_k\|_{\infty} \leq \tau \right\}.$$

Thus, it holds that

$$E_{q,\bar{q},G,\tau}^c \subseteq \bigcup_{j=1}^G \bigcup_{k=1}^{\bar{q}} \left\{ \forall i \in \{[(j-1)\bar{q}+k-1]J, \dots, [(j-1)\bar{q}+k]J\} : \|(\omega_i)^{(0)} - \bar{\omega}_k\|_\infty > \tau \right\}.$$

This implies

$$\mathbb{P}(E_{q,\bar{q},G,\tau}^c) \leq \bar{q}G \mathbb{P}(E_{q,1,1,\tau}^c) \leq \bar{q}G \left[ 1 - \tau^{\bar{H}(\bar{H}+1)\bar{D}} (2R_{\bar{\omega}})^{-\bar{H}(\bar{H}+1)\bar{D}} \right]^J.$$

Therefore, we have

$$\mathbb{P}(E_{q,\bar{q},G,\tau}) = 1 - \mathbb{P}(E_{q,\bar{q},G,\tau}^c) \geq 1 - \bar{q}G \left[ 1 - \tau^{\bar{H}(\bar{H}+1)\bar{D}} (2R_{\bar{\omega}})^{-\bar{H}(\bar{H}+1)\bar{D}} \right]^J.$$

Based on  $E_{q,\bar{q},G,\tau}$ , we introduce a sequence of random indices  $r_{k,u}(\omega)$  which serve to identify the ‘well-initialized’ parameters: If  $\omega \in E_{q,\bar{q},G,\tau}$ , let  $r_{k,u}$  be the index of the  $u$ -th sub-network among  $q$  satisfying  $\|(\omega_{\cdot})^{(0)} - \bar{\omega}_k\|_\infty \leq \tau$ . We also need  $r_{k,u} \neq r_{k',u'}$  when  $k \neq k'$  or  $u \neq u'$ ; If  $\omega \notin E_{q,\bar{q},G,\tau}$ , we set  $r_{k,u} = (k-1)G + u$ .

A key subsequent step is to define a set of ‘transition parameters’ that leverage  $r_{k,u}$  and  $(\bar{\delta}_1, \dots, \bar{\delta}_{\bar{q}})$  to establish a connection between  $v_{\mathcal{A}}$  and the target  $v_{\bar{q},\bar{\omega}}$ . Specifically, we define  $\omega_{\text{all}}^{q,*}$  as follows:

$$\omega_{\text{all}}^{q,*} := (\omega_{\text{int}}^{q,*}, \omega_{\text{ext}}^{q,*}), \quad \text{where } \omega_{\text{int}}^{q,*} := (\omega_{\text{int}}^q)^{(0)}.$$

For  $\omega_{\text{ext}}^{q,*} := (\delta_1^*, \dots, \delta_{\bar{q}}^*)$ , we set

$$\delta_i^* := \begin{cases} \bar{\delta}_k / G & \text{if } i = r_{k,u} \text{ for } k = 1, \dots, \bar{q}, u = 1, \dots, G, \\ 0 & \text{if } i \notin \{r_{k,u} : k = 1, \dots, \bar{q}, u = 1, \dots, G\}. \end{cases}$$

Denoting  $v_{q,\omega}$  parameterized by  $\omega_{\text{all}}^{q,*}$  as  $v_q^*$ , we express it as:

$$v_q^*(\mathbf{s}) = \sum_{i=1}^q \delta_i^* \cdot (\Phi_{\omega}^i)^{(0)}(\mathbf{s}) = \sum_{k,u} \delta_{r_{k,u}}^* \cdot (\Phi_{\omega}^{r_{k,u}})^{(0)}(\mathbf{s}) = \sum_{k,u} \frac{\bar{\delta}_k}{G} \cdot (\Phi_{\omega}^{r_{k,u}})^{(0)}(\mathbf{s}), \quad (3.3)$$

with  $(\Phi_{\omega}^i)^{(0)}$  indicating the  $i$ -th sub-network  $\Phi_{\omega}^i$  parameterized by  $(\omega_i)^{(0)}$ . To clarify the connection between  $v_q^*$  and  $v_{\bar{q},\bar{\omega}}$ , we expand the latter as:

$$v_{\bar{q},\bar{\omega}}(\mathbf{s}) = \sum_{k=1}^{\bar{q}} \bar{\delta}_k \cdot \Phi_{\bar{\omega}}^k(\mathbf{s}) = \sum_{k=1}^{\bar{q}} \sum_{u=1}^G \frac{\bar{\delta}_k}{G} \cdot \Phi_{\bar{\omega}}^k(\mathbf{s}), \quad (3.4)$$

where  $\Phi_{\bar{\omega}}^k$  represents the  $k$ -th sub-network in  $v_{\bar{q},\bar{\omega}}$  with weights  $\bar{\omega}_k$ . When the weights  $\bar{\omega}_k$  and  $(\omega_{r_{k,u}})^{(0)}$  are close,  $v_q^*$  approximates  $v_{\bar{q},\bar{\omega}}$  well — a property guaranteed by event  $E_{q,\bar{q},G,\tau}$ . The optimization error is then decomposed into two parts through  $\hat{\mathcal{J}}(v_q^*)$ :

$$\varepsilon_{opt} = \underbrace{\hat{\mathcal{J}}(v_q^*) - \hat{\mathcal{J}}(v_{\bar{q},\bar{\omega}})}_{\text{initialization error}} + \underbrace{\mathbb{E}_{\text{SGD}}[\hat{\mathcal{J}}(v_{\mathcal{A}})] - \hat{\mathcal{J}}(v_q^*)}_{\text{iteration error}}. \quad (3.5)$$

As we will see, the two error components exhibit both independence and coupling. The independence is reflected in the fact that regardless of the initial sub-network weights  $(\omega_{\text{int}}^q)^{(0)}$  (and consequently, the initialization error), the iteration error can be rendered arbitrarily small through convex optimization techniques, provided that the iteration steps  $T_{\text{iter}}$  and over-parameterization index  $G$  are sufficiently large. Their coupling manifests in that controlling both errors at the same prescribed order with high probability necessitates specific conditions on the over-parameterization index  $J$ , which in turn constrains the selection of iteration steps  $T_{\text{iter}}$ .

**Initialization error:** As shown in (3.3) and (3.4), we have

$$v_q^*(s) = \sum_{k,u} \frac{\bar{\delta}_k}{G} (\Phi_{\omega}^{r_{k,u}})^{(0)}(s), \quad v_{\bar{q},\bar{\omega}}(s) = \sum_{k=1}^{\bar{q}} \bar{\delta}_k \Phi_{\bar{\omega}}^k(s) = \sum_{k,u} \frac{\bar{\delta}_k}{G} \Phi_{\bar{\omega}}^k(s).$$

Denote by  $e_i(\cdot; \omega)$  a function in  $\mathcal{E}_i^{\text{sub}}$ ,  $i = 1, 2$ . Using Cauchy-Schwartz inequality and Lemma A.4, it holds that

$$\begin{aligned} & \left| \widehat{\mathcal{J}}(v_q^*) - \widehat{\mathcal{J}}(v_{\bar{q},\bar{\omega}}) \right| \\ & \leq \frac{|\mathcal{D}|}{N} \sum_{k=1}^N \left| \left[ \sum_{k,u} \frac{\bar{\delta}_k}{G} \cdot e_1(S_k; \omega_{r_{k,u}}^{(0)}) - f(S_k) \right]^2 - \left[ \sum_{k,u} \frac{\bar{\delta}_k}{G} \cdot e_1(S_k; \bar{\omega}_k) - f(S_k) \right]^2 \right| \\ & \quad + \frac{|\partial \mathcal{D}|}{N} \sum_{k=1}^N \left| \left[ \sum_{k,u} \frac{\bar{\delta}_k}{G} \cdot e_2(T_k; \omega_{r_{k,u}}^{(0)}) - \psi(T_k) \right]^2 - \left[ \sum_{k,u} \frac{\bar{\delta}_k}{G} \cdot e_2(T_k; \bar{\omega}_k) - \psi(T_k) \right]^2 \right| \\ & \leq \frac{2|\mathcal{D}|}{N} \sum_{k=1}^N (B_1 \bar{K} + R_0) \cdot \left[ \sum_{k,u} \left| \frac{\bar{\delta}_k}{G} \right| \cdot \left| e_1(S_k; \omega_{r_{k,u}}^{(0)}) - e_1(S_k; \bar{\omega}_k) \right| \right] \\ & \quad + \frac{2|\partial \mathcal{D}|}{N} \sum_{k=1}^N (B_2 \bar{K} + R_0) \cdot \left[ \sum_{k,u} \left| \frac{\bar{\delta}_k}{G} \right| \cdot \left| e_2(T_k; \omega_{r_{k,u}}^{(0)}) - e_2(T_k; \bar{\omega}_k) \right| \right] \\ & \leq 2(|\mathcal{D}| + |\partial \mathcal{D}|)(L_1 + L_2)(B_1 \bar{K} + B_2 \bar{K} + 2R_0) \bar{K} \cdot \max_{k,u} \|(\omega_{r_{k,u}})^{(0)} - \bar{\omega}_k\|_2. \end{aligned}$$

Since

$$\|(\omega_{r_{k,u}})^{(0)} - \bar{\omega}_k\|_2 \leq [\bar{H}(\bar{H} + 1)\bar{D}]^{1/2} \cdot \max_{k,u} \|(\omega_{r_{k,u}})^{(0)} - \bar{\omega}_k\|_{\infty},$$

we get

$$\widehat{\mathcal{J}}(v_q^*) - \widehat{\mathcal{J}}(v_{\bar{q},\bar{\omega}}) \leq C_4 \cdot \bar{K}^2 \cdot R_{\bar{\omega}}^{5\bar{D}} \cdot \max_{k,u} \|(\omega_{r_{k,u}})^{(0)} - \bar{\omega}_k\|_{\infty}.$$

Therefore, with probability at least

$$1 - \bar{q}G \left[ 1 - \tau^{\bar{H}(\bar{H}+1)\bar{D}} (2R_{\bar{\omega}})^{-\bar{H}(\bar{H}+1)\bar{D}} \right]^J,$$

the initialization error in (3.5) is bounded by

$$\widehat{\mathcal{J}}(v_q^*) - \widehat{\mathcal{J}}(v_{\bar{q}, \bar{\omega}}) \leq C_4 \cdot \bar{K}^2 \cdot R_{\bar{\omega}}^{5\bar{D}} \cdot \tau.$$

**Iteration error:** According to Section 2.3, we have  $(\boldsymbol{\omega}_{\text{int}}^q, \boldsymbol{\omega}_{\text{ext}}^q)^{(t)} \in X_{\bar{\zeta}} \times Y_{\rho}$ , which means  $\|(\boldsymbol{\omega}_{\text{ext}}^q)^{(t)}\|_1 \leq \rho$  and  $\|(\boldsymbol{\omega}_{\text{int}}^q)^{(t)} - (\boldsymbol{\omega}_{\text{int}}^q)^{(0)}\|_2 \leq \bar{\zeta}$ . Also, by

$$\begin{aligned} \|(\boldsymbol{\omega}_{\text{int}}^q)^{(t)}\|_{\infty} &\leq \|(\boldsymbol{\omega}_{\text{int}}^q)^{(0)}\|_{\infty} + \|(\boldsymbol{\omega}_{\text{int}}^q)^{(t)} - (\boldsymbol{\omega}_{\text{int}}^q)^{(0)}\|_{\infty} \\ &\leq \|(\boldsymbol{\omega}_{\text{int}}^q)^{(0)}\|_{\infty} + \|(\boldsymbol{\omega}_{\text{int}}^q)^{(t)} - (\boldsymbol{\omega}_{\text{int}}^q)^{(0)}\|_2, \end{aligned}$$

we have  $\|(\boldsymbol{\omega}_{\text{int}}^q)^{(t)}\|_{\infty} \leq R_{\bar{\omega}} + \bar{\zeta}$ . Thus,

$$v_{\mathcal{A}} \in \mathcal{PF}_{q, \rho}(\bar{H}, \bar{D}, R_{\bar{\omega}} + \bar{\zeta}).$$

The following lemma, whose proof is referred to Appendix C, is useful for bounding the iteration error.

