

GRADIENT-ENHANCED DEEP NEURAL NETWORK APPROXIMATIONS

Xiaodong Feng & Li Zeng*

LSEC, Institute of Computational Mathematics and Scientific/Engineering Computing, AMSS, Chinese Academy of Sciences, Beijing, China

*Address all correspondence to: Li Zeng, LSEC, Institute of Computational Mathematics and Scientific/Engineering Computing, AMSS, Chinese Academy of Sciences, Beijing, China, E-mail: zengli@lsec.cc.ac.cn

Original Manuscript Submitted: 11/8/2022; Final Draft Received: 12/10/2022

We propose in this work the gradient-enhanced deep neural network (DNN) approach for function approximations and uncertainty quantification. More precisely, the proposed approach adopts both the function evaluations and the associated gradient information to yield enhanced approximation accuracy. In particular, the gradient information is included as a regularization term in the gradient-enhanced DNN approach, for which we present posterior estimates (by the two-layer neural networks) similar to those in the path-norm regularized DNN approximations. We also discuss the application of this approach to gradient-enhanced uncertainty quantification, and present several numerical experiments to show that the proposed approach can outperform the traditional DNN approach in many cases of interest.

KEY WORDS: *deep neural networks, two-layer neural network, Barron space, uncertainty quantification*

1. INTRODUCTION

In recent years, deep neural networks (DNNs) have been widely used for dealing with scientific and engineering problems, such as function approximations (E et al., 2022; Schwab and Zech, 2019; Siegel and Xu, 2020), numerical partial differential equations (PDEs) (E and Yu, 2018; Raissi et al., 2019; Sirignano and Spiliopoulos, 2018), image classification (He et al., 2016; Litjens et al., 2017), and uncertainty quantification (Meng and Karniadakis, 2020; Qin et al., 2021; Yang et al., 2021), to name a few. Compared to traditional tools such as polynomials (DeVore and Lorentz, 1993), radial basis functions (Majdisova and Skala, 2017), and kernel methods (Liu et al., 2020), one of the main advantages of DNNs is their potential approximation capacity for high-dimensional problems. Unlike classic tools such as polynomial approximations (for which the relevant theoretical analysis results have been well studied), the associated theoretical analysis for DNNs is still in its infancy. Among others, we mention the seminal works by Barron (1993) and E et al. (2022) where the concept of “Barron space” was proposed, and some approximation results for DNNs were presented.

In this work, we shall present the gradient-enhanced DNN approach. More precisely, our approach adopts both the function evaluations and the associated gradient information. This is similar to the classic Hermite-type interpolation. Our main contributions are summarized as follows:

- We present the gradient-enhanced DNN approach, where the gradient information is included as a regularization term.
- For the gradient-enhanced DNN approach, we present posterior estimates (via a two-layer neural network) similar to those in the path-norm regularized DNN approximations. More precisely, we show that the posterior generalization error can be bounded by

$$\mathcal{O}\left(d\|\boldsymbol{\theta}\|_{\mathcal{P}}\sqrt{\frac{\ln(2d)}{n}} + \|\boldsymbol{\theta}\|_{\mathcal{P}}\frac{\ln(\|\boldsymbol{\theta}\|_{\mathcal{P}})}{\sqrt{n}}\right),$$

where n is the number of training points, $\|\boldsymbol{\theta}\|_{\mathcal{P}}$ is the path norm, and d is the dimension.

- We discuss the application of our approach to gradient-enhanced uncertainty quantification and present several numerical experiments to show that the gradient-enhanced DNN approach can outperform the traditional DNNs in many cases of interest.

We remark that gradient-enhanced polynomial approximations have been proposed for uncertainty quantification (Guo et al., 2018; Jakeman et al., 2015; Li et al., 2011; Lockwood and Mavriplis, 2013; Peng et al., 2016). Under some circumstances of uncertainty quantification, the cost of calculating derivatives can be inexpensive, e.g., by solving adjoint equations. These additional derivatives actually increase the existing data information, which may bring a better approximation. On the other hand, there are already many gradient-enhanced methods based on deep learning. Many researchers consider enhancing the gradient value back-propagated by the network to avoid the phenomenon of gradient vanishing or gradient explosion during training (He et al., 2016; Schmidhuber and Hochreiter, 1997; Yan et al., 2022). A gradient-enhanced damage model was used in Zhuang et al. (2020) to ensure the well-posedness of the boundary value and yields mesh-independent results in computational methods. Some also regularize the derivatives of the output of neural networks to improve the adversarial robustness and interpretability (Ross and Doshi-Velez, 2018). In addition, a gradient-enhanced physics-informed neural network (PINN) was proposed to improve the accuracy and training efficiency of PINNs in Yu et al. (2022), where the gradient information (i.e., the associated adjoint equation) is included to yield a modified PINN loss function. Different from these methods, our method is data driven, tackling the case that derivatives are given along with function evaluations. Moreover, theoretical analysis is provided to illustrate the feasibility and effectiveness of this method.

The rest of this paper is organized as follows. In Section 2, we set up the problem and present some preliminaries. In Section 3 we present the error estimations in two-layer neural networks for the gradient regularized approximation problem. Applications to uncertainty quantification are discussed in Section 4. Finally, we give some concluding remarks in Section 5.

2. PRELIMINARIES

2.1 Problem Setup

We begin the discussion by considering function approximations via DNNs with labeled data. We consider the target function $f^* : \mathbb{R}^d \rightarrow \mathbb{R}$, and we assume that the following data are available: $\{\mathbf{x}_i, y_i, \mathbf{y}'_i\}_{i=1}^n$. Here y_i and \mathbf{y}'_i are the functional evaluations and gradient evaluations, respectively. Namely,

$$\begin{cases} y_i = f^*(\mathbf{x}_i), & i = 1, \dots, n. \\ \mathbf{y}'_i = \nabla f^*(\mathbf{x}_i), & i = 1, \dots, n. \end{cases} \quad (1)$$

For simplicity, we assume that the data $\{\mathbf{x}_i\}_i$ lie in $\mathcal{X} = [-1, 1]^d$ and $0 \leq f^* \leq 1$. We shall show our analysis results via a two-layer neural network, for which the nonlinear function can be defined as

$$f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{k=1}^m a_k \sigma(\mathbf{w}_k^T \mathbf{x}), \quad (2)$$

where $\mathbf{w}_k \in \mathbb{R}^d$, $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is a nonlinear activation function, and $\boldsymbol{\theta} = \{(a_k, \mathbf{w}_k)\}_{k=1}^m$ is the unknown parameter. We define a truncated form of $f(\mathbf{x}; \boldsymbol{\theta})$ through

$$Tf(\mathbf{x}; \boldsymbol{\theta}) = \max\{\min\{f(\mathbf{x}; \boldsymbol{\theta}), 1\}, 0\}.$$

By an abuse of notation, in the following we still use $f(\mathbf{x}; \boldsymbol{\theta})$ to denote $Tf(\mathbf{x}; \boldsymbol{\theta})$. For the training $\{\mathbf{x}_i, y_i\}_{i=1}^n$, the population risk can be defined by

$$L(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, y} [\ell(f(\mathbf{x}; \boldsymbol{\theta}), y)]. \quad (3)$$