**Lemma 3.1.** *Let  $d_1, d_2 \in \mathbb{N}$ , and  $U, V \geq 0$ . Let  $X \subset \mathbb{R}^{d_1}$  and  $Y \subseteq \mathbb{R}^{d_2}$  be closed and convex sets. Consider the function*

$$F(\mathbf{x}, \mathbf{y}) = \frac{1}{N} \sum_{i=1}^N F_i(\mathbf{x}, \mathbf{y}): \quad \mathbb{R}^{d_1} \times \mathbb{R}^{d_2} \rightarrow \mathbb{R}_+,$$

where for each  $i \in [N]$ ,  $F_i(\mathbf{x}, \mathbf{y})$  is differentiable and  $\mathbf{y} \mapsto F_i(\mathbf{x}, \mathbf{y})$  is convex for all  $\mathbf{x} \in \mathbb{R}^{d_1}$ . Assume that for any  $i \in [N]$ ,

$$\|\nabla_{\mathbf{y}} F_i(\mathbf{x}, \mathbf{y})\|_2 \leq V, \quad \forall (\mathbf{x}, \mathbf{y}) \in X \times Y. \quad (3.6)$$

Starting from  $(\mathbf{x}_0, \mathbf{y}_0) \in X \times Y$ , consider the iteration

$$(\mathbf{x}_t, \mathbf{y}_t) = \text{Proj}_{X \times Y} \{ (\mathbf{x}_{t-1}, \mathbf{y}_{t-1}) - \lambda \nabla F_{i_t}(\mathbf{x}_{t-1}, \mathbf{y}_{t-1}) \},$$

where for each time step  $t \in [T]$ , the index  $i_t$  is randomly sampled from  $[N]$  with uniform probability, and  $\lambda = T^{-1}$ . Let  $\mathbf{y}^* \in Y$  satisfy

$$|F(\mathbf{x}_t, \mathbf{y}^*) - F(\mathbf{x}_0, \mathbf{y}^*)| \leq U \|\mathbf{y}^*\|_2 \|\mathbf{x}_t - \mathbf{x}_0\|_2 \quad (3.7)$$

for all  $t \in [T]$ . Then, it holds that

$$\min_{t=0, \dots, T} \mathbb{E}[F(\mathbf{x}_t, \mathbf{y}_t)] - F(\mathbf{x}_0, \mathbf{y}^*) \leq U \|\mathbf{y}^*\|_2 \text{diam}(X) + \frac{\|\mathbf{y}^* - \mathbf{y}_0\|_2^2}{2} + \frac{V^2}{2T}, \quad (3.8)$$

where the expectation is taken with respect to the random indices  $i_1, \dots, i_T$ .

Note that  $\|\omega_{\text{ext}}^{q,*}\|_1 \leq \bar{K}$ . If we set  $\rho = \bar{K}$ , the 'transition parameter'  $\omega_{\text{all}}^{q,*} = (\omega_{\text{int}}^{q,(0)}, \omega_{\text{ext}}^{q,*})$  will fall in  $X_{\bar{\zeta}} \times Y_{\rho}$ . Letting  $(x_t, y_t)$  to be  $(\omega_{\text{int}}^q, \omega_{\text{ext}}^q)^{(t)}$ ,  $y^*$  to be  $\omega_{\text{ext}}^{q,*}$  and applying Lemma 3.1, we could directly obtain an estimate for the iteration error in (3.5). However, such application is valid only when conditions (3.6) and (3.7) are satisfied.

We first obtain an upper bound of  $\nabla_{\omega_{\text{ext}}^q} \mathcal{L}_k(\omega_{\text{int}}^q, \omega_{\text{ext}}^q)$  to meet (3.6). Denote by  $e_i(\cdot; \omega)$  a function in  $\mathcal{E}_i^{\text{sub}}$ ,  $i=1,2$ . Then, for any  $k \in [N]$ , we have

$$\begin{aligned}
& \|\nabla_{\omega_{\text{ext}}^q} \mathcal{L}_k(\omega_{\text{int}}^q, \omega_{\text{ext}}^q)\|_2^2 \\
&= 4 \sum_{i=1}^q \left\{ |\mathcal{D}| \cdot e_1(S_k; \omega_i) \cdot \left[ \sum_{j=1}^q \delta_j e_1(S_k; \omega_j) - f(S_k) \right] \right. \\
&\quad \left. + |\partial \mathcal{D}| \cdot e_2(T_k; \omega_i) \cdot \left[ \sum_{j=1}^q \delta_j e_2(T_k; \omega_j) - \psi(T_k) \right] \right\}^2 \\
&\leq 8 \sum_{i=1}^q \left\{ |\mathcal{D}| \cdot e_1(S_k; \omega_i) \cdot \left[ \sum_{j=1}^q \delta_j e_1(S_k; \omega_j) - f(S_k) \right] \right\}^2 \\
&\quad + 8 \sum_{i=1}^q \left\{ |\partial \mathcal{D}| \cdot e_2(T_k; \omega_i) \cdot \left[ \sum_{j=1}^q \delta_j e_2(T_k; \omega_j) - \psi(T_k) \right] \right\}^2 \\
&\leq 16q |\mathcal{D}|^2 \cdot \sup_{e_1 \in \mathcal{E}_1^{\text{sub}}} \|e_1\|_{\infty}^2 \cdot \left( \|\omega_{\text{ext}}^q\|_1^2 \cdot \sup_{e_1 \in \mathcal{E}_1^{\text{sub}}} \|e_1\|_{\infty}^2 + R_0^2 \right) \\
&\quad + 16q |\partial \mathcal{D}|^2 \cdot \sup_{e_2 \in \mathcal{E}_2^{\text{sub}}} \|e_2\|_{\infty}^2 \cdot \left( \|\omega_{\text{ext}}^q\|_1^2 \cdot \sup_{e_2 \in \mathcal{E}_2^{\text{sub}}} \|e_2\|_{\infty}^2 + R_0^2 \right) \\
&\leq 16 \cdot q \cdot \bar{K}^2 (|\mathcal{D}|^2 + |\partial \mathcal{D}|^2) [B_1^2 (\bar{K}^2 B_1^2 + R_0^2) + B_2^2 (\bar{K}^2 B_2^2 + R_0^2)] \\
&= C_5 \cdot q \cdot \bar{K}^2 \cdot (R_{\bar{\omega}} + \bar{\zeta})^{8D},
\end{aligned}$$

where we have used the properties outlined in Lemma A.4 with  $R_{\omega} = R_{\bar{\omega}} + \bar{\zeta}$ .

Then, we assure that  $\mathcal{L}(\omega_{\text{int}}^q, \omega_{\text{ext}}^q)$  satisfies the condition in (3.7):

$$\begin{aligned}
& \left| \mathcal{L}(\omega_{\text{int}}^{q,(t)}, \omega_{\text{ext}}^{q,*}) - \mathcal{L}(\omega_{\text{int}}^{q,(0)}, \omega_{\text{ext}}^{q,*}) \right| \\
&\leq \frac{|\mathcal{D}|}{N} \sum_{k=1}^N \left| \left[ \sum_{i=1}^q \delta_i^* e_1(S_k; \omega_i^{(t)}) - f(S_k) \right]^2 - \left[ \sum_{i=1}^q \delta_i^* e_1(S_k; \omega_i^{(0)}) - f(S_k) \right]^2 \right| \\
&\quad + \frac{|\partial \mathcal{D}|}{N} \sum_{k=1}^N \left| \left[ \sum_{i=1}^q \delta_i^* e_1(T_k; \omega_i^{(t)}) - \psi(T_k) \right]^2 - \left[ \sum_{i=1}^q \delta_i^* e_1(T_k; \omega_i^{(0)}) - \psi(T_k) \right]^2 \right|.
\end{aligned}$$

By Cauchy-Schwarz inequality and Lemma A.4, we further deduce

$$\begin{aligned}
& \left| \mathcal{L}(\boldsymbol{\omega}_{\text{int}}^{q,(t)}, \boldsymbol{\omega}_{\text{ext}}^{q,*}) - \mathcal{L}(\boldsymbol{\omega}_{\text{int}}^{q,(0)}, \boldsymbol{\omega}_{\text{ext}}^{q,*}) \right| \\
& \leq \frac{2|\mathcal{D}|}{N} \sum_{k=1}^N (B_1 \bar{K} + R_0) \cdot \left\{ \sum_{i=1}^q |\delta_i^*|^2 \cdot \sum_{i=1}^q [e_1(S_k; \boldsymbol{\omega}_i^{(t)}) - e_1(S_k; \boldsymbol{\omega}_i^{(0)})]^2 \right\}^{1/2} \\
& \quad + \frac{2|\partial \mathcal{D}|}{N} \sum_{k=1}^N (B_2 \bar{K} + R_0) \cdot \left\{ \sum_{i=1}^q |\delta_i^*|^2 \cdot \sum_{i=1}^q [e_2(S_k; \boldsymbol{\omega}_i^{(t)}) - e_2(S_k; \boldsymbol{\omega}_i^{(0)})]^2 \right\}^{1/2} \\
& \leq 2(|\mathcal{D}| + |\partial \mathcal{D}|)(L_1 + L_2)[\bar{K}(B_1 + B_2) + 2R_0] \cdot \|\boldsymbol{\omega}_{\text{ext}}^{q,*}\|_2 \cdot \left( \sum_{i=1}^q \|\boldsymbol{\omega}_i^{(t)} - \boldsymbol{\omega}_i^{(0)}\|_2^2 \right)^{1/2} \\
& = C_6 \cdot \bar{K} \cdot (R_{\bar{\omega}} + \xi)^{5\bar{D}} \cdot \|\boldsymbol{\omega}_{\text{ext}}^{q,*}\|_2 \cdot \|\boldsymbol{\omega}_{\text{int}}^{q,(t)} - \boldsymbol{\omega}_{\text{int}}^{q,(0)}\|_2.
\end{aligned}$$

Since

$$\|\boldsymbol{\omega}_{\text{ext}}^{q,*}\|_2 = \sqrt{\sum_{i=1}^q |\delta_i^*|^2} = \frac{1}{\sqrt{G}} \sqrt{\sum_{k=1}^{\bar{q}} |\bar{\delta}_k|^2} \leq \frac{1}{\sqrt{G}} \sum_{k=1}^{\bar{q}} |\bar{\delta}_k| \leq \frac{\bar{K}}{\sqrt{G}},$$

according to (3.8), the iteration error in (3.5) is bounded by

$$\begin{aligned}
& \mathbb{E}_{\text{SGD}}[\hat{\mathcal{J}}(v_{\mathcal{A}})] - \hat{\mathcal{J}}(v_q^*) \\
& \leq C_6 \cdot \bar{K} \cdot \xi \cdot (R_{\bar{\omega}} + \xi)^{5\bar{D}} \cdot \|\boldsymbol{\omega}_{\text{ext}}^{q,*}\|_2 + \frac{1}{2} \|\boldsymbol{\omega}_{\text{ext}}^{q,*}\|_2^2 + \frac{C_5 \cdot q \cdot \bar{K}^2 \cdot (R_{\bar{\omega}} + \xi)^{8\bar{D}}}{2T_{\text{iter}}} \\
& \leq \frac{C_6 \cdot \bar{K}^2 \cdot (R_{\bar{\omega}} + \xi)^{5\bar{D}} \cdot \xi}{\sqrt{G}} + \frac{\bar{K}^2}{2G} + \frac{C_5 \cdot q \cdot \bar{K}^2 \cdot (R_{\bar{\omega}} + \xi)^{8\bar{D}}}{2T_{\text{iter}}} \\
& \leq \frac{C_7 \cdot \bar{K}^2 \cdot (R_{\bar{\omega}} + \xi)^{5\bar{D}} \cdot \xi}{\sqrt{G}} + \frac{C_5 \cdot q \cdot \bar{K}^2 \cdot (R_{\bar{\omega}} + \xi)^{8\bar{D}}}{2T_{\text{iter}}}.
\end{aligned}$$

Finally, we achieve the following result.

**Theorem 3.3.** *Let  $\tau > 0$ , and let  $G, J \in \mathbb{N}$  with  $J$  sufficiently large. Consider the solution  $v_{\mathcal{A}} \in \mathcal{PF}_{q, \bar{K}}(\bar{H}, \bar{D}, R_{\bar{\omega}} + \xi)$  obtained by the PSGD algorithm after  $T_{\text{iter}}$  steps, under the  $\ell_1$  linear coefficient constraint  $\rho = \bar{K}$  and step size  $\gamma = T_{\text{iter}}^{-1}$ . If the number of sub-networks is set to  $q = \bar{q} \cdot G \cdot J$ , then with probability at least*

$$1 - \bar{q}G \left[ 1 - \tau^{\bar{H}(\bar{H}+1)\bar{D}} (2R_{\bar{\omega}})^{-\bar{H}(\bar{H}+1)\bar{D}} \right]^J,$$

the optimization error  $\varepsilon_{\text{opt}}$  in (3.2) is upper bounded by

$$\varepsilon_{\text{opt}} \leq C_4 \cdot \bar{K}^2 \cdot R_{\bar{\omega}}^{5\bar{D}} \cdot \tau + \frac{C_7 \cdot \bar{K}^2 \cdot (R_{\bar{\omega}} + \xi)^{5\bar{D}} \cdot \xi}{\sqrt{G}} + \frac{C_5 \cdot q \cdot \bar{K}^2 \cdot (R_{\bar{\omega}} + \xi)^{8\bar{D}}}{2T_{\text{iter}}},$$

where  $\xi$  is the projection radius of sub-network weights, while  $C_4$ ,  $C_5$ , and  $C_7$  are universal constants depending only on  $\mathcal{D}$ ,  $\bar{H}$ ,  $\bar{D}$ ,  $d$ , and  $R_0$ .

### 3.4 Statistical error

We define the statistical error as

$$\varepsilon_{sta} := \sup_{v \in \mathcal{P}\mathcal{F}} |\mathcal{J}(v) - \widehat{\mathcal{J}}(v)|.$$

Our goal is to bound  $\varepsilon_{sta}$  with high probability, for which we first analyze its expectation  $\mathbb{E}[\varepsilon_{sta}]$ . The following lemma is direct.