The empirical risk with the training data yields

$$L_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i), \quad (4)$$

where $\ell(f(\mathbf{x}), y) = (1/2)(f(\mathbf{x}) - y)^2$. For the gradient information, we consider

$$L'(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, y} [\tilde{\ell}(\nabla f(\mathbf{x}; \boldsymbol{\theta}), \mathbf{y}')], \quad (5)$$

where $\tilde{\ell}(\nabla f(\mathbf{x}; \boldsymbol{\theta}), \mathbf{y}') = \|\nabla f(\mathbf{x}; \boldsymbol{\theta}) - \mathbf{y}'\|_2$ and $\|\cdot\|_q$ indicates the ℓ_q norm of a vector. Similarly, the empirical risk with the training data $\{\mathbf{x}_i, y_i, \mathbf{y}'_i\}_{i=1}^n$ yields

$$L'_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n [\tilde{\ell}(\nabla f(\mathbf{x}_i; \boldsymbol{\theta}), \mathbf{y}'_i)]^2. \quad (6)$$

Moreover, the path-norm of the two-layer neural network is defined as

$$\|\boldsymbol{\theta}\|_{\mathcal{P}} := \sum_{k=1}^m |a_k| \|\mathbf{w}_k\|_1. \quad (7)$$

We are now ready to present the so-called gradient-enhanced DNN approach.

Definition 2.1 (Gradient-enhanced DNNs model). For a two-layer neural network $f(\cdot; \boldsymbol{\theta})$ of width m , the gradient regularized risk is defined as follows:

$$J_{n, \beta}(\boldsymbol{\theta}) := L_n(\boldsymbol{\theta}) + \beta \cdot L'_n(\boldsymbol{\theta}).$$

The corresponding regularized estimator is defined as

$$\boldsymbol{\theta}_{n, \beta} = \arg \min J_{n, \beta}(\boldsymbol{\theta}).$$

Note that the minimizers are not necessarily unique, and $\boldsymbol{\theta}_{n, \beta}$ should be understood as any of the minimizers.

The above approach can be viewed as an extension of the classic Hermite interpolation (Spitzbart, 1960; Wu, 1992), and is motivated by applications such as gradient-enhanced uncertainty quantification (Guo et al., 2018; Jakeman et al., 2015; Li et al., 2011; Lockwood and Mavriplis, 2013; Peng et al., 2016).

2.2 Barron Space

We now provide a brief overview of Barron space (E et al., 2019). Let $\mathbb{S}^{d-1} := \{\mathbf{w} \in \mathbb{R}^d \mid \|\mathbf{w}\|_1 = 1\}$. Let \mathcal{F} be the Borel σ -algebra on \mathbb{S}^{d-1} and let $\mathbb{P}(\mathbb{S}^{d-1})$ be the collection of the probability measures on $(\mathbb{S}^{d-1}, \mathcal{F})$. Let $\mathcal{B}(\mathcal{X})$ be the collection of functions that admit the following integral representation:

$$f(\mathbf{x}) = \int_{\mathbb{S}^{d-1}} a(\mathbf{w}) \sigma(\langle \mathbf{w}, \mathbf{x} \rangle) d\pi(\mathbf{w}), \quad \forall \mathbf{x} \in \mathcal{X}, \quad (8)$$

where $\pi \in \mathbb{P}(\mathbb{S}^{d-1})$, and $a(\cdot)$ is a measurable function with respect to $(\mathbb{S}^d, \mathcal{F})$. For any $f \in \mathcal{B}(\mathcal{X})$ and $p \geq 1$, we define the following norm:

$$\gamma_p(f) := \inf_{(a, \pi) \in \Theta_f} \left(\int_{\mathbb{S}^{d-1}} |a(\mathbf{w})|^p d\pi(\mathbf{w}) \right)^{1/p}, \quad (9)$$

where

$$\Theta_f = \left\{ (a, \pi) \mid f(\mathbf{x}) = \int_{\mathbb{S}^{d-1}} a(\mathbf{w}) \sigma(\langle \mathbf{w}, \mathbf{x} \rangle) d\pi(\mathbf{w}) \right\}.$$

Definition 2.2 [Barron space (E et al., 2019)]. The Barron space is defined as

$$\mathcal{B}_p(\mathcal{X}) := \left\{ f \in \mathcal{B}(\mathcal{X}) \mid \gamma_p(f) < \infty \right\}.$$

We next present several assumptions.

Assumption 2.1. *Throughout the paper we assume that*

- $\mathcal{X} = [-1, 1]^d$ and $0 \leq f^* \leq 1$.
- The derivative of $f^*(\mathbf{x})$ is bounded by a constant D .
- The activation function σ is scaling invariant, namely, $\sigma(kz) = k\sigma(z)$, and satisfies $|\sigma(z)| \leq C_1|z|$, $|\sigma'| \leq C_2$ and σ' is Lipschitz continuous with a positive constant C_3 . Specifically, we use rectified linear unit (ReLU), i.e., $\sigma(z) = \max\{z, 0\}$ as activation function with constant $C_1 = C_2 = C_3 = 1$.
- $\ln(2d) \geq 1$; here d is the dimension of input data.

Proposition 2.1 [Proposition of E et al. (2022)]. Let $f \in C(\mathcal{X})$, the space of continuous functions on \mathcal{X} , and assume that f satisfies

$$\gamma(f) := \inf_{\hat{f}} \int_{\mathbb{R}^d} \|\mathbf{w}\|_1^2 |\hat{f}(\mathbf{w})| d\mathbf{w} \leq \infty, \quad (10)$$

where \hat{f} is the Fourier transform of an extension of f to \mathbb{R}^d . Then f admits an integral representation. Moreover,

$$\gamma_p(f) \leq 2\gamma(f) + 2\|\nabla f(0)\|_1 + 2|f(0)|. \quad (11)$$

It is usually difficult to check condition (10). For further understanding, we recall the following relationship of Barron space and more classical function spaces, which can be found in E and Stephan (2021) and references therein for details.

Proposition 2.2 [Theorem 3.1 of E and Stephan (2021)]. If $f \in H^s(\mathbb{R}^d)$ for $s > d/2 + 2$, then $f \in \mathcal{B}(\mathbb{R}^d)$.

Proposition 2.2 states that every sufficient smooth function admits an integral representation.

3. ERROR ESTIMATES FOR THE GRADIENT-ENHANCED DNN APPROACH

In this section, we present the error estimates for the gradient-enhanced DNN approach. To this aim, we first present the following approximation theorem, which combines the relationship between Barron spaces and two-layer neural networks.

Theorem 3.1. *For any $f \in \mathcal{B}_2(\mathcal{X})$, there exists a two-layer neural network $f(\cdot; \tilde{\theta})$ of width m , such that*

$$\mathbb{E}_{\mathbf{x}} \left[(f(\mathbf{x}) - f(\mathbf{x}; \tilde{\theta}))^2 \right] \leq \frac{3\gamma_2^2(f)}{m}, \quad (12)$$

$$\mathbb{E}_{\mathbf{x}} \left[\left\| \nabla f(\mathbf{x}) - \nabla f(\mathbf{x}; \tilde{\theta}) \right\|_2^2 \right] \leq \frac{7\gamma_2^2(f)}{m}, \quad (13)$$

$$\|\tilde{\theta}\|_{\mathcal{D}} \leq 2\gamma_2(f). \quad (14)$$

Proof. Let (a, π) be the best representation of f , i.e., $\gamma_2^2(f) = \mathbb{E}_{\pi}[|a(\mathbf{w})|^2]$, and let $U = \{\mathbf{w}_j\}_{j=1}^m$ be i.i.d. random variables drawn from $\pi(\cdot)$. Define

$$\hat{f}_U(\mathbf{x}) := \frac{1}{m} \sum_{j=1}^m a(\mathbf{w}_j) \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle).$$

Then the derivative of f with respect to \mathbf{x} is

$$\nabla \hat{f}_U(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m a(\mathbf{w}_j) \sigma'(\langle \mathbf{w}_j, \mathbf{x} \rangle) \mathbf{w}_j^T.$$

Let $L_U^1 = \mathbb{E}_{\mathbf{x}} \left[|\hat{f}_U(\mathbf{x}) - f(\mathbf{x})|^2 \right]$, $L_U^2 = \mathbb{E}_{\mathbf{x}} \left[\left\| \nabla \hat{f}_U(\mathbf{x}) - \nabla f(\mathbf{x}) \right\|_2^2 \right]$. Then we have

$$\begin{aligned} \mathbb{E}_U [L_U^1] &= \mathbb{E}_{\mathbf{x}} \mathbb{E}_U \left[|\hat{f}_U(\mathbf{x}) - f(\mathbf{x})|^2 \right] \\ &= \mathbb{E}_{\mathbf{x}} \mathbb{E}_U \left[\left| \frac{1}{m} \sum_{j=1}^m a(\mathbf{w}_j) \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle) - f(\mathbf{x}) \right|^2 \right] \\ &= \frac{1}{m^2} \mathbb{E}_{\mathbf{x}} \mathbb{E}_U \left[\left(\sum_{j=1}^m (a(\mathbf{w}_j) \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle) - f(\mathbf{x})) \right) \left(\sum_{i=1}^m (a(\mathbf{w}_i) \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) - f(\mathbf{x})) \right) \right] \\ &= \frac{1}{m^2} \mathbb{E}_{\mathbf{x}} \sum_{i,j=1}^m \mathbb{E}_{\mathbf{w}_i, \mathbf{w}_j} \left[(a(\mathbf{w}_j) \sigma(\langle \mathbf{w}_j, \mathbf{x} \rangle) - f(\mathbf{x})) (a(\mathbf{w}_i) \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) - f(\mathbf{x})) \right] \\ &= \frac{1}{m} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{w}} \left[(a(\mathbf{w}) \sigma(\langle \mathbf{w}, \mathbf{x} \rangle) - f(\mathbf{x}))^2 \right] \\ &\leq \frac{1}{m} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{w}} \left[(a(\mathbf{w}) \sigma(\langle \mathbf{w}, \mathbf{x} \rangle))^2 \right] \\ &\leq \frac{1}{m} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{w}} \left[|a(\mathbf{w})|^2 |\sigma(\langle \mathbf{w}, \mathbf{x} \rangle)|^2 \right] \\ &\leq \frac{1}{m} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{w}} \left[|a(\mathbf{w})|^2 \right] = \frac{1}{m} \mathbb{E}_{\mathbf{w}} \left[|a(\mathbf{w})|^2 \right]. \end{aligned}$$

Since $\|\mathbf{w}_j\|_1 = 1$ and $\mathbf{x} \in [-1, 1]^d$, we have $\langle \mathbf{w}_j, \mathbf{x} \rangle \leq \langle \mathbf{w}_j, \mathbf{x}_0 \rangle \leq \|\mathbf{w}_j\|_1 = 1$ where $\mathbf{x}_0 = (\mathbf{x}_{01}, \mathbf{x}_{02}, \dots, \mathbf{x}_{0d})$ and $\mathbf{x}_{0i} = \text{sgn}(\mathbf{w}_{ji})$, which implies the last inequality. Meanwhile,

$$\begin{aligned}
\mathbb{E}_U[L_U^2] &= \mathbb{E}_{\mathbf{x}} \mathbb{E}_U \left[\|\nabla \hat{f}_U(\mathbf{x}) - \nabla f(\mathbf{x})\|_2^2 \right] \\
&= \mathbb{E}_{\mathbf{x}} \mathbb{E}_U \left[\left\| \frac{1}{m} \sum_{j=1}^m a(\mathbf{w}_j) \sigma'(\langle \mathbf{w}_j, \mathbf{x} \rangle) \mathbf{w}_j^T - \nabla f(\mathbf{x}) \right\|_2^2 \right] \\
&= \frac{1}{m^2} \mathbb{E}_{\mathbf{x}} \mathbb{E}_U \left[\left(\sum_{j=1}^m (a(\mathbf{w}_j) \sigma'(\langle \mathbf{w}_j, \mathbf{x} \rangle) \mathbf{w}_j^T - \nabla f(\mathbf{x})) \right) \right. \\
&\quad \times \left. \left(\sum_{i=1}^m (a(\mathbf{w}_i) \sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle) \mathbf{w}_i^T - \nabla f(\mathbf{x})) \right) \right] \\
&= \frac{1}{m^2} \mathbb{E}_{\mathbf{x}} \sum_{i,j=1}^m \mathbb{E}_{\mathbf{w}_i, \mathbf{w}_j} \left[(a(\mathbf{w}_j) \sigma'(\langle \mathbf{w}_j, \mathbf{x} \rangle) \mathbf{w}_j^T - \nabla f(\mathbf{x}))^T \right. \\
&\quad \times \left. (a(\mathbf{w}_i) \sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle) \mathbf{w}_i^T - \nabla f(\mathbf{x})) \right] \\
&= \frac{1}{m} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{w}} \left[\|a(\mathbf{w}) \sigma'(\langle \mathbf{w}, \mathbf{x} \rangle) \mathbf{w}^T - \nabla f(\mathbf{x})\|_2^2 \right] \\
&\leq \frac{1}{m} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{w}} \left[\|a(\mathbf{w}) \sigma'(\langle \mathbf{w}, \mathbf{x} \rangle) \mathbf{w}^T\|_2^2 \right] \\
&\leq \frac{1}{m} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{w}} \left[|a(\mathbf{w})|^2 \|\mathbf{w}\|_2^2 \right] \\
&\leq \frac{1}{m} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{w}} \left[|a(\mathbf{w})|^2 \|\mathbf{w}\|_1^2 \right] = \frac{1}{m} \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{w}} \left[|a(\mathbf{w})|^2 \right] = \frac{1}{m} \mathbb{E}_{\mathbf{w}} \left[|a(\mathbf{w})|^2 \right].
\end{aligned}$$

Denote the path-norm of $\hat{f}_U(\mathbf{x})$ by A_U ; we have $\mathbb{E}[A_U] = \gamma_1(f) \leq \gamma_2(f)$.