**Lemma 3.2.** *The expected statistical error decomposes as:*

$$\begin{aligned} \mathbb{E}[\varepsilon_{sta}] &= \mathbb{E}_{\{S_k, T_k\}_{k=1}^N} \left[ \sup_{v_{q,\omega} \in \mathcal{P}\mathcal{F}} |\mathcal{J}(v_{q,\omega}) - \widehat{\mathcal{J}}(v_{q,\omega})| \right] \\ &\leq \sum_{i=1}^2 \mathbb{E}_{\{S_k, T_k\}_{k=1}^N} \sup_{v_{q,\omega} \in \mathcal{P}\mathcal{F}} |\mathcal{J}_i(v_{q,\omega}) - \widehat{\mathcal{J}}_i(v_{q,\omega})| = \sum_{i=1}^2 \mathbb{E}[\varepsilon_{sta}^i], \end{aligned}$$

where  $\mathcal{J}_1$  quantifies the PDE residual over  $\mathcal{D}$ :

$$\begin{aligned} \mathcal{J}_1(v_{q,\omega}) &= |\mathcal{D}| \mathbb{E}_{S \sim \mathcal{U}(\mathcal{D})} \left[ - \sum_{i,j=1}^d \kappa_{ij}(S) \partial_{s_i} \partial_{s_j} v_{q,\omega}(S) + \sum_{i=1}^d v_i(S) \partial_{s_i} v_{q,\omega}(S) \right. \\ &\quad \left. + \mu(S) v_{q,\omega}(S) - f(S) \right]^2, \end{aligned}$$

$\mathcal{J}_2$  measures the boundary condition mismatch:

$$\mathcal{J}_2(v_{q,\omega}) = |\partial\mathcal{D}| \mathbb{E}_{T \sim \mathcal{U}(\partial\mathcal{D})} [v_{q,\omega}(T) - \psi(T)]^2,$$

and their discrete counterparts  $\widehat{\mathcal{J}}_i$  are denote by:

$$\begin{aligned} \widehat{\mathcal{J}}_1(v_{q,\omega}) &= \frac{|\mathcal{D}|}{N} \sum_{k=1}^N \left[ - \sum_{i,j=1}^d \kappa_{ij}(S_k) \partial_{s_i} \partial_{s_j} v_{q,\omega}(S_k) + \sum_{i=1}^d v_i(S_k) \partial_{s_i} v_{q,\omega}(S_k) \right. \\ &\quad \left. + \mu(S_k) v_{q,\omega}(S_k) - f(S_k) \right]^2, \\ \widehat{\mathcal{J}}_2(v_{q,\omega}) &= \frac{|\partial\mathcal{D}|}{N} \sum_{k=1}^N [v_{q,\omega}(T_k) - \psi(T_k)]^2. \end{aligned}$$

To control each  $\mathbb{E}[\varepsilon_{sta}^i]$ , we use Rademacher complexity as our main analytical tool.

**Definition 3.1.** For a function class  $\mathcal{F}$  and a random sample  $\{Z_k\}_{k=1}^N$ , we characterize two variants of Rademacher complexity:

$$\mathcal{R}_N(\mathcal{F}) = \mathbb{E}_{\{Z_k, \eta_k\}_{k=1}^N} \left[ \sup_{f \in \mathcal{F}} \frac{1}{N} \sum_{k=1}^N \eta_k f(Z_k) \right], \quad \bar{\mathcal{R}}_N(\mathcal{F}) = \mathbb{E}_{\{Z_k, \eta_k\}_{k=1}^N} \left[ \sup_{f \in \mathcal{F}} \frac{1}{N} \left| \sum_{k=1}^N \eta_k f(Z_k) \right| \right],$$

where  $\{\eta_k\}_{k=1}^N$  are i.i.d. Rademacher random variables.

Our analysis relies on standard tools from statistical learning theory, such as Talagrand's contraction inequality (see, e.g., [40]). The following lemma connects  $\mathbb{E}[\varepsilon_{sta}^i]$  to the Rademacher complexity of specific function classes derived from our model  $\mathcal{PF}$ .

**Lemma 3.3.** *Let  $\mathcal{PF} = \mathcal{PF}_{q,K}(H, D, R_\omega)$ . Recall the auxiliary function classes defined in Section 2.2. Then, it holds that*

$$\begin{aligned}\mathbb{E}[\varepsilon_{sta}^1] &\leq 16|\mathcal{D}|d^2R_0K^2DH^{2D-2}R_\omega^{2D} \cdot \bar{\mathcal{R}}_N(\mathcal{E}_1^{\text{sub}}), \\ \mathbb{E}[\varepsilon_{sta}^2] &\leq 2|\partial\mathcal{D}|[K(H+1)R_\omega + 2R_0]K \cdot \bar{\mathcal{R}}_N(\mathcal{E}_2^{\text{sub}}).\end{aligned}$$

*Proof.* Note that  $\mathcal{E}'_1$  and  $\mathcal{E}'_2$  are exactly the function classes corresponding to the squared loss terms in  $\mathcal{J}_1$  and  $\mathcal{J}_2$ . Using standard symmetrization arguments (see, e.g., [48], Chapter 4), we conclude that

$$\mathbb{E}[\varepsilon_{sta}^1] \leq 2|\mathcal{D}| \cdot \mathcal{R}_N(\mathcal{E}'_1) \quad \text{and} \quad \mathbb{E}[\varepsilon_{sta}^2] \leq 2|\partial\mathcal{D}| \cdot \mathcal{R}_N(\mathcal{E}'_2).$$

The core of the proof is to relate  $\mathcal{R}_N(\mathcal{E}'_i)$  to the complexity of the much simpler sub-network class  $\mathcal{E}_i^{\text{sub}}$ . We show the derivation for the first term.

**Step 1: From squared loss to linear loss.** Let  $e(\cdot; \omega) \in \mathcal{E}_1$  and  $\tilde{e}(\cdot; \omega) = e(\cdot; \omega)^2 \in \mathcal{E}'_1$ . The function  $\Phi(x) = x^2$  is Lipschitz on the range of  $e$ . Specifically, for any two functions  $e_1, e_2 \in \mathcal{E}_1$ , we have  $|e_1^2 - e_2^2| = |e_1 + e_2||e_1 - e_2|$ . The term  $|e_1 + e_2|$  can be bounded using the properties of the function class  $\mathcal{PF}$  and Lemma A.4. This leads to a Lipschitz constant for  $\Phi$ , which by Talagrand's contraction inequality (see [40], Lemma 5.7) yields:

$$\mathcal{R}_N(\mathcal{E}'_1) \leq 8d^2R_0DKR_\omega^{2D}H^{2D-2} \cdot \mathcal{R}_N(\mathcal{E}_1).$$

A similar bound holds for the boundary term:

$$\mathcal{R}_N(\mathcal{E}'_2) \leq 2[K(H+1)R_\omega + R_0] \cdot \mathcal{R}_N(\mathcal{E}_2).$$

**Step 2: From the full network class to the sub-network class.** Finally, we relate  $\mathcal{R}_N(\mathcal{E}_i)$  to  $\bar{\mathcal{R}}_N(\mathcal{E}_i^{\text{sub}})$ . A function  $e \in \mathcal{E}_1$  has the form  $e(\mathbf{s}) = \sum_{j=1}^q \delta_j e_*(\mathbf{s}; \omega_j) - f(\mathbf{s})$ , where  $e_* \in \mathcal{E}_1^{\text{sub}}$  and  $\sum_j |\delta_j| \leq K$ . By the definition of Rademacher complexity and properties of the supremum, we get:

$$\begin{aligned}\mathcal{R}_N(\mathcal{E}_1) &= \mathbb{E}_{\{S_k, \eta_k\}} \left[ \sup_{\omega_{\text{all}}^q} \frac{1}{N} \sum_{k=1}^N \eta_k \left( \sum_{j=1}^q \delta_j e_*(S_k; \omega_j) - f(S_k) \right) \right] \\ &\leq \mathbb{E}_{\{S_k, \eta_k\}} \left[ \sup_{\omega_{\text{all}}^q} \sum_{j=1}^q |\delta_j| \cdot \left| \frac{1}{N} \sum_{k=1}^N \eta_k e_*(S_k; \omega_j) \right| \right] \\ &\leq K \cdot \mathbb{E}_{\{S_k, \eta_k\}} \left[ \sup_{e_* \in \mathcal{E}_1^{\text{sub}}} \left| \frac{1}{N} \sum_{k=1}^N \eta_k e_*(S_k) \right| \right] = K \cdot \bar{\mathcal{R}}_N(\mathcal{E}_1^{\text{sub}}).\end{aligned}$$

Similarly,  $\mathcal{R}_N(\mathcal{E}_2) \leq K \cdot \bar{\mathcal{R}}_N(\mathcal{E}_2^{\text{sub}})$ . Combining the results from both steps yields the statement of the lemma.  $\square$

**Remark 3.1.** This conclusion reveals an important fact: the Rademacher complexity of the full loss function class is controlled by  $K$  and the complexity of the elementary sub-network class  $\mathcal{E}_i^{\text{sub}}$ . This suggests that the network's overall complexity may not be affected by the number of sub-networks  $q$ , aiding us in managing the statistical error within the over-parameterized setting, where  $q$  can grow arbitrarily large.

Then, we introduce the definition of 'covering number' and Proposition 3.1, known as Dudley's entropy integral theorem (see [15]), which provides an effective tool to control  $\bar{\mathcal{R}}_N(\mathcal{E}_i^{\text{sub}})$ ,  $i=1,2$ .

**Definition 3.2.** An  $\epsilon$ -cover of a set  $T$  in a metric space  $(S, \tau)$  is a subset  $T_c \subset S$  such that for each  $t \in T$ , there exists a  $t_c \in T_c$  such that  $\tau(t, t_c) \leq \epsilon$ . The  $\epsilon$ -covering number of  $T$ , denoted as  $\mathcal{C}(\epsilon, T, \tau)$  is defined to be the minimum cardinality among all  $\epsilon$ -cover of  $T$  with respect to the metric  $\tau$ .

**Proposition 3.1.** For any function class  $\mathcal{F}$  mapping from  $\mathcal{D}$  to  $\mathbb{R}$  that contains the zero function and  $\|u\|_{L^\infty(\mathcal{D})} \leq \mathcal{B}$  for all  $u \in \mathcal{F}$ , it holds that

$$\bar{\mathcal{R}}_N(\mathcal{F}) \leq \inf_{0 < \delta < \mathcal{B}} \left( 4\delta + \frac{12}{\sqrt{N}} \int_\delta^{\mathcal{B}} \sqrt{\log \mathcal{C}(\epsilon, \mathcal{F}, \|\cdot\|_{L^\infty})} d\epsilon \right).$$

Based on above results, we present the following lemma, providing an upper bound for each  $\bar{\mathcal{R}}_N(\mathcal{E}_i^{\text{sub}})$ ,  $i=1,2$  in terms of the covering number of  $\mathcal{E}_i^{\text{sub}}$ .

**Lemma 3.4.** For  $i=1,2$ , we have

$$\bar{\mathcal{R}}_N(\mathcal{E}_i^{\text{sub}}) \leq C(H, D, d, R_0) R_\omega^{2D} \cdot N^{-1/2} \cdot \sqrt{\log(R_\omega H D N)}.$$

*Proof.* Using Proposition 3.1, for  $i=1,2$ ,

$$\bar{\mathcal{R}}_N(\mathcal{E}_i^{\text{sub}}) \leq \inf_{0 < \delta < B_i} \left[ 4\delta + \frac{12}{\sqrt{N}} \int_\delta^{B_i} \sqrt{\log \mathcal{C}(\epsilon, \mathcal{E}_i^{\text{sub}}, \|\cdot\|_{L^\infty})} d\epsilon \right].$$

Applying Lemma 5.6 in [27], we have

$$\mathcal{C}(\epsilon, \mathcal{E}_i^{\text{sub}}, \|\cdot\|_{L^\infty}) \leq \mathcal{C}(\epsilon L_i^{-1}, \Omega, \|\cdot\|_2).$$

By Lemma 5.5 in [27], it further holds that

$$\mathcal{C}(\epsilon L_i^{-1}, \Omega, \|\cdot\|_2) \leq \left[ \frac{2R_\omega \sqrt{H(H+1)D} \cdot L_i}{\epsilon} \right]^{H(H+1)D}.$$

The above  $B_i$  and  $L_i$  denote the boundedness and smoothness indices of  $\mathcal{E}_i^{\text{sub}}$  with respect to model parameters (see Lemma A.4 for details). For instance, to upper bound

$\bar{\mathcal{R}}_N(\mathcal{E}_1^{\text{sub}})$ , we substitute  $B_1$  and  $L_1$  from Lemma A.4 and obtain

$$\begin{aligned} \bar{\mathcal{R}}_N(\mathcal{E}_1^{\text{sub}}) &\leq \inf_{0 < \delta < B_1} \left[ 4\delta + \frac{12}{\sqrt{N}} \int_{\delta}^{B_1} \sqrt{\log \mathcal{C}(\epsilon, \mathcal{E}_1^{\text{sub}}, \|\cdot\|_{L^\infty})} d\epsilon \right] \\ &\leq \inf_{0 < \delta < B_1} \left[ 4\delta + 12H\sqrt{D} \cdot B_1 N^{-1/2} \log^{1/2}(44d^2 D^3 R_0 H^{3D+2} R_\omega^{2D} \delta^{-1}) \right]. \end{aligned}$$

Choosing  $\delta = N^{-1/2} < B_1/2$ , we have

$$\bar{\mathcal{R}}_N(\mathcal{E}_1^{\text{sub}}) \leq 36R_0 H^{2D} d^2 D^2 (3D+2) \cdot R_\omega^{2D} \cdot N^{-1/2} \cdot \sqrt{\log(R_\omega HDN)}.$$

Likewise,  $\bar{\mathcal{R}}_N(\mathcal{E}_2^{\text{sub}})$  is bounded similarly to  $\bar{\mathcal{R}}_N(\mathcal{E}_1^{\text{sub}})$ . □

Leveraging all preceding analyses, we can ultimately establish probabilistic bounds on  $\varepsilon_{sta}$  with high confidence:

**Theorem 3.4.** Let  $\mathcal{PF} = \mathcal{PF}_{q,K}(H, D, R_\omega)$ . For  $0 < \zeta < 1$ , with probability at least  $1 - \zeta$ , it holds that

$$\begin{aligned} \varepsilon_{sta} &= \sup_{v_{q,\omega} \in \mathcal{PF}} |\mathcal{J}(v_{q,\omega}) - \hat{\mathcal{J}}(v_{q,\omega})| \\ &\leq C_8 \cdot K^2 R_\omega^{4D} N^{-1/2} (\sqrt{\log(R_\omega HDN)} + \sqrt{\log \zeta^{-1}}), \end{aligned}$$

where  $C_8$  is a universal constant which only depends on  $\mathcal{D}, H, D, d$  and  $R_0$ .

*Proof.* By Lemma 3.2 and Lemma 3.3, we have that

$$\mathbb{E}[\varepsilon_{sta}] \leq 16|\mathcal{D}|d^2 R_0 K^2 D R_\omega^{2D} H^{2D-2} \cdot \bar{\mathcal{R}}_N(\mathcal{E}_1^{\text{sub}}) + 2|\partial\mathcal{D}|[K(H+1)R_\omega + 2R_0]K \cdot \bar{\mathcal{R}}_N(\mathcal{E}_2^{\text{sub}}).$$

By Lemma 3.4, for  $i = 1, 2$ ,

$$\bar{\mathcal{R}}_N(\mathcal{E}_i^{\text{sub}}) \leq C(H, D, d, R_0) R_\omega^{2D} \cdot N^{-1/2} \cdot \sqrt{\log(R_\omega HDN)},$$

then we can get

$$\mathbb{E}[\varepsilon_{sta}] \leq C_9(\mathcal{D}, R_0, d, H, D) \cdot K^2 R_\omega^{4D} N^{-1/2} \sqrt{\log(R_\omega HDN)}. \tag{3.9}$$

Now we define

$$\varphi(S_1, \dots, S_N, T_1, \dots, T_N) := \sup_{v_{q,\omega} \in \mathcal{PF}} |\mathcal{J}(v_{q,\omega}) - \hat{\mathcal{J}}(v_{q,\omega})| = \varepsilon_{sta}.$$

Further, we expand  $\mathcal{J}(v_{q,\omega})$  as follows

$$\mathcal{J}(v_{q,\omega}) = |\mathcal{D}| \mathbb{E}_{S \sim \mathcal{U}(\mathcal{D})} [\tilde{\varepsilon}(S; \omega_{\text{all}}^q)] + |\partial\mathcal{D}| \mathbb{E}_{T \sim \mathcal{U}(\partial\mathcal{D})} [\hat{\varepsilon}(T; \omega_{\text{all}}^q)],$$

where  $\tilde{\ell}(\cdot; \omega)$  and  $\hat{\ell}(\cdot; \omega)$  denote a function in  $\mathcal{E}'_1$  and  $\mathcal{E}'_2$ , respectively. Also, we expand  $\hat{\mathcal{J}}(v_{q,\omega})$  as follows

$$\hat{\mathcal{J}}(v_{q,\omega}) = \frac{|\mathcal{D}|}{N} \sum_{k=1}^N \tilde{\ell}(S_k; \omega_{\text{all}}^q) + \frac{|\partial\mathcal{D}|}{N} \sum_{k=1}^N \hat{\ell}(T_k; \omega_{\text{all}}^q).$$

We then examine the difference of  $\varphi(S_1, \dots, S_N, T_1, \dots, T_N)$ :

$$\begin{aligned} & |\varphi(S_1, \dots, S_i, \dots, T_N) - \varphi(S_1, \dots, S'_i, \dots, T_N)| \\ & \leq \frac{|\mathcal{D}|}{N} \sup_{\omega_{\text{all}}^q \in \Omega^q} |\tilde{\ell}(S_i; \omega_{\text{all}}^q) - \tilde{\ell}(S'_i; \omega_{\text{all}}^q)| \\ & = \frac{|\mathcal{D}|}{N} \sup_{\omega_{\text{all}}^q \in \Omega^q} \left| \left[ \sum_{j=1}^q \delta_j e_*(S_i; \omega_j) - f(S_i) \right]^2 - \left[ \sum_{j=1}^q \delta_j e_*(S'_i; \omega_j) - f(S'_i) \right]^2 \right| \\ & \leq 32 |\mathcal{D}| N^{-1} d^4 R_0^2 K^2 D^2 H^{4D-4} R_\omega^{4D}, \end{aligned}$$

where  $e_*(\cdot; \omega)$  denotes a function in  $\mathcal{E}_1^{\text{sub}}$ , and we have used the boundedness properties outlined in Lemma A.4. Similarly, we have

$$|\varphi(S_1, \dots, T_j, \dots, T_N) - \varphi(S_1, \dots, T'_j, \dots, T_N)| \leq 8 |\partial\mathcal{D}| N^{-1} R_0^2 K^2 (H+1)^2 R_\omega^2.$$

Then, by McDiarmid's inequality (see [37]), we have

$$\begin{aligned} \varepsilon_{sta} & \leq \mathbb{E}[\varepsilon_{sta}] + \epsilon \\ & \leq C_9(\mathcal{D}, R_0, d, H, D) \cdot K^2 R_\omega^{4D} N^{-1/2} \sqrt{\log(R_\omega H D N)} + \epsilon \end{aligned}$$

with probability as least

$$1 - 2 \exp \left\{ - \frac{N \epsilon^2}{C_{10} d^4 (|\partial\mathcal{D}|^4 + |\mathcal{D}|^2) R_0^4 H^{8D-8} R_\omega^{8D} K^4 R_0^4 D^4} \right\}.$$

This implies that with probability at least  $1 - \zeta$ , it holds that

$$\varepsilon_{sta} \leq C_8(\mathcal{D}, R_0, d, H, D) \cdot K^2 R_\omega^{4D} N^{-1/2} (\sqrt{\log(R_\omega H D N)} + \sqrt{\log \zeta^{-1}}).$$

This completes the proof.  $\square$

**Remark 3.2.** Theorem 3.4 provides a statistical error analysis for general  $\mathcal{PF}$  classes with a notable property: the error bound remains independent of the sub-network count  $q$ . This feature proves valuable in over-parameterized settings where  $q$  may grow arbitrarily large. Since Theorem 3.3 establishes that  $v_A \in \mathcal{PF}_{q, \bar{K}}(\bar{H}, \bar{D}, R_\omega + \zeta)$ , combining these results yields more refined bounds on  $\mathcal{E}_{sta}$ .

### 3.5 Main result

After analyzing the approximation, optimization, and statistical errors individually, we present our main theorem, which combines these components to establish the total error bound for PINNs in solving (2.1). This result addresses the question of how to determine the appropriate number of training samples, key architectural parameters of the neural networks, step size for the projected stochastic gradient descent optimization, and the required number of iterations to ensure that the gradient descent process closely approximates the true solution of (2.1) to a specified precision.

**Theorem 3.5.** *To solve the elliptic PDE (2.1) using PINNs, we employ the  $\mathcal{PF}$  architecture from Section 2.2 which consists of  $q$  sub-networks, each with width  $H$  and depth  $D$ . Sub-network parameters are initialized uniformly on  $[-B, B]$  via (2.5). We optimize using the PSGD algorithm from Section 2.3, with projection radii  $\zeta$  and  $\rho$  defined in (2.6) and (2.7). Let the Monte Carlo sample size in (2.4) be  $N$ . Let  $v_{\mathcal{A}}$  be the solution of the PSGD algorithm with iteration steps  $T_{\text{iter}}$  and step size  $\gamma$ . For any  $0 < \epsilon \ll 1$ , set*

$$\begin{aligned} q &= \lceil C \cdot \epsilon^{-\tilde{C}_1(\iota, d, \zeta_0, m)} \rceil, & H &= 2^{\lceil \log_2(d+m-1) \rceil + 1}, & D &= \lceil \log_2(d+m-1) \rceil + 2, \\ B &= C \cdot \epsilon^{-\frac{2d+2m}{m-2\iota-2}}, & \zeta &= \epsilon^{-\zeta}, & \rho &= C \cdot \epsilon^{-\frac{d}{m-2\iota-2}}, \\ T_{\text{iter}} &= C \cdot \epsilon^{-\tilde{C}_2(\iota, d, \zeta_0, m)}, & \gamma &= C \cdot \epsilon^{\tilde{C}_2(\iota, d, \zeta_0, m)}, & N &= \lceil C \cdot \epsilon^{-\tilde{C}_3(\iota, d, \zeta_0, m)} \rceil. \end{aligned}$$

Suppose that  $v_* \in C^m(\bar{\mathcal{D}})$  for  $m \geq 3$  is the target solution of Eq. (2.1). Then, with probability at least  $1 - 2 \cdot \epsilon^{\tilde{C}_3(\iota, d, \zeta_0, m)}$ , the total error

$$\mathbb{E}_{\text{SGD}} [\|v_{\mathcal{A}} - v_*\|_{H^{1/2}(\mathcal{D})}^2] \leq C\epsilon \log^{1/2}(C\epsilon^{-1}) = \tilde{\mathcal{O}}(\epsilon),$$

where

$$\begin{aligned} \zeta &> 0, & \zeta_0 &= \max\{\zeta, (2d+2m)(m-2\iota-2)^{-1}\}, \\ \tilde{C}_1(\iota, d, \zeta_0, m) &= \frac{C' \cdot (d+m)^3 \log^2(d+m-1)}{m-2\iota-2} + 19\zeta_0 \log_2(d+m-1) + 60\zeta_0, \\ \tilde{C}_2(\iota, d, \zeta_0, m) &= \frac{C'' \cdot (d+m)^3 \log^2(d+m-1)}{m-2\iota-2} + 27\zeta_0 \log_2(d+m-1) + 84\zeta_0, \\ \tilde{C}_3(\iota, d, \zeta_0, m) &= 8\zeta_0 \log_2(d+m-1) + \frac{4d}{m-2\iota-2} + 24\zeta_0 + 2. \end{aligned}$$

Meanwhile,  $C$  denotes a universal constant which is defined place by place and only depends on  $\mathcal{D}, H, D, d, R_0$  and  $m$ ;  $C'$  and  $C''$  are positive constants;  $0 < \iota < 1$  is an arbitrarily small positive number.

*Proof.* Reviewing the error decomposition in Theorem 3.1, we have

$$\begin{aligned} & \mathbb{E}_{\text{SGD}} [\|v_{\mathcal{A}} - v_*\|_{H^{1/2}(\mathcal{D})}^2] \\ & \leq C(\mathcal{D}, d, R_0) \left\{ \underbrace{\sup_{v \in \mathcal{P}\mathcal{F}} |\mathcal{J}(v) - \widehat{\mathcal{J}}(v)|}_{\varepsilon_{sta}} + \underbrace{\mathbb{E}_{\text{SGD}} [\widehat{\mathcal{J}}(v_{\mathcal{A}})] - \widehat{\mathcal{J}}(\bar{v})}_{\varepsilon_{opt}} + \underbrace{\|\bar{v} - v_*\|_{\mathcal{C}^2(\mathcal{D})}^2}_{\varepsilon_{app}} \right\}. \end{aligned}$$

We can divide the proof of the theorem into the following three steps:

**Step 1:** According to Theorem 3.2, we know that for any  $\epsilon > 0$ , there exists a neural network  $v_{\bar{q}, \bar{\omega}} \in \mathcal{P}\mathcal{F}_{\bar{q}, \bar{K}}(\bar{H}, \bar{D}, R_{\bar{\omega}})$ , with

$$\begin{aligned} \bar{q} &= \lceil C_1 \cdot \epsilon^{-\frac{d}{m-2l-2}} \rceil, \quad \bar{K} = C_2 \cdot \epsilon^{-\frac{d}{m-2l-2}}, \quad \bar{H} = 2^{\lceil \log_2(d+m-1) \rceil + 1}, \\ \bar{D} &= \lceil \log_2(d+m-1) \rceil + 2, \quad R_{\bar{\omega}} = C_3 \cdot \epsilon^{-\frac{2d+2m}{m-2l-2}}, \end{aligned}$$

such that the approximation error  $\varepsilon_{app} \leq C\epsilon^2 \leq C\epsilon$ .