Define events $E_1 = \left\{ L_U^1 < \frac{3\gamma_2^2(f)}{m} \right\}$, $E_2 = \left\{ L_U^2 \leq \frac{7\gamma_2^2(f)}{m} \right\}$, $E_3 = \{A_U < 2\gamma_2(f)\}$.

Using Markov's inequality, we have

$$\begin{aligned}
\mathbb{P}(E_1) &= 1 - \mathbb{P} \left(\left\{ L_U^1 \geq \frac{3\gamma_2^2(f)}{m} \right\} \right) \geq 1 - \frac{\mathbb{E}_U[L_U^1]}{3\gamma_2^2(f)/m} \geq 1 - \frac{\gamma_2(f)^2/m}{3\gamma_2^2(f)/m} = \frac{2}{3}, \\
\mathbb{P}(E_2) &= 1 - \mathbb{P} \left(\left\{ L_U^2 \geq \frac{7\gamma_2^2(f)}{m} \right\} \right) \geq 1 - \frac{\mathbb{E}_U[L_U^2]}{7\gamma_2^2(f)/m} \geq 1 - \frac{\gamma_2(f)^2/m}{7\gamma_2^2(f)/m} = \frac{6}{7}, \\
\mathbb{P}(E_3) &= 1 - \mathbb{P}(\{A_U \geq 2\gamma_2(f)\}) \geq 1 - \frac{\mathbb{E}_U[A_U]}{2\gamma_2(f)} \geq 1 - \frac{\gamma_2(f)}{2\gamma_2(f)} = \frac{1}{2}.
\end{aligned}$$

Therefore, let $E_4 = E_1 \cap E_3$, then

$$\mathbb{P}(E_1 \cap E_3) \geq \mathbb{P}(E_1) + \mathbb{P}(E_3) - 1 = \frac{1}{6}.$$

Hence,

$$\mathbb{P}(E_1 \cap E_2 \cap E_3) = \mathbb{P}(E_4 \cap E_2) \geq \mathbb{P}(E_4) + \mathbb{P}(E_2) - 1 \geq \frac{1}{6} + \frac{6}{7} - 1 = \frac{1}{42} > 0.$$

That is, $E_1 \cap E_2 \cap E_3 \neq \emptyset$, which completes the proof. \square

According to Theorem 3.1, for a function in Barron space, we can use a two-layer neural network to simultaneously approximate its function value and derivative. Next we will investigate the posterior error estimation of the gradient-enhanced method.

3.1 A Posteriori Error Estimation

To give a *a posteriori* error estimation, we need to introduce the definition of the Rademacher complexity and review some related conclusions.

Definition 3.1 [Rademacher complexity (Shalev-Shwartz and Ben-David, 2014)]. Let $S = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ be n i.i.d samples, and let $\mathcal{F} \circ S$ be the set of all possible evaluations a function $f \in \mathcal{F}$ can achieve on a sample S , namely,

$$\mathcal{F} \circ S = \left\{ (f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n)) \mid f \in \mathcal{F} \right\}.$$

Let each component of random variable ξ be i.i.d. according to $\mathbb{P}[\xi_i = 1] = \mathbb{P}[\xi_i = -1] = 1/2$. Then, the Rademacher complexity of \mathcal{F} with respect to S is defined as follows:

$$\mathcal{R}_n(\mathcal{F} \circ S) := \frac{1}{n} \mathbb{E}_{\xi \sim \{\pm 1\}^n} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^n \xi_i f(\mathbf{x}_i) \right]. \quad (15)$$

Lemma 3.1 [Lemma 26.11 of Shalev-Shwartz and Ben-David (2014)]. Let $S = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ be vectors in \mathbb{R}^d , $\mathcal{H}_1 = \{\mathbf{x} \mapsto \langle \mathbf{w}, \mathbf{x} \rangle \mid \|\mathbf{w}\|_1 \leq 1\}$. Then,

$$\mathcal{R}_n(\mathcal{H}_1 \circ S) \leq \max_i \|\mathbf{x}_i\|_\infty \sqrt{\frac{2 \ln(2d)}{n}}.$$

Lemma 3.2 [Lemma 26.9 of Shalev-Shwartz and Ben-David (2014)]. For each $i \in [n] = \{1, 2, \dots, n\}$, let $\phi_i : \mathbb{R} \mapsto \mathbb{R}$ be a ρ -Lipschitz continuous, namely, for all $\alpha, \beta \in \mathbb{R}$ we have $|\phi_i(\alpha) - \phi_i(\beta)| \leq \rho |\alpha - \beta|$. For $\mathbf{a} \in \mathbb{R}^n$, let $\Phi(\mathbf{a})$ denote the vector $(\phi_1(a_1), \dots, \phi_n(a_n))$. Let $\Phi \circ A = \{\Phi(\mathbf{a}) : \mathbf{a} \in A\}$. Then,

$$\mathcal{R}_n(\Phi \circ A) \leq \rho \mathcal{R}_n(A).$$

Lemma 3.3 [Lemma 26.6 of Shalev-Shwartz and Ben-David (2014)]. *For any $A \subset \mathbb{R}^n$, scale $c \in \mathbb{R}$, and vector $\mathbf{a}_0 \in \mathbb{R}^n$, we have*

$$\mathcal{R}_n(\{c\mathbf{a} + \mathbf{a}_0 \mid \mathbf{a} \in A\}) \leq |c| \mathcal{R}_n(A).$$

Lemma 3.4. *For any $A, B \subset \mathbb{R}^n$, we have*

$$\mathcal{R}_n(A + B) \leq \mathcal{R}_n(A) + \mathcal{R}_n(B).$$

Here $A + B = \{a + b \mid a \in A, b \in B\}$.

Lemma 3.5 [Lemma B.3 of E et al. (2019)]. *Let $\mathcal{F}_Q = \{f(\mathbf{x}; \boldsymbol{\theta}) \mid \|\boldsymbol{\theta}\|_{\mathcal{D}} \leq Q\}$ be the set of two-layer networks with path-norm bounded by Q ; then we have*

$$\mathcal{R}_n(\mathcal{F}_Q \circ S) \leq 2Q \sqrt{\frac{2 \ln(2d)}{n}}.$$

Lemma 3.6 [Vector valued Rademacher complexity (Maurer, 2016)]. *Let \mathcal{X} be any set, $(\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathcal{X}^n$, let \mathcal{F} be a class of functions $\mathbf{f} : \mathcal{X} \rightarrow \ell_2$, and let $h_i : \ell_2 \rightarrow \mathbb{R}$ have Lipschitz norm Lip , where ℓ_2 denote the Hilbert space of square summable sequences of real numbers. Then*

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^n \xi_i h_i(\mathbf{f}(\mathbf{x}_i)) \right] \leq \sqrt{2} \text{Lip} \cdot \mathbb{E} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^n \sum_{k=1}^d \xi_{ik} f_k(\mathbf{x}_i) \right], \quad (16)$$

where ξ_{ik} is an independent doubly indexed Rademacher sequence according to probability distribution $\mathbb{P}[\xi_{ik} = -1] = \mathbb{P}[\xi_{ik} = 1] = 1/2$ and $f_k(\mathbf{x}_i)$ is the k th component of $\mathbf{f}(\mathbf{x}_i)$.