**Step 2:** By Theorem 3.4, with probability at least  $1 - \zeta$ , the statistical error satisfies

$$\varepsilon_{sta} \leq C_8 \cdot \bar{K}^2 (R_{\bar{\omega}} + \zeta)^{4\bar{D}} N^{-1/2} \cdot \left\{ \log^{1/2}[(R_{\bar{\omega}} + \zeta)\bar{H}\bar{D}N] + \log^{1/2}(\zeta^{-1}) \right\}.$$

Setting

$$N = \lceil C \cdot \epsilon^{-\tilde{C}_3(t, d, \zeta_0, m)} \rceil, \quad \zeta = C \cdot \epsilon^{\tilde{C}_3(t, d, \zeta_0, m)}$$

with

$$\begin{aligned} \zeta_0 &= \max\{\zeta, (2d+2m)(m-2l-2)^{-1}\}, \\ \tilde{C}_3 &= 8\zeta_0 \log_2(d+m-1) + \frac{4d}{m-2l-2} + 24\zeta_0 + 2, \end{aligned}$$

it follows that, with probability at least  $1 - \zeta$ ,  $\varepsilon_{sta} \leq C\epsilon \log^{1/2}(C\epsilon^{-1}) = \tilde{\mathcal{O}}(\epsilon)$ .

**Step 3:** In order to bound the optimization error  $\varepsilon_{opt} \leq C\epsilon$  with probability at least  $1 - \zeta$ , we need to determine parameters  $G, J, \tau, T_{\text{iter}}, \rho$  such that the following inequalities hold

$$1 - \bar{q}G \left[ 1 - \tau^{\bar{H}(\bar{H}+1)\bar{D}} (2R_{\bar{\omega}})^{-\bar{H}(\bar{H}+1)\bar{D}} \right]^J \geq 1 - \zeta, \quad (3.10)$$

$$C_4 \bar{K}^2 R_{\bar{\omega}}^{5\bar{D}} \tau + \frac{C_7 \bar{K}^2 (R_{\bar{\omega}} + \zeta)^{5\bar{D}} \zeta}{\sqrt{G}} + \frac{C_5 \bar{q} \bar{K}^2 (R_{\bar{\omega}} + \zeta)^{8\bar{D}}}{2T_{\text{iter}}} \leq C\epsilon. \quad (3.11)$$

To ensure the first and third terms in (3.11) are bounded by  $C\epsilon$ , we require

$$\begin{aligned} \tau &\leq C \cdot \epsilon^{\frac{10(d+m)\log_2(d+m-1) + 23d + 20m}{m-2l-2} + 1}, \\ G &\geq \lceil C \cdot \epsilon^{-10\log_2(d+m-1)\zeta_0 - 32\zeta_0 - \frac{6d}{m-2l-2} - 2} \rceil. \end{aligned}$$

Since  $1 - x \leq \exp(-x)$ , and  $\exp(-x) \leq x^{-2}$  for  $x \geq 0$ , we could solve for  $J$  as follows:

$$J \geq \sqrt{\frac{\bar{q}G}{\zeta}} \left( \frac{\tau}{2R_{\bar{\omega}}} \right)^{-\bar{H}(\bar{H}+1)\bar{D}} \Rightarrow J \geq \lceil C\epsilon^{-\frac{C_0(d+m)^3 \log^2(d+m-1)}{m-2l-2} - 9\zeta_0 \log_2(d+m-1) - 28\zeta_0} \rceil,$$

where  $C_0$  is a positive constant. Then,  $q = \bar{q} \cdot G \cdot J \geq \lceil C \cdot \epsilon^{-\tilde{C}_1(l,d,\zeta_0,m)} \rceil$ . Finally, we turn to the second term of (3.11). To ensure that

$$\frac{C_5 \cdot q \cdot \bar{K}^2 \cdot (R_{\bar{\omega}} + \zeta)^{8\bar{D}}}{2T_{\text{iter}}} \leq C\epsilon,$$

we require

$$T_{\text{iter}} \geq \lceil C \cdot \epsilon^{-\tilde{C}_2(l,d,\zeta_0,m)} \rceil.$$

By Theorem 3.3,  $\gamma = T_{\text{iter}}^{-1}$ . Thus, we complete the proof.  $\square$

**Remark 3.3.** The projection radii in our method,  $\zeta$  and  $\rho$ , diverge with the target accuracy  $\epsilon$ , permitting significant parameter updates during optimization. This behavior is distinct from both the minimal parameter changes inherent to NTK theory [24] and the static, frozen inner-layer parameters of random feature models [10].

**Remark 3.4.** Our convergence analysis relies on a sufficiently large number of subnetworks,  $q$  (over-parameterization), to guarantee that with high probability, a randomly initialized subnetwork is close enough to the target function for optimization to succeed. This leads to a theoretical requirement where  $q$  scales polynomially with the inverse accuracy,  $\epsilon^{-1}$ . It is crucial to recognize this as a ‘worst-case’ sufficiency condition that may be looser than what is needed in practice. In line with the broader deep learning literature, practitioners should view this bound as a foundational guideline for model scaling, not a strict prescription. A dedicated empirical study of this scaling relationship remains a valuable direction for our future research.

## 4 Conclusion

In this paper, we provide the first complete error analysis for PINNs that includes the approximation error, statistical error, and optimization error in the scenario of over-parameterization. Our analysis is based on the PSGD algorithm and does not require constraining the neural network weights near their initial values during the optimization process. This marks a milestone in the field of theoretical understanding of solving PDEs via PINNs.

## Acknowledgments

This work is supported by the National Natural Science Foundation of China (No. 123B2019, No. 12125103, No. U24A2002, No. 12371441), and by the Fundamental Research Funds for the Central Universities.

## A Some properties of the neural networks

Lemmas A.1–A.4 establish several properties of  $\mathcal{F}(H, D, R_\omega)$  and  $\mathcal{E}_i^{\text{sub}}$  defined in Section 2.2. Lemmas A.1 and A.2 can be found in [27], and the proof of Lemma A.3 is provided in Appendix B. Below, we always assume that  $\omega$  and  $\tilde{\omega}$  are two arbitrary parameter vectors satisfying

$$\omega, \tilde{\omega} \in \Omega = [-R_\omega, R_\omega]^{\mathfrak{G}(H, D, d)}, \quad R_\omega \geq 1.$$

**Lemma A.1.** *For any  $\Phi_\omega \in \mathcal{F}(H, D, R_\omega)$  and  $\mathbf{s} \in \mathcal{D}$ , we have  $|\Phi_\omega(\mathbf{s})| \leq (H+1)R_\omega$ . Moreover, for any  $\Phi_\omega, \Phi_{\tilde{\omega}} \in \mathcal{F}(H, D, R_\omega)$  and  $\mathbf{s} \in \mathcal{D}$ ,*

$$|\Phi_\omega(\mathbf{s}) - \Phi_{\tilde{\omega}}(\mathbf{s})| \leq 2H^D \sqrt{D} \cdot R_\omega^{D-1} \|\omega - \tilde{\omega}\|_2.$$

**Lemma A.2.** *For any  $\Phi_\omega \in \mathcal{F}(H, D, R_\omega)$  and  $\mathbf{s} \in \mathcal{D}$ , we have*

$$|\partial_{s_\alpha} \Phi_\omega(\mathbf{s})| \leq H^{D-1} R_\omega^D.$$

Moreover, for any  $\Phi_\omega, \Phi_{\tilde{\omega}} \in \mathcal{F}(H, D, R_\omega)$  and  $\mathbf{s} \in \mathcal{D}$ ,  $\alpha \in [d]$ ,

$$|\partial_{s_\alpha} \Phi_\omega(\mathbf{s}) - \partial_{s_\alpha} \Phi_{\tilde{\omega}}(\mathbf{s})| \leq 4H^{2D-1} D^{3/2} \cdot R_\omega^{2D} \|\omega - \tilde{\omega}\|_2.$$

**Lemma A.3.** *For any  $\Phi_\omega \in \mathcal{F}(H, D, R_\omega)$  and  $\mathbf{s} \in \mathcal{D}$ , we have*

$$|\partial_{s_\alpha} \partial_{s_\beta} \Phi_\omega(\mathbf{s})| \leq DH^{2D-2} R_\omega^{2D}.$$

Moreover, for any  $\Phi_\omega, \Phi_{\tilde{\omega}} \in \mathcal{F}(H, D, R_\omega)$  and  $\mathbf{s} \in \mathcal{D}$ ,  $\alpha, \beta \in [d]$ ,

$$|\partial_{s_\alpha} \partial_{s_\beta} \Phi_\omega(\mathbf{s}) - \partial_{s_\alpha} \partial_{s_\beta} \Phi_{\tilde{\omega}}(\mathbf{s})| \leq 18H^{3D+1} D^{5/2} \cdot R_\omega^{3D} \|\omega - \tilde{\omega}\|_2.$$

As a direct result of Lemmas A.1 to A.3, we have

**Lemma A.4.** *Let  $e_i(\cdot; \omega), e_i(\cdot; \tilde{\omega}) \in \mathcal{E}_i^{\text{sub}}$  for  $i=1, 2$ . Then for all  $\mathbf{s} \in \mathcal{D}$ , they satisfy both boundedness and Lipschitz conditions with respect to parameters:*

$$|e_i(\mathbf{s}; \omega)| \leq B_i, \quad |e_i(\mathbf{s}; \omega) - e_i(\mathbf{s}; \tilde{\omega})| \leq L_i \|\omega - \tilde{\omega}\|_2,$$

where the constants are given by:

$$\begin{aligned} B_1 &= 3d^2 R_0 D H^{2D-2} \cdot R_\omega^{2D}, & B_2 &= (H+1) \cdot R_\omega, \\ L_1 &= 22d^2 R_0 D^{5/2} H^{3D+1} \cdot R_\omega^{3D}, & L_2 &= 2H^D \sqrt{D} \cdot R_\omega^{D-1}. \end{aligned}$$

## B Proof of Lemma A.3

Note that the activation function  $\sigma = \tanh$ , which is 1-Lipschitz and has a 1-Lipschitz continuous gradient. For  $\ell = 1, 2, \dots, D$ ,

$$\Phi_i^{[\ell]} = \sigma \left( \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \Phi_j^{[\ell-1]} + c_i^{[\ell-1]} \right).$$

By Lemma A.2, we have

$$|\partial_{s_\alpha} \Phi_i^{[\ell]}| = \left| \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\alpha} \Phi_j^{[\ell-1]} \sigma' \left( \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \Phi_j^{[\ell-1]} + c_i^{[\ell-1]} \right) \right| \leq H^{\ell-1} R_\omega^\ell.$$

Then, it holds that

$$\begin{aligned} |\partial_{s_\alpha} \partial_{s_\beta} \Phi_i^{[\ell]}| &= \left| \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\alpha} \partial_{s_\beta} \Phi_j^{[\ell-1]} \sigma' \left( \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \Phi_j^{[\ell-1]} + c_i^{[\ell-1]} \right) \right. \\ &\quad \left. + \left[ \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\alpha} \Phi_j^{[\ell-1]} \right] \left[ \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\beta} \Phi_j^{[\ell-1]} \right] \sigma'' \left( \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \Phi_j^{[\ell-1]} + c_i^{[\ell-1]} \right) \right| \\ &\leq R_\omega \sum_{j=1}^{H_{\ell-1}} \partial_{s_\alpha} \partial_{s_\beta} \Phi_j^{[\ell-1]} + \left[ R_\omega \sum_{j=1}^{H_{\ell-1}} \partial_{s_\alpha} \Phi_j^{[\ell-1]} \right]^2 \\ &= R_\omega \sum_{j=1}^{H_{\ell-1}} \partial_{s_\alpha} \partial_{s_\beta} \Phi_j^{[\ell-1]} + R_\omega^2 H^{2\ell-2} \\ &\leq \dots \leq R_\omega^{\ell+1} H^{\ell-1} + R_\omega^{\ell+2} H^\ell + \dots + R_\omega^{2\ell-1} H^{2\ell-3} + R_\omega^{2\ell} H^{2\ell-2} \\ &\leq \ell R_\omega^{2\ell} H^{2\ell-2}. \end{aligned}$$