Lemma 3.7. *Let $S = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$. $\{\mathbf{y}'_i\}_{i=1}^n$ denotes the gradient information, where $\mathbf{y}'_i \in \mathbb{R}^d \forall i \in [n]$.*

$$\mathcal{F}'_{Q,j} = \{\nabla_j f(\mathbf{x}; \boldsymbol{\theta}) \mid \|\boldsymbol{\theta}\|_{\mathcal{D}} \leq Q\}, \quad j = 1, \dots, d, \quad \mathcal{F}'_Q = \Pi_{j=1}^d \mathcal{F}'_{Q,j}.$$

Define $\mathbf{g} : \mathcal{X} \rightarrow \ell_2$, $\mathbf{g}(\mathbf{x}; \boldsymbol{\theta}) = (\partial_1 f(\mathbf{x}; \boldsymbol{\theta}), \partial_2 f(\mathbf{x}; \boldsymbol{\theta}), \dots, \partial_d f(\mathbf{x}; \boldsymbol{\theta}))^T$ and $\tilde{\ell}_j : \ell_2 \rightarrow \mathbb{R}$, $\tilde{\ell}_j(\mathbf{z}) = \|\mathbf{z} - \mathbf{y}'_j\|_2$ for each $j \in \{1, 2, \dots, n\}$, where ℓ_2 denote the Hilbert space of square summable sequences of real numbers. Note that $\tilde{\ell}_j$ is 1-Lipschitz function; then we have

$$\mathbb{E} \left[\sup_{\|\boldsymbol{\theta}\|_{\mathcal{D}} \leq Q} \sum_{i=1}^n \xi_i \tilde{\ell}_i(\mathbf{g}(\mathbf{x}_i; \boldsymbol{\theta})) \right] \leq \sqrt{2} \cdot \mathbb{E} \left[\sup_{\|\boldsymbol{\theta}\|_{\mathcal{D}} \leq Q} \sum_{i=1}^n \sum_{k=1}^d \xi_{ik} g_k(\mathbf{x}_i; \boldsymbol{\theta}) \right],$$

where $g_k(\mathbf{x}_i; \boldsymbol{\theta})$ is the k th component of $\mathbf{g}(\mathbf{x}_i; \boldsymbol{\theta})$. Moreover, denote $\tilde{\ell} = (\tilde{\ell}_1, \tilde{\ell}_2, \dots, \tilde{\ell}_n)$; we have

$$\mathcal{R}_n(\tilde{\ell} \circ \mathcal{F}'_Q \circ S) \leq 2\sqrt{2} Q d \sqrt{\frac{2 \ln(2d)}{n}}.$$

Proof. Without loss of generality, let $\|\mathbf{w}_j\|_1 = 1, \forall j = \{1, \dots, m\}$. Applying Lemma 3.6, we immediately obtain

$$\begin{aligned}
\mathcal{R}_n(\tilde{\ell} \circ \mathcal{F}'_Q \circ S) &= \frac{1}{n} \mathbb{E}_{\xi} \left[\sup_{\|\boldsymbol{\theta}\|_{\mathcal{S}} \leq Q} \sum_{i=1}^n \xi_i \|g(\mathbf{x}_i; \boldsymbol{\theta}) - \mathbf{y}'_i\|_2 \right] \\
&\leq \frac{\sqrt{2}}{n} \mathbb{E}_{\xi} \left[\sup_{\|\boldsymbol{\theta}\|_{\mathcal{S}} \leq Q} \sum_{k=1}^d \sum_{i=1}^n \xi_{ik} g_k(\mathbf{x}_i; \boldsymbol{\theta}) \right] \\
&= \frac{\sqrt{2}}{n} \mathbb{E}_{\xi} \left[\sup_{\|\boldsymbol{\theta}\|_{\mathcal{S}} \leq Q} \sum_{k=1}^d \sum_{i=1}^n \xi_{ik} \sum_{j=1}^m a_j \sigma'(\mathbf{w}_j^T \mathbf{x}_i) \mathbf{w}_{jk} \right] \\
&= \frac{\sqrt{2}}{n} \mathbb{E}_{\xi} \left[\sup_{\|\boldsymbol{\theta}\|_{\mathcal{S}} \leq Q} \sum_{j=1}^m a_j \|\mathbf{w}_j\|_1 \sum_{k=1}^d \sum_{i=1}^n \xi_{ik} \sigma'(\mathbf{w}_j^T \mathbf{x}_i) \mathbf{w}_{jk} \right] \\
&\leq \frac{\sqrt{2}}{n} \mathbb{E}_{\xi} \left[\sup_{\|\boldsymbol{\theta}\|_{\mathcal{S}} \leq Q} \sum_{j=1}^m |a_j| \|\mathbf{w}_j\|_1 \sup_{\|\mathbf{v}\|_1=1} \left| \sum_{k=1}^d \sum_{i=1}^n \xi_{ik} \sigma'(\mathbf{v}^T \mathbf{x}_i) \mathbf{v}_k \right| \right] \\
&\leq \frac{\sqrt{2}Q}{n} \mathbb{E}_{\xi} \left[\sup_{\|\mathbf{v}\|_1=1} \left| \sum_{k=1}^d \sum_{i=1}^n \xi_{ik} \sigma'(\mathbf{v}^T \mathbf{x}_i) \mathbf{v}_k \right| \right] \\
&= \frac{\sqrt{2}Q}{n} \mathbb{E}_{\boldsymbol{\eta}} \left[\sup_{\|\mathbf{v}\|_1=1} \left| \left\langle \sum_{i=1}^n \boldsymbol{\eta}_i \sigma'(\mathbf{v}^T \mathbf{x}_i), \mathbf{v} \right\rangle \right| \right] \quad (\boldsymbol{\eta}_i = (\xi_{i1}, \xi_{i2}, \dots, \xi_{id})) \\
&\leq \frac{\sqrt{2}Q}{n} \mathbb{E}_{\boldsymbol{\eta}} \left[\sup_{\|\mathbf{v}\|_1=1} \left\| \sum_{i=1}^n \boldsymbol{\eta}_i \sigma'(\mathbf{v}^T \mathbf{x}_i) \right\|_1 \right] \\
&\leq \frac{\sqrt{2}Qd}{n} \mathbb{E}_{\xi} \left[\sup_{\|\mathbf{v}\|_1 \leq 1} \sum_{i=1}^n \xi_i \sigma'(\mathbf{v}^T \mathbf{x}_i) \right],
\end{aligned}$$

where $\xi \in [-1, 1]$ and $\boldsymbol{\eta} \in [-1, 1]^d$. Due to the symmetry, we have

$$\begin{aligned}
\mathbb{E}_{\xi} \left[\sup_{\|\mathbf{u}\|_1 \leq 1} \left| \sum_{i=1}^n \xi_i \sigma'(\mathbf{u}^T \mathbf{x}_i) \right| \right] &\leq \mathbb{E}_{\xi} \left[\sup_{\|\mathbf{u}\|_1 \leq 1} \sum_{i=1}^n \xi_i \sigma'(\mathbf{u}^T \mathbf{x}_i) + \sup_{\|\mathbf{u}\|_1 \leq 1} \sum_{i=1}^n -\xi_i \sigma'(\mathbf{u}^T \mathbf{x}_i) \right] \\
&\leq 2 \mathbb{E}_{\xi} \left[\sup_{\|\mathbf{u}\|_1 \leq 1} \sum_{i=1}^n \xi_i \sigma'(\mathbf{u}^T \mathbf{x}_i) \right].
\end{aligned}$$

Then applying Lemma 3.2 and Lemma 3.1, we obtain

$$\mathcal{R}_n(\tilde{\ell} \circ \mathcal{F}'_Q \circ S) \leq 2\sqrt{2}Qd \sqrt{\frac{2 \ln(2d)}{n}}.$$

□

Theorem 3.2 (Shalev-Shwartz and Ben-David, 2014). *Assume that for all samples \mathbf{x} and h in hypothesis space \mathcal{H} , the loss function $\ell : \mathcal{H} \times \mathcal{Y} \rightarrow \mathbb{R}$ satisfies $|\ell(h(\mathbf{x}), y)| \leq B$. Then, with probability of at least $1 - \delta$, for all $h \in \mathcal{H}$, and $S = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$,*

$$\left| \frac{1}{n} \sum_{i=1}^n \ell(h(\mathbf{x}_i), y_i) - \mathbb{E}_{\mathbf{x}, y} [\ell(h(\mathbf{x}), y)] \right| \leq 2\mathbb{E}_S [\mathcal{R}_n(\ell \circ \mathcal{H} \circ S)] + B\sqrt{\frac{2\ln(2/\delta)}{n}}. \quad (17)$$

We are now ready to present the main results of this section.

Theorem 3.3. *If Assumption 2.1 holds, then with probability at least $1 - \delta$ we have,*

$$\begin{aligned} & \sup_{\|\boldsymbol{\theta}\|_{\mathcal{Q}} \leq Q} \left| L(\boldsymbol{\theta}) + \beta L'(\boldsymbol{\theta}) - (L_n(\boldsymbol{\theta}) + \beta L'_n(\boldsymbol{\theta})) \right| \\ & \leq 4(1 + \sqrt{2}\beta d)Q\sqrt{\frac{2\ln(2d)}{n}} + \left(\frac{1}{2} + \beta(Q + D) \right) \sqrt{\frac{2\ln(2/\delta)}{n}}. \end{aligned} \quad (18)$$

Proof. Using the triangular inequality, we have

$$\begin{aligned} & \sup_{\|\boldsymbol{\theta}\|_{\mathcal{Q}} \leq Q} \left| L(\boldsymbol{\theta}) + \beta L'(\boldsymbol{\theta}) - (L_n(\boldsymbol{\theta}) + \beta L'_n(\boldsymbol{\theta})) \right| \\ & \leq \sup_{\|\boldsymbol{\theta}\|_{\mathcal{Q}} \leq Q} |L(\boldsymbol{\theta}) - L_n(\boldsymbol{\theta})| + \sup_{\|\boldsymbol{\theta}\|_{\mathcal{Q}} \leq Q} \beta |L'(\boldsymbol{\theta}) - L'_n(\boldsymbol{\theta})|. \end{aligned} \quad (19)$$

Define $\mathcal{H}_Q = \{ \ell(f(\mathbf{x}; \boldsymbol{\theta}), y) \mid f(\mathbf{x}; \boldsymbol{\theta}) \in \mathcal{F}_Q \}$ and $\mathcal{H}'_Q = \{ \tilde{\ell}(\nabla f(\mathbf{x}; \boldsymbol{\theta}), \mathbf{y}') \mid f(\mathbf{x}; \boldsymbol{\theta}) \in \mathcal{F}_Q \}$, where $\ell(f(\mathbf{x}; \boldsymbol{\theta}), y) = (1/2)(f(\mathbf{x}; \boldsymbol{\theta}) - y)^2$ and $\tilde{\ell}(\nabla f(\mathbf{x}; \boldsymbol{\theta}), \mathbf{y}') = \|\nabla f(\mathbf{x}; \boldsymbol{\theta}) - \mathbf{y}'\|_2$.

Note that $f(\mathbf{x}; \boldsymbol{\theta}) \in [0, 1]$, $\ell(\cdot, y)$ is 1-Lipschitz continuous,

$$\begin{aligned} & \frac{1}{2} \left| (f(\mathbf{x}_i; \boldsymbol{\theta}_f) - y_i)^2 - (g(\mathbf{x}_i; \boldsymbol{\theta}_g) - y_i)^2 \right| \\ & = \frac{1}{2} |f(\mathbf{x}_i; \boldsymbol{\theta}_f) - g(\mathbf{x}_i; \boldsymbol{\theta}_g)| \cdot |f(\mathbf{x}_i; \boldsymbol{\theta}_f) + g(\mathbf{x}_i; \boldsymbol{\theta}_g) - 2y_i| \leq |f(\mathbf{x}_i; \boldsymbol{\theta}_f) - g(\mathbf{x}_i; \boldsymbol{\theta}_g)|. \end{aligned}$$

Following Lemma 3.5 and Lemma 3.2, we have

$$\mathcal{R}_n(\mathcal{H}_Q) = \frac{1}{n} \mathbb{E}_{\mathbf{x}} \left[\sup_{f \in \mathcal{F}_Q} \sum_{i=1}^n \xi_i \frac{1}{2} |f(\mathbf{x}_i; \boldsymbol{\theta}) - y_i|^2 \right] \leq 2Q\sqrt{\frac{2\ln(2d)}{n}}, \quad (20)$$

and ℓ is bounded, $\ell(f(\mathbf{x}; \boldsymbol{\theta}), y) = (1/2)(f(\mathbf{x}; \boldsymbol{\theta}) - y)^2 \leq 1/2$. Hence applying Theorem 3.2, we can obtain,

$$\sup_{\|\boldsymbol{\theta}\|_{\mathcal{Q}} \leq Q} |L(\boldsymbol{\theta}) - L_n(\boldsymbol{\theta})| \leq 4Q\sqrt{\frac{2\ln(2d)}{n}} + \frac{1}{2}\sqrt{\frac{2\ln(2/\delta)}{n}}. \quad (21)$$

Similarly, $\tilde{\ell}(\nabla f(\mathbf{x}; \boldsymbol{\theta}), \mathbf{y}') = \|\nabla f(\mathbf{x}; \boldsymbol{\theta}) - \mathbf{y}'\|_2$ is 1-Lipschitz with respect to $\nabla f(\mathbf{x}; \boldsymbol{\theta})$,

$$\left| \|\nabla f(\mathbf{x}_i; \boldsymbol{\theta}_f) - \mathbf{y}'_i\|_2 - \|\nabla g(\mathbf{x}_i; \boldsymbol{\theta}_g) - \mathbf{y}'_i\|_2 \right| \leq \|\nabla f(\mathbf{x}_i; \boldsymbol{\theta}_f) - \nabla g(\mathbf{x}_i; \boldsymbol{\theta}_g)\|_2.$$