Note that  $H_0 = d$ ,  $H_D = 1$ . For  $\ell = 1$ , we have

$$\begin{aligned} &|\partial_{s_\alpha} \partial_{s_\beta} \Phi_i^{[1]} - \partial_{s_\alpha} \partial_{s_\beta} \tilde{\Phi}_i^{[1]}| \\ &= \left| w_{i\alpha}^{[0]} w_{i\beta}^{[0]} \sigma'' \left( \sum_{j=1}^{H_0} w_{ij}^{[0]} s_j + c_i^{[0]} \right) - \tilde{w}_{i\alpha}^{[0]} \tilde{w}_{i\beta}^{[0]} \sigma'' \left( \sum_{j=1}^{H_0} \tilde{w}_{ij}^{[0]} s_j + \tilde{c}_i^{[0]} \right) \right| \\ &\leq |w_{i\alpha}^{[0]} - \tilde{w}_{i\alpha}^{[0]}| |\tilde{w}_{i\beta}^{[0]}| \left| \sigma'' \left( \sum_{j=1}^{H_0} w_{ij}^{[0]} s_j + c_i^{[0]} \right) \right| \\ &\quad + |w_{i\beta}^{[0]} - \tilde{w}_{i\beta}^{[0]}| |w_{i\alpha}^{[0]}| \left| \sigma'' \left( \sum_{j=1}^{H_0} w_{ij}^{[0]} s_j + c_i^{[0]} \right) \right| \\ &\quad + |\tilde{w}_{i\alpha}^{[0]}| |\tilde{w}_{i\beta}^{[0]}| \left| \sigma'' \left( \sum_{j=1}^{H_0} w_{ij}^{[0]} s_j + c_i^{[0]} \right) - \sigma'' \left( \sum_{j=1}^{H_0} \tilde{w}_{ij}^{[0]} s_j + \tilde{c}_i^{[0]} \right) \right| \end{aligned}$$

$$\begin{aligned}
&\leq |w_{i\alpha}^{[0]} - \tilde{w}_{i\alpha}^{[0]}| R_\omega + |w_{i\beta}^{[0]} - \tilde{w}_{i\beta}^{[0]}| R_\omega \\
&\quad + \left| \sigma'' \left( \sum_{j=1}^{H_0} w_{ij}^{[0]} s_j + c_i^{[0]} \right) - \sigma'' \left( \sum_{j=1}^{H_0} \tilde{w}_{ij}^{[0]} s_j + \tilde{c}_i^{[0]} \right) \right| R_\omega^2 \\
&\leq (2R_\omega + 2R_\omega^2) \sum_{k=1}^{n_1} |\omega_k - \tilde{\omega}_k|,
\end{aligned}$$

where we use  $\{\omega_k\}_{k=1}^{n_1}$  to denote all the parameters  $w_{ij}^{[0]}, c_i^{[0]}$ . For  $2 \leq \ell \leq D$ , it holds that

$$\begin{aligned}
&|\partial_{s_\alpha} \partial_{s_\beta} \Phi_i^{[\ell]} - \partial_{s_\alpha} \partial_{s_\beta} \tilde{\Phi}_i^{[\ell]}| \\
&\leq \left| \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\alpha} \partial_{s_\beta} \Phi_j^{[\ell-1]} \sigma' \left( \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \Phi_j^{[\ell-1]} + c_i^{[\ell-1]} \right) \right. \\
&\quad \left. - \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \partial_{s_\alpha} \partial_{s_\beta} \tilde{\Phi}_j^{[\ell-1]} \sigma' \left( \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \tilde{\Phi}_j^{[\ell-1]} + \tilde{c}_i^{[\ell-1]} \right) \right| \\
&\quad + \left| \left[ \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\alpha} \Phi_j^{[\ell-1]} \right] \left[ \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\beta} \Phi_j^{[\ell-1]} \right] \sigma'' \left( \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \Phi_j^{[\ell-1]} + c_i^{[\ell-1]} \right) \right. \\
&\quad \left. - \left[ \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \partial_{s_\alpha} \tilde{\Phi}_j^{[\ell-1]} \right] \left[ \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \partial_{s_\beta} \tilde{\Phi}_j^{[\ell-1]} \right] \sigma'' \left( \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \tilde{\Phi}_j^{[\ell-1]} + \tilde{c}_i^{[\ell-1]} \right) \right|.
\end{aligned}$$

For the first term, we have

$$\begin{aligned}
&\left| \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\alpha} \partial_{s_\beta} \Phi_j^{[\ell-1]} \sigma' \left( \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \Phi_j^{[\ell-1]} + c_i^{[\ell-1]} \right) \right. \\
&\quad \left. - \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \partial_{s_\alpha} \partial_{s_\beta} \tilde{\Phi}_j^{[\ell-1]} \sigma' \left( \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \tilde{\Phi}_j^{[\ell-1]} + \tilde{c}_i^{[\ell-1]} \right) \right| \\
&\leq \left| \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\alpha} \partial_{s_\beta} \Phi_j^{[\ell-1]} \sigma' \left( \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \Phi_j^{[\ell-1]} + c_i^{[\ell-1]} \right) \right. \\
&\quad \left. - \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\alpha} \partial_{s_\beta} \Phi_j^{[\ell-1]} \sigma' \left( \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \tilde{\Phi}_j^{[\ell-1]} + \tilde{c}_i^{[\ell-1]} \right) \right| \\
&\quad + \left| \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\alpha} \partial_{s_\beta} \Phi_j^{[\ell-1]} \sigma' \left( \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \tilde{\Phi}_j^{[\ell-1]} + \tilde{c}_i^{[\ell-1]} \right) \right. \\
&\quad \left. - \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\alpha} \partial_{s_\beta} \tilde{\Phi}_j^{[\ell-1]} \sigma' \left( \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \tilde{\Phi}_j^{[\ell-1]} + \tilde{c}_i^{[\ell-1]} \right) \right|
\end{aligned}$$

$$\begin{aligned}
& + \left| \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\alpha} \partial_{s_\beta} \tilde{\Phi}_j^{[\ell-1]} \sigma' \left( \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \tilde{\Phi}_j^{[\ell-1]} + \tilde{c}_i^{[\ell-1]} \right) \right. \\
& \quad \left. - \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \partial_{s_\alpha} \partial_{s_\beta} \tilde{\Phi}_j^{[\ell-1]} \sigma' \left( \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \tilde{\Phi}_j^{[\ell-1]} + \tilde{c}_i^{[\ell-1]} \right) \right| \\
& \leq (\ell-1) R_\omega^{2\ell-1} H^{2\ell-3} \left( \left| \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} - \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \right| + |c_i^{[\ell-1]} - \tilde{c}_i^{[\ell-1]}| \right) \\
& \quad + (\ell-1) R_\omega^{2\ell} H^{2\ell-2} \max_j |\Phi_j^{[\ell-1]} - \tilde{\Phi}_j^{[\ell-1]}| + R_\omega H \cdot \max_j |\partial_{s_\alpha} \partial_{s_\beta} \Phi_j^{[\ell-1]} - \partial_{s_\alpha} \partial_{s_\beta} \tilde{\Phi}_j^{[\ell-1]}| \\
& \quad + \left| \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} - \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \right| \cdot \max_j |\partial_{s_\alpha} \partial_{s_\beta} \tilde{\Phi}_j^{[\ell-1]}| \\
& \leq (\ell-1) (R_\omega^{2\ell-1} H^{2\ell-3} + R_\omega^{3\ell-2} H^{3\ell-4} + R_\omega^{2\ell-2} H^{2\ell-4}) \sum_{k=1}^{n_{\ell-1}} |\omega_k - \tilde{\omega}_k| \\
& \quad + R_\omega H \cdot \max_j |\partial_{s_\alpha} \partial_{s_\beta} \Phi_j^{[\ell-1]} - \partial_{s_\alpha} \partial_{s_\beta} \tilde{\Phi}_j^{[\ell-1]}|,
\end{aligned}$$

where in the final step above, we apply the natural extension of Lemma A.1 to the intermediate layer, and we use  $\{\omega_k\}_{k=1}^{n_{\ell-1}}$  to denote all the parameters from the first  $\ell-1$  layers. Similarly,

$$\begin{aligned}
& \left| \left[ \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\alpha} \Phi_j^{[\ell-1]} \right] \left[ \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \partial_{s_\beta} \Phi_j^{[\ell-1]} \right] \sigma'' \left( \sum_{j=1}^{H_{\ell-1}} w_{ij}^{[\ell-1]} \Phi_j^{[\ell-1]} + c_i^{[\ell-1]} \right) \right. \\
& \quad \left. - \left[ \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \partial_{s_\alpha} \tilde{\Phi}_j^{[\ell-1]} \right] \left[ \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \partial_{s_\beta} \tilde{\Phi}_j^{[\ell-1]} \right] \sigma'' \left( \sum_{j=1}^{H_{\ell-1}} \tilde{w}_{ij}^{[\ell-1]} \tilde{\Phi}_j^{[\ell-1]} + \tilde{c}_i^{[\ell-1]} \right) \right| \\
& \leq (R_\omega^{2\ell} H^{2\ell-2} + R_\omega^{3\ell-1} H^{3\ell-3} + 2\ell R_\omega^{3\ell-1} H^{3\ell-4} + 2R_\omega^{2\ell-1} H^{2\ell-3}) \sum_{k=1}^{n_{\ell-1}} |\omega_k - \tilde{\omega}_k| \\
& \leq 9\ell H^{3\ell} R_\omega^{3\ell} \sum_{k=1}^{n_{\ell-1}} |\omega_k - \tilde{\omega}_k|,
\end{aligned}$$

then, it holds that

$$\begin{aligned}
|\partial_{s_\alpha} \partial_{s_\beta} \Phi_i^{[\ell]} - \partial_{s_\alpha} \partial_{s_\beta} \tilde{\Phi}_i^{[\ell]}| & \leq R_\omega H \cdot \max_j |\partial_{s_\alpha} \partial_{s_\beta} \Phi_j^{[\ell-1]} - \partial_{s_\alpha} \partial_{s_\beta} \tilde{\Phi}_j^{[\ell-1]}| \\
& \quad + 9\ell H^{3\ell} R_\omega^{3\ell} \sum_{k=1}^{n_{\ell-1}} |\omega_k - \tilde{\omega}_k| \leq \dots \leq 9\ell^2 H^{3\ell} R_\omega^{3\ell} \sum_{k=1}^{n_\ell} |\omega_k - \tilde{\omega}_k|.
\end{aligned}$$

And by  $\sum_{k=1}^{n_D} |\omega_k - \tilde{\omega}_k| \leq \sqrt{H(H+1)D} \|\omega - \tilde{\omega}\|_2$ ,

$$|\partial_{s_\alpha} \partial_{s_\beta} \Phi_i^{[D]} - \partial_{s_\alpha} \partial_{s_\beta} \tilde{\Phi}_i^{[D]}| \leq 18H^{3D+1} \cdot R_\omega^{3D} D^2 \sqrt{D} \|\omega - \tilde{\omega}\|_2, \quad \forall s \in \mathcal{D}.$$

### C Proof of Lemma 3.1

Note that  $\mathbf{y}^* \in Y$ . For any  $i \in [N]$ , by convexity of  $\mathbf{y} \mapsto F_i(\mathbf{x}_t, \mathbf{y})$ , and (3.6), we have

$$\begin{aligned} F_{i_{t+1}}(\mathbf{x}_t, \mathbf{y}_t) - F_{i_{t+1}}(\mathbf{x}_t, \mathbf{y}^*) &\leq \langle \nabla_{\mathbf{y}} F_{i_{t+1}}(\mathbf{x}_t, \mathbf{y}_t), \mathbf{y}_t - \mathbf{y}^* \rangle \\ &= \frac{1}{2\lambda} \cdot 2 \cdot \langle \lambda \nabla_{\mathbf{y}} F_{i_{t+1}}(\mathbf{x}_t, \mathbf{y}_t), \mathbf{y}_t - \mathbf{y}^* \rangle \\ &= \frac{1}{2\lambda} \left[ -\|\mathbf{y}_t - \lambda \nabla_{\mathbf{y}} F_{i_{t+1}}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{y}^*\|_2^2 + \|\mathbf{y}_t - \mathbf{y}^*\|_2^2 + \|\lambda \nabla_{\mathbf{y}} F_{i_{t+1}}(\mathbf{x}_t, \mathbf{y}_t)\|_2^2 \right] \\ &\leq \frac{1}{2\lambda} \left[ -\|\mathbf{y}_{t+1} - \mathbf{y}^*\|_2^2 + \|\mathbf{y}_t - \mathbf{y}^*\|_2^2 + \lambda^2 V^2 \right], \quad \forall t = 0, \dots, T-1. \end{aligned}$$

Taking conditional expectation  $\mathbb{E}[\cdot | i_1, \dots, i_t]$  on both sides, we get

$$F(\mathbf{x}_t, \mathbf{y}_t) - F(\mathbf{x}_t, \mathbf{y}^*) \leq \frac{1}{2\lambda} \left\{ \|\mathbf{y}_t - \mathbf{y}^*\|_2^2 - \mathbb{E}[\|\mathbf{y}_{t+1} - \mathbf{y}^*\|_2^2 | i_1, \dots, i_t] + \lambda^2 V^2 \right\}.$$

By further taking the full expectation, we obtain

$$\mathbb{E}[F(\mathbf{x}_t, \mathbf{y}_t)] - \mathbb{E}[F(\mathbf{x}_t, \mathbf{y}^*)] \leq \frac{1}{2\lambda} \left\{ \mathbb{E}[\|\mathbf{y}_t - \mathbf{y}^*\|_2^2] - \mathbb{E}[\|\mathbf{y}_{t+1} - \mathbf{y}^*\|_2^2] + \lambda^2 V^2 \right\}.$$

Since  $\lambda T = 1$ , this implies

$$\frac{1}{T} \sum_{t=0}^{T-1} \left\{ \mathbb{E}[F(\mathbf{x}_t, \mathbf{y}_t)] - \mathbb{E}[F(\mathbf{x}_t, \mathbf{y}^*)] \right\} \leq \frac{\|\mathbf{y}_0 - \mathbf{y}^*\|_2^2}{2} + \frac{V^2}{2T}.$$

Using above results, we have

$$\begin{aligned} &\min_{t=0, \dots, T} \mathbb{E}[F(\mathbf{x}_t, \mathbf{y}_t)] \\ &\leq \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[F(\mathbf{x}_t, \mathbf{y}^*)] + \frac{\|\mathbf{y}^* - \mathbf{y}_0\|_2^2}{2} + \frac{V^2}{2T} \\ &\leq F(\mathbf{x}_0, \mathbf{y}^*) + \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[|F(\mathbf{x}_t, \mathbf{y}^*) - F(\mathbf{x}_0, \mathbf{y}^*)|] + \frac{\|\mathbf{y}^* - \mathbf{y}_0\|_2^2}{2} + \frac{V^2}{2T}. \end{aligned}$$

Finally, by (3.7) we get

$$\min_{t=0, \dots, T} \mathbb{E}[F(\mathbf{x}_t, \mathbf{y}_t)] \leq F(\mathbf{x}_0, \mathbf{y}^*) + U \|\mathbf{y}^*\|_2 \text{diam}(X) + \frac{\|\mathbf{y}^* - \mathbf{y}_0\|_2^2}{2} + \frac{V^2}{2T}.$$

### D Proof of Theorem 3.2

The proof follows the constructive approach detailed in Theorem 3.2 of [28]. We establish an approximation bound in the  $C^2(\bar{\mathcal{D}})$  norm for a target function  $v_* \in C^m(\bar{\mathcal{D}})$ .

**Step 1: Preliminary Lemmas.**

We first state two preliminary lemmas that are essential for the error analysis in the  $W^{2,\infty}$  space. The first lemma establishes that shallow neural networks can approximate basic monomials in the  $W^{2,\infty}$  norm, see Proposition 4.7 in [18].

**Lemma D.1.** *Let  $M > 0$ , and let  $\sigma(x) = \tanh(x)$ . Suppose  $x_0 \in \mathbb{R}$  is a point such that  $\sigma^{(r)}(x_0) \neq 0$  for  $r \in \{1,2,3\}$ . For any  $\epsilon \in (0,1)$ , there exist  $\Phi_{\omega_1} \in \mathcal{F}(1,2,C_1(M)\epsilon^{-1})$  and  $\Phi_{\omega_2} \in \mathcal{F}(2,2,C_2(M)\epsilon^{-1})$  such that*

$$\|x - \Phi_{\omega_1}(x)\|_{W^{2,\infty}([-M,M])} \leq \epsilon \quad \text{and} \quad \|x^2 - \Phi_{\omega_2}(x)\|_{W^{2,\infty}([-M,M])} \leq \epsilon.$$

The second lemma provides necessary calculus rules for functions in the  $W^{2,\infty}$  space.

**Lemma D.2.** *Let  $d_1, d_2 \in \mathbb{N}$ , and let  $\mathcal{D}_1 \subset \mathbb{R}^{d_1}$ ,  $\mathcal{D}_2 \subset \mathbb{R}^{d_2}$  be open, bounded, and convex domains. Then there exists  $C_1 = C(d_1, d_2) > 0$  and  $C_2 = C(d_1) > 0$  such that:*

1. **(Chain rule).** *Let  $f \in W^{2,\infty}(\mathcal{D}_1; \mathbb{R}^{d_2})$  and  $g \in W^{2,\infty}(\mathcal{D}_2)$  be Lipschitz functions with  $\text{range}(f) \subset \mathcal{D}_2$ . Then  $g \circ f \in W^{2,\infty}(\mathcal{D}_1)$ , and*

$$\|g \circ f\|_{W^{2,\infty}(\mathcal{D}_1)} \leq C_1 \max \left\{ \|g\|_{W^{1,\infty}(\mathcal{D}_2)} \|f\|_{W^{2,\infty}(\mathcal{D}_1; \mathbb{R}^{d_2})}, \|g\|_{W^{2,\infty}(\mathcal{D}_2)} \cdot \|f\|_{W^{1,\infty}(\mathcal{D}_1; \mathbb{R}^{d_2})}^2 \right\}. \tag{D.1}$$

2. **(Product rule).** *Let  $u, v \in W^{2,\infty}(\mathcal{D}_1)$ . Then  $uv \in W^{2,\infty}(\mathcal{D}_1)$ , and*

$$\|uv\|_{W^{2,\infty}(\mathcal{D}_1)} \leq C_2 \|u\|_{W^{2,\infty}(\mathcal{D}_1)} \|v\|_{W^{2,\infty}(\mathcal{D}_1)}. \tag{D.2}$$

**Proof of the Chain Rule.** Let  $h(x) = g(f(x))$ . The zeroth-order term is bounded by  $\|h\|_{L^\infty(\mathcal{D}_1)} \leq \|g\|_{L^\infty(\mathcal{D}_2)}$ . The first-order derivatives, given by  $\partial_{x_i} h = \sum_j ((\partial_{y_j} g) \circ f) \cdot (\partial_{x_i} f_j)$ , are bounded by  $\|\partial_{x_i} h\|_{L^\infty} \leq d_2 \|g\|_{W^{1,\infty}(\mathcal{D}_2)} \|f\|_{W^{1,\infty}(\mathcal{D}_1)}$ . Further, we have

$$\partial_{x_k} \partial_{x_i} h = \sum_{j=1}^{d_2} \left[ \left( \sum_{l=1}^{d_2} (\partial_{y_l} \partial_{y_j} g) \circ f \cdot \partial_{x_k} f_l \right) \cdot (\partial_{x_i} f_j) \right] + \sum_{j=1}^{d_2} \left[ ((\partial_{y_j} g) \circ f) \cdot (\partial_{x_k} \partial_{x_i} f_j) \right],$$

which leads to

$$\begin{aligned} \|\partial_{x_k} \partial_{x_i} h\|_{L^\infty} &\leq \sum_{j,l} \|\partial_{y_l} \partial_{y_j} g\|_{L^\infty(\mathcal{D}_2)} \|\partial_{x_k} f_l\|_{L^\infty(\mathcal{D}_1)} \|\partial_{x_i} f_j\|_{L^\infty(\mathcal{D}_1)} \\ &\quad + \sum_j \|\partial_{y_j} g\|_{L^\infty(\mathcal{D}_2)} \|\partial_{x_k} \partial_{x_i} f_j\|_{L^\infty(\mathcal{D}_1)} \\ &\leq d_2^2 \|g\|_{W^{2,\infty}(\mathcal{D}_2)} \|f\|_{W^{1,\infty}(\mathcal{D}_1)}^2 + d_2 \|g\|_{W^{1,\infty}(\mathcal{D}_2)} \|f\|_{W^{2,\infty}(\mathcal{D}_1)}. \end{aligned}$$

Note that the bound for the second derivative dominates the lower-order terms. This proves the chain rule equation D.1.

**Proof of the Product Rule.** Let  $h(x) = u(x)v(x)$ . Directly, we have

$$\|h\|_{L^\infty} \leq \|u\|_{L^\infty} \|v\|_{L^\infty} \leq \|u\|_{W^{2,\infty}(\mathcal{D}_1)} \|v\|_{W^{2,\infty}(\mathcal{D}_1)}.$$

For the first-order derivatives,  $\partial_i h = (\partial_i u)v + u(\partial_i v)$ , it holds that

$$\|\partial_i h\|_{L^\infty} \leq \|\partial_i u\|_{L^\infty} \|v\|_{L^\infty} + \|u\|_{L^\infty} \|\partial_i v\|_{L^\infty} \leq 2\|u\|_{W^{2,\infty}(\mathcal{D}_1)} \|v\|_{W^{2,\infty}(\mathcal{D}_1)}.$$

For the second-order derivatives,  $\partial_j \partial_i h = (\partial_j \partial_i u)v + (\partial_i u)(\partial_j v) + (\partial_j u)(\partial_i v) + u(\partial_j \partial_i v)$ . Then, we have

$$\|\partial_j \partial_i h\|_{L^\infty} \leq \|\partial_j \partial_i u\|_{L^\infty} \|v\|_{L^\infty} + \dots + \|u\|_{L^\infty} \|\partial_j \partial_i v\|_{L^\infty} \leq 4\|u\|_{W^{2,\infty}(\mathcal{D}_1)} \|v\|_{W^{2,\infty}(\mathcal{D}_1)}.$$

Since  $\|uv\|_{W^{2,\infty}(\mathcal{D}_1)} = \max_{\|\alpha\|_1 \leq 2} \|\partial^\alpha (uv)\|_{L^\infty}$ , taking the maximum over all derivative orders allows us to choose  $C_2 = 4$  to satisfy the inequality equation (D.2).

### Step 2: Local Polynomial Approximation.

We first approximate the target function  $v_*$  with a function  $v_N$  composed of localized polynomials, in direct analogy to Proposition 4.1 of [18]. The function  $v_N$  has the form:

$$v_N := \sum_{m \in \{0, \dots, N\}^d} \sum_{\|\alpha\|_1 \leq m-1} c_{v_*, m, \alpha} \Psi_m^s \mathbf{x}^\alpha,$$

where the functions  $\Psi_m^s$  constitute an tanh-related exponential partition of unity as specified in Lemma 1 of [28]. Also, we have

$$\|v_* - v_N\|_{W^{2,\infty}(\mathcal{D})} \leq C \|v_*\|_{C^m(\bar{\mathcal{D}})} \cdot N^{-(m-2-2\iota)}$$

for a sufficiently large  $N$  and  $s = N^\iota$ , where  $\iota > 0$  is an arbitrarily small constant. The coefficients are bounded by  $|c_{v_*, m, \alpha}| \leq C \|v_*\|_{C^m(\bar{\mathcal{D}})}$ .

### Step 3: Neural Network Approximation of Localized Polynomials.

With the help of Lemmas D.1 and D.2, the construction of sub-networks to approximate  $\Psi_m^s(\mathbf{x})\mathbf{x}^\alpha$  and its error analysis can be established directly following the methodology of Lemmas 4-6 in [28]. While the network architecture remains unchanged, the error analysis must be adapted to the  $W^{2,\infty}$  norm, which introduces a factor  $(sN)^2$  in the error bound attributable to the second-order derivative terms. This result is formalized in the following lemma:

**Lemma D.3.** *Let  $d, N, s \geq 1$ ,  $m \geq 2$ . For any  $0 < \epsilon_{NN} < \epsilon^*$ , where  $\epsilon^* > 0$  is sufficiently small, there exists a neural network  $\Phi_\omega^{m,\alpha} \in \mathcal{F}(H, D, R_\omega)$  with architecture parameters*

$$H = 2^{\lceil \log_2(d + \|\alpha\|_1) \rceil + 1}, \quad D = \lceil \log_2(d + \|\alpha\|_1) \rceil + 2, \quad R_\omega = \max\{3Ns, (3d + 3/2)s, C(d)\epsilon_{NN}^{-2}\},$$

such that for some constant  $C(m, d) > 0$ ,

$$\|\Psi_m^s \mathbf{x}^\alpha - \Phi_\omega^{m,\alpha}(\mathbf{x})\|_{W^{2,\infty}(\mathcal{D})} \leq C(m, d) (sN)^2 \epsilon_{NN},$$

for all  $m \in \{0, \dots, N\}^d$  and  $\alpha \in \mathbb{N}^d$  with  $\|\alpha\|_1 \leq m-1$ .