Let $\tilde{\ell} = (\tilde{\ell}_1, \dots, \tilde{\ell}_n)$, in which $\tilde{\ell}_i(\cdot) = \|\cdot - \mathbf{y}'_i\|_2$ for $i = 1, \dots, n$. We can obtain the estimation of $\mathcal{R}_n(\mathcal{H}'_Q)$ directly from Lemma 3.7,

$$\mathcal{R}_n(\mathcal{H}'_Q) = \frac{1}{n} \mathbb{E}_{\mathbf{x}} \left[\sup_{f \in \mathcal{F}_Q} \sum_{i=1}^n \xi_i \|\nabla f(\mathbf{x}_i; \boldsymbol{\theta}) - \mathbf{y}'_i\|_2 \right] \leq 2\sqrt{2} Q d \sqrt{\frac{2\ln(2d)}{n}}, \quad (22)$$

and $\tilde{\ell}(\nabla f(\mathbf{x}; \boldsymbol{\theta}), \mathbf{y}')$ is bounded since we assume that the gradient of the objective function is bounded by a positive constant D :

$$f(\mathbf{x}; \boldsymbol{\theta}) = \sum_{k=1}^m a_k \sigma(\mathbf{w}_k^T \mathbf{x}),$$

$$\|\nabla f(\mathbf{x}; \boldsymbol{\theta})\|_2^2 = \left\| \sum_{k=1}^m a_k \sigma' \mathbf{w}_k^T \right\|_2^2 \leq \left\| \sum_{k=1}^m a_k \sigma' \mathbf{w}_k^T \right\|_1^2 \leq \left\| \sum_{k=1}^m a_k \mathbf{w}_k^T \right\|_1^2 \leq \|\boldsymbol{\theta}_f\|_{\mathcal{D}}^2, \quad (23)$$

$$\tilde{\ell}(\nabla f(\mathbf{x}; \boldsymbol{\theta}), \mathbf{y}') = \|\nabla f(\mathbf{x}; \boldsymbol{\theta}) - \mathbf{y}'\|_2 \leq \|\nabla f(\mathbf{x}; \boldsymbol{\theta})\|_2 + \|\mathbf{y}'\| \leq \|\boldsymbol{\theta}_f\|_{\mathcal{D}} + D \leq Q + D. \quad (24)$$

Thus,

$$\sup_{\|\boldsymbol{\theta}\|_{\mathcal{D}} \leq Q} |L'(\boldsymbol{\theta}) - L'_n(\boldsymbol{\theta})| \leq 4\sqrt{2} Q d \sqrt{\frac{2 \ln(2d)}{n}} + (Q + D) \sqrt{\frac{2 \ln(2/\delta)}{n}}. \quad (25)$$

The desired result follows by plugging Eqs. (21) and (25) into Eq. (19), and the proof is completed. \square

We next present a posterior generalization bound by relaxing such restrictions.

Theorem 3.4 (A posterior generalization bound). *Assume that Assumption 2.1 holds, then for any $\delta > 0$, with probability at least $1 - \delta$ over the choice of the training set S , we have, for any two-layer network $f_m(\cdot, \boldsymbol{\theta})$,*

$$\begin{aligned} \left| L(\boldsymbol{\theta}) + \beta L'(\boldsymbol{\theta}) - (L_n(\boldsymbol{\theta}) + \beta L'_n(\boldsymbol{\theta})) \right| &\leq 4(1 + \sqrt{2}\beta d) (\|\boldsymbol{\theta}\|_{\mathcal{D}} + 1) \sqrt{\frac{2 \ln(2d)}{n}} \\ &\quad + \left(\frac{1}{2} + \beta (\|\boldsymbol{\theta}\|_{\mathcal{D}} + 1 + D) \right) \sqrt{\frac{2 \ln(2c(\|\boldsymbol{\theta}\|_{\mathcal{D}} + 1)^2/\delta)}{n}}, \end{aligned} \quad (26)$$

where $c = \sum_{k=1}^{\infty} 1/k^2$.

Proof. Consider the decomposition $\mathcal{F} = \cup_{k=1}^{\infty} \mathcal{F}_k$, where $\mathcal{F}_k = \{f(\mathbf{x}; \boldsymbol{\theta}) \mid \|\boldsymbol{\theta}\|_{\mathcal{D}} \leq k\}$. Let $\delta_k = \delta/(ck^2)$ where $c = \sum_{k=1}^{\infty} 1/k^2$. According to Theorem 3.3, if we fix k in advance, then with probability at least $1 - \delta_k$ over the choice of S , we have

$$\begin{aligned} &\left| L(\boldsymbol{\theta}) + \beta L'(\boldsymbol{\theta}) - (L_n(\boldsymbol{\theta}) + \beta L'_n(\boldsymbol{\theta})) \right| \\ &\leq 4(1 + \sqrt{2}\beta d) k \sqrt{\frac{2 \ln(2d)}{n}} + \left(\frac{1}{2} + \beta(k + D) \right) \sqrt{\frac{2 \ln(2/\delta_k)}{n}}. \end{aligned} \quad (27)$$

Therefore $\mathbb{P}(\{\text{inequality (27) unholds}\}) \leq \sum_{k=1}^{\infty} \delta_k =: \delta$, namely, $\mathbb{P}(\{\text{inequality (27) holds for all } k\}) \geq 1 - \delta$. In other words, with probability at least $1 - \delta$, the inequality (27) holds for all k .

Given an arbitrary set of parameters $\boldsymbol{\theta}$, denote $k_0 = \min\{k \mid \|\boldsymbol{\theta}\|_{\mathcal{D}} \leq k\}$, then $k_0 \leq \|\boldsymbol{\theta}\|_{\mathcal{D}} + 1$. Inequality (27) implies that

$$\begin{aligned} & \left| L(\boldsymbol{\theta}) + \beta L'(\boldsymbol{\theta}) - (L_n(\boldsymbol{\theta}) + \beta L'_n(\boldsymbol{\theta})) \right| \\ & \leq 4(1 + \sqrt{2}\beta d)k_0 \sqrt{\frac{2\ln(2d)}{n}} + \left(\frac{1}{2} + \beta(k_0 + D) \right) \sqrt{\frac{2\ln(2ck_0^2/\delta)}{n}}. \\ & \leq 4(1 + \sqrt{2}\beta d)(\|\boldsymbol{\theta}\|_{\mathcal{D}} + 1) \sqrt{\frac{2\ln(2d)}{n}} \\ & \quad + \left(\frac{1}{2} + \beta(\|\boldsymbol{\theta}\|_{\mathcal{D}} + 1 + D) \right) \sqrt{\frac{2\ln(2c(\|\boldsymbol{\theta}\|_{\mathcal{D}} + 1)^2/\delta)}{n}}. \end{aligned}$$

This completes the proof. \square

We note that the generalization gap is bounded roughly by

$$d\|\boldsymbol{\theta}\|_{\mathcal{D}} \sqrt{\frac{\ln(2d)}{n}} + \|\boldsymbol{\theta}\|_{\mathcal{D}} \frac{\ln(\|\boldsymbol{\theta}\|_{\mathcal{D}})}{\sqrt{n}},$$

which shares a similar convergence rate with the method without gradient information (E et al., 2019), indicating that the gradient-enhanced method would not destroy the original function approximation algorithm.

3.2 An Upper Bound for the Empirical Risk

Recalling the approximation property, there exists a two-layer neural network $f(\cdot; \tilde{\boldsymbol{\theta}})$ whose path-norm is independent of the network width, while achieving the optimal approximation error. Furthermore, this path-norm can also be used to bound the generalization gap (Theorem 3.4). In this section, we want to estimate the gradient regularized risk of $\tilde{\boldsymbol{\theta}}$. To this end, we first assume that the norm of the gradient of the target function can be bounded by a constant D .

Let $\hat{\gamma}_p(f) = \max\{1, \gamma_p(f)\}$, where d is the dimension of input data.

Theorem 3.5. *Let $\tilde{\boldsymbol{\theta}}$ be the network mentioned in Theorem 3.1; then with probability at least $1 - \delta$, we have*

$$J_{n,\beta}(\tilde{\boldsymbol{\theta}}) \lesssim \frac{\gamma_2^2(f^*)}{m} + \beta \frac{1}{\sqrt{m}} \gamma_2(f^*) + \beta \hat{\gamma}_2(f^*) (\hat{\gamma}_2(f^*) + d) \sqrt{\frac{2\ln(2d)}{n}} + \beta \hat{\gamma}_2(f^*) \sqrt{\frac{\ln(2c/\delta)}{n}}, \quad (28)$$

where

$$J_{n,\beta}(\tilde{\boldsymbol{\theta}}) = L_n(\tilde{\boldsymbol{\theta}}) + \beta L'_n(\tilde{\boldsymbol{\theta}}) \quad \text{and} \quad c = \sum_{k=1}^{\infty} \frac{1}{k^2}.$$

Proof. According to the definition of a regularized model, the properties that

$$\|\tilde{\boldsymbol{\theta}}\|_{\mathcal{D}} \leq 2\gamma_2(f^*), \quad L(\tilde{\boldsymbol{\theta}}) \leq \frac{3\gamma_2^2(f^*)}{m}, \quad (L'(\tilde{\boldsymbol{\theta}}))^2 \leq \mathbb{E}_x [\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}; \tilde{\boldsymbol{\theta}})\|_2^2] \leq \frac{7\gamma_2^2(f^*)}{m},$$

and the posteriori error bound, the regularized risk of $\tilde{\Theta}$ satisfies

$$\begin{aligned}
 J_{n,\beta}(\tilde{\Theta}) &= L_n(\tilde{\Theta}) + \beta L'_n(\tilde{\Theta}) \\
 &\leq L(\tilde{\Theta}) + \beta L'(\tilde{\Theta}) + 4(1 + \sqrt{2}\beta d)(\|\tilde{\Theta}\|_{\mathcal{D}} + 1)\sqrt{\frac{2\ln(2d)}{n}} \\
 &\quad + \left(\frac{1}{2} + \beta(\|\tilde{\Theta}\|_{\mathcal{D}} + 1 + D)\right)\sqrt{\frac{2\ln(2c(\|\tilde{\Theta}\|_{\mathcal{D}} + 1)^2/\delta)}{n}} \\
 &\leq L(\tilde{\Theta}) + \beta L'(\tilde{\Theta}) + 4(1 + \sqrt{2}\beta d)(2\gamma_2(f^*) + 1)\sqrt{\frac{2\ln(2d)}{n}} \tag{29}
 \end{aligned}$$

$$+ \left(\frac{1}{2} + \beta(2\gamma_2(f^*) + 1 + D)\right)\sqrt{\frac{2\ln(2c(1 + 2\gamma_2(f^*))^2/\delta)}{n}}. \tag{30}$$

The last term can be simplified by using $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ and $\ln(a) \leq a$ for $a \geq 0, b \geq 0$. Thus we have

$$\begin{aligned}
 \sqrt{2\ln(2c(1 + 2\gamma_2(f^*))^2/\delta)} &= \sqrt{2\ln(2c/\delta) + 2\ln(1 + 2\gamma_2(f^*))^2} \\
 &\leq \sqrt{2\ln(2c/\delta)} + \sqrt{2\ln(1 + 2\gamma_2(f^*))^2} \\
 &\leq \sqrt{2\ln(2c/\delta)} + \sqrt{2\ln(3\hat{\gamma}_2(f^*))^2} \\
 &\leq \sqrt{2\ln(2c/\delta)} + 3\sqrt{2}\hat{\gamma}_2(f^*).
 \end{aligned}$$

Plugging it into Eq. (30), we obtain

$$\begin{aligned}
 J_{n,\beta}(\tilde{\Theta}) &\leq L(\tilde{\Theta}) + \beta L'(\tilde{\Theta}) + 4(1 + \sqrt{2}\beta d)(2\gamma_2(f^*) + 1)\sqrt{\frac{2\ln(2d)}{n}} \\
 &\quad + \left(\frac{1}{2} + \beta(2\gamma_2(f^*) + 1 + D)\right)\left(\sqrt{\frac{2\ln(2c/\delta)}{n}} + 3\sqrt{\frac{2}{n}}\hat{\gamma}_2(f^*)\right).
 \end{aligned}$$

Thus after some simplifications, we have

$$J_{n,\beta}(\tilde{\Theta}) \lesssim \frac{\gamma_2^2(f^*)}{m} + \beta \frac{1}{\sqrt{m}}\gamma_2(f^*) + \beta \hat{\gamma}_2(f^*)(\gamma_2(f^*) + d)\sqrt{\frac{2\ln(2d)}{n}} + \beta \hat{\gamma}_2(f^*)\sqrt{\frac{\ln(2c/\delta)}{n}}.$$

This completes the proof. \square

According to the definition of $\Theta_{n,\beta}$, we have $J_{n,\beta}(\Theta_{n,\beta}) \leq J_{n,\beta}(\tilde{\Theta})$. Thus the above theorem gives an upper bound for $J_{n,\beta}(\Theta_{n,\beta})$. Such an upper bound for the empirical risk verifies the feasibility of our method that we can obtain a good approximation by increasing the number of neurons in the hidden layer and samples.

4. APPLICATIONS TO UNCERTAINTY QUANTIFICATION

We now consider the application of the gradient-enhanced DNN approach for uncertainty quantification.

4.1 Gradient-Enhanced Uncertainty Quantification

In complex engineering systems, mathematical models can only serve as simplified and reduced representations of true physics, and the effect of some uncertainties, such as boundary/initial conditions and parameter values, can be significant. Uncertainty quantification (UQ) aims to develop numerical tools that can accurately