Following Theorem 6 of [28], we construct the final parallel network by combining these sub-networks with their corresponding coefficients  $c_{v_*, m, \alpha}$ , as shown below:

**Lemma D.4.** *For any  $0 < \epsilon_{NN} < \epsilon^*$ , where  $\epsilon^* > 0$  is sufficiently small, there exists a parallel neural network  $v_{\bar{q}, \bar{\omega}} \in \mathcal{PF}_{\bar{q}, \bar{K}}(\bar{H}, \bar{D}, R_{\bar{\omega}})$  with parameters*

$$\begin{aligned} \bar{q} &= C_1(m, d)(N+1)^d, \quad \bar{K} = C_2(m, d)(N+1)^d, \quad \bar{H} = 2^{\lceil \log_2(d+m-1) \rceil + 1}, \\ \bar{D} &= \lceil \log_2(d+m-1) \rceil + 2, \quad R_{\bar{\omega}} = \max\{3Ns, (3d+3/2)s, C(d)\epsilon_{NN}^{-2}\}, \end{aligned}$$

such that

$$\|v_N - v_{\bar{q}, \bar{\omega}}\|_{W^{2,\infty}(\mathcal{D})} \leq C(m, d)(N+1)^d (sN)^2 \epsilon_{NN}.$$

#### Step 4: Final Error Bound.

To achieve a total approximation error  $\|v_* - v_{\bar{q}, \bar{\omega}}\|_{C^2(\bar{\mathcal{D}})} \leq \epsilon$ , we first have

$$\|v_* - v_{\bar{q}, \bar{\omega}}\|_{C^2(\bar{\mathcal{D}})} \leq \|v_* - v_N\|_{W^{2,\infty}(\mathcal{D})} + \|v_N - v_{\bar{q}, \bar{\omega}}\|_{W^{2,\infty}(\mathcal{D})}.$$

Then, from Step 1, to ensure  $\|v_* - v_N\|_{W^{2,\infty}(\mathcal{D})} \leq \epsilon/2$ , we choose  $N$  such that  $N^{-(m-2-2\iota)} = \mathcal{O}(\epsilon)$ . This gives:

$$N = \mathcal{O}\left(\epsilon^{-\frac{1}{m-2-2\iota}}\right) \quad \text{and} \quad s = N^\iota = \mathcal{O}\left(\epsilon^{-\frac{\iota}{m-2-2\iota}}\right).$$

From Step 2, to ensure  $\|v_N - v_{\bar{q}, \bar{\omega}}\|_{W^{2,\infty}(\mathcal{D})} \leq \epsilon/2$ , we need to choose the sub-network accuracy  $\epsilon_{NN}$  such that  $C(m, d)(N+1)^d (sN)^2 \epsilon_{NN} \leq \epsilon/2$ , which yields:

$$\epsilon_{NN} = \mathcal{O}\left(\frac{\epsilon}{s^2 N^{d+2}}\right) = \mathcal{O}\left(\epsilon^{1 + \frac{d+2+2\iota}{m-2-2\iota}}\right).$$

These choices directly determine the complexity of the resulting parallel network  $v_{\bar{q}, \bar{\omega}}$ , as summarized in the following result which is exactly Theorem 3.2 in the main text:

**Theorem D.1.** *Given any  $v_* \in C^m(\bar{\mathcal{D}})$  with  $m \geq 3$ , for some small  $\epsilon^* > 0$  and any  $0 < \epsilon < \epsilon^*$ , there exists  $v_{\bar{q}, \bar{\omega}} \in \mathcal{PF}_{\bar{q}, \bar{K}}(\bar{H}, \bar{D}, R_{\bar{\omega}})$  with*

$$\begin{aligned} \bar{q} &= \lceil C_1 \cdot \epsilon^{-\frac{d}{m-2\iota-2}} \rceil, \quad \bar{K} = C_2 \cdot \epsilon^{-\frac{d}{m-2\iota-2}}, \quad \bar{H} = 2^{\lceil \log_2(d+m-1) \rceil + 1}, \\ \bar{D} &= \lceil \log_2(d+m-1) \rceil + 2, \quad R_{\bar{\omega}} = C_3 \cdot \epsilon^{-\frac{2d+2m}{m-2\iota-2}}, \end{aligned}$$

such that

$$\|v_{\bar{q}, \bar{\omega}} - v_*\|_{C^2(\bar{\mathcal{D}})} \leq \epsilon,$$

where  $0 < \iota < 1$ ; Constants  $C_1$ ,  $C_2$ , and  $C_3$  depend exclusively on  $d$  and  $m$ .

## References

- [1] Cosmin Anitescu, Elena Atroshchenko, Naif Alajlan, and Timon Rabczuk. Artificial neural network methods for the solution of second order boundary value problems. *Computers, Materials and Continua*, 59(1):345–359, 2019.
- [2] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- [3] Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin. Deep learning: A statistical viewpoint. *Acta Numerica*, 30:87–201, 2021.
- [4] Mikhail Belkin. Fit without fear: Remarkable mathematical phenomena of deep learning through the prism of interpolation. *Acta Numerica*, 30:203–248, 2021.
- [5] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.
- [6] Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to understand kernel learning. In *International Conference on Machine Learning*, pages 541–549. PMLR, 2018.
- [7] Mikhail Belkin, Alexander Rakhlin, and Alexandre B Tsybakov. Does data interpolation contradict statistical optimality? In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1611–1619. PMLR, 2019.
- [8] Julius Berner, Markus Dablander, and Philipp Grohs. Numerically solving parametric families of high-dimensional Kolmogorov partial differential equations via deep learning. In *Advances in Neural Information Processing Systems*, volume 33, pages 16615–16627. Curran Associates, Inc., 2020.
- [9] Susanne Brenner and Ridgway Scott. *The Mathematical Theory of Finite Element Methods*, volume 15. Springer Science & Business Media, 2007.
- [10] Jingrun Chen, Xurong Chi, Zhouwang Yang, et al. Bridging traditional and machine learning-based algorithms for solving PDEs: The random feature method. *Journal of Machine Learning*, 1:268–298, 2022.
- [11] Mo Chen, Zhao Ding, Yuling Jiao, Xiliang Lu, Peiyang Wu, and Jerry Zhijian Yang. Convergence analysis of PINNs with over-parameterization. *Communications in Computational Physics*, 37(4):942–974, 2025.
- [12] Philippe G Ciarlet. *The Finite Element Method for Elliptic Problems*. SIAM, 2002.
- [13] Yongcheng Dai, Bangti Jin, Ramesh Chandra Sau, and Zhi Zhou. Solving elliptic optimal control problems via neural networks and optimality system. *Advances in Computational Mathematics*, 51(4):31, 2025.
- [14] John Duchi, Shai Shalev-Shwartz, Yoram Singer, and Tushar Chandra. Efficient projections onto the  $l_1$ -ball for learning in high dimensions. In *Proceedings of the 25th international Conference on Machine Learning*, pages 272–279, 2008.
- [15] R.M Dudley. The sizes of compact subsets of hilbert space and continuity of Gaussian processes. *Journal of Functional Analysis*, 1(3):290–330, 1967.
- [16] Weinan E and Bing Yu. The deep Ritz method: A deep learning-based numerical algorithm for solving variational problems. *Communications in Mathematics and Statistics*, 6(1):1–12, 2018.
- [17] Fengjiang Fu and Xiaoqun Wang. Convergence analysis of a quasi-Monte Carlo-based deep learning algorithm for solving partial differential equations. *Numerical Mathematics: Theory, Methods and Applications*, 16(3):668–700, 2023.

- [18] Ingo Gühring and Mones Raslan. Approximation rates for neural networks with encodable weights in smoothness spaces. *Neural Networks*, 134:107–130, 2021.
- [19] Jiequn Han, Arnulf Jentzen, and E Weinan. Solving high-dimensional partial differential equations using deep learning. *Proceedings of the National Academy of Sciences*, 115(34):8505–8510, 2018.
- [20] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *Annals of Statistics*, 50(2):949, 2022.
- [21] Qingguo Hong, Jonathan W Siegel, and Jinchao Xu. Rademacher complexity and numerical quadrature analysis of stable neural networks with applications to numerical PDEs. *arXiv preprint arXiv:2104.02903*, 2021.
- [22] Tianhao Hu, Bangti Jin, and Zhi Zhou. Solving Poisson problems in polygonal domains with singularity enriched physics informed neural networks. *SIAM Journal on Scientific Computing*, 46(4):C369–C398, 2024.
- [23] Martin Hutzenthaler, Arnulf Jentzen, Thomas Kruse, Tuan Anh Nguyen, and Philippe von Wurstemberger. Overcoming the curse of dimensionality in the numerical approximation of semilinear parabolic partial differential equations. *Proceedings of the Royal Society A*, 476(2244):20190630, 2020.
- [24] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in Neural Information Processing Systems*, 31, 2018.
- [25] Ameya D Jagtap, Ehsan Kharazmi, and George Em Karniadakis. Conservative physics-informed neural networks on discrete domains for conservation laws: Applications to forward and inverse problems. *Computer Methods in Applied Mechanics and Engineering*, 365:113028, 2020.
- [26] Yuling Jiao, Yanming Lai, Dingwei Li, Xiliang Lu, Fengru Wang, Yang Wang, and Jerry Zhijian Yang. A rate of convergence of physics informed neural networks for the linear second order elliptic PDEs. *Communications in Computational Physics*, 31(4):1272–1295, 2022.
- [27] Yuling Jiao, Yanming Lai, Yisu Lo, Yang Wang, and Yunfei Yang. Error analysis of deep Ritz methods for elliptic equations. *Analysis and Applications*, 2023.
- [28] Yuling Jiao, Ruoxuan Li, Peiying Wu, Jerry Zhijian Yang, and Pingwen Zhang. DRM revisited: A complete error analysis. *Journal of Machine Learning Research*, 26(115):1–76, 2025.
- [29] Yuling Jiao, Xiliang Lu, Peiying Wu, and Jerry Zhijian Yang. Convergence analysis for over-parameterized deep learning. *Communications in Computational Physics*, 36(1):71–103, 2024.
- [30] Yuling Jiao, Yang Wang, and Yunfei Yang. Approximation bounds for norm constrained neural networks with applications to regression and gans. *Applied and Computational Harmonic Analysis*, 65:249–278, 2023.
- [31] Yuling Jiao, Jerry Zhijian Yang, Cheng Yuan, and Junyu Zhou. A rate of convergence of weak adversarial neural networks for the second order parabolic PDEs. *Communications in Computational Physics*, 34(3):813–836, 2023.
- [32] Gitta Kutyniok, Philipp Petersen, Mones Raslan, and Reinhold Schneider. A theoretical analysis of deep neural networks and parametric PDEs. *Constructive Approximation*, 55(1):73–125, 2022.
- [33] Tengyuan Liang and Alexander Rakhlin. Just interpolate: Kernel “ridgeless” regression can generalize. *The Annals of Statistics*, 48(3):1329–1347, 2020.
- [34] Zichao Long, Yiping Lu, Xianzhong Ma, and Bin Dong. PDE-Net: Learning PDEs from data. In *International Conference on Machine Learning*, pages 3208–3216. PMLR, 2018.
- [35] Lu Lu, Xuhui Meng, Zhiping Mao, and George Em Karniadakis. DeepXDE: A deep learning

- library for solving differential equations. *SIAM Review*, 63(1):208–228, 2021.
- [36] Yiping Lu, Haoxuan Chen, Jianfeng Lu, Lexing Ying, and Jose Blanchet. Machine learning for elliptic PDEs: Fast rate generalization bound, neural scaling law and minimax optimality. *ICLR*, 2021.
- [37] Colin McDiarmid. *On the method of bounded differences*, page 148–188. London Mathematical Society Lecture Note Series. Cambridge University Press, 1989.
- [38] Siddhartha Mishra and Roberto Molinaro. Estimates on the generalization error of physics-informed neural networks for approximating a class of inverse problems for PDEs. *IMA Journal of Numerical Analysis*, 42(2):981–1022, 2022.
- [39] Siddhartha Mishra and T Konstantin Rusch. Enhancing accuracy of deep learning algorithms by training with low-discrepancy sequences. *SIAM Journal on Numerical Analysis*, 59(3):1811–1834, 2021.
- [40] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- [41] Johannes Müller and Marius Zeinhofer. Error estimates for the variational training of neural networks with boundary penalty. *arXiv preprint arXiv:2103.01007*, 2021.
- [42] Preetum Nakkiran, Prayaag Venkat, Sham M Kakade, and Tengyu Ma. Optimal regularization can mitigate double descent. In *International Conference on Learning Representations*, 2020.
- [43] Guofei Pang, Marta D’Elia, Michael Parks, and George Karniadakis. nPINNs: Nonlocal physics-informed neural networks for a parametrized nonlocal universal Laplacian operator. Algorithms and applications. *Journal of Computational Physics*, 422:109760, 08 2020.
- [44] Guofei Pang, Lu Lu, and George Em Karniadakis. fPINNs: Fractional physics-informed neural networks. *SIAM Journal on Scientific Computing*, 41(4):A2603–A2626, 2019.
- [45] Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- [46] Yeonjong Shin. On the convergence of physics informed neural networks for linear second-order elliptic and parabolic type PDEs. *Communications in Computational Physics*, 28(5):2042–2074, 2020.
- [47] Justin A. Sirignano and K. Spiliopoulos. DGM: A deep learning algorithm for solving partial differential equations. *Journal of Computational Physics*, 375:1339–1364, 2018.
- [48] Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019.
- [49] Sifan Wang, Xinling Yu, and Paris Perdikaris. When and why PINNs fail to train: A neural tangent kernel perspective. *Journal of Computational Physics*, 449:110768, 2022.
- [50] Yahong Yang and Juncai He. Deeper or wider: A perspective from optimal generalization error with sobolev loss. In *Proceedings of the 41st International Conference on Machine Learning*, pages 56109–56138, 2024.
- [51] Yunfei Yang and Ding-Xuan Zhou. Optimal rates of approximation by shallow  $\text{relu}^k$  neural networks and applications to nonparametric regression. *Constructive Approximation*, pages 1–32, 2024.
- [52] Yaohua Zang, Gang Bao, Xiaojing Ye, and Haomin Zhou. Weak adversarial networks for high-dimensional partial differential equations. *Journal of Computational Physics*, 411:109409, 2020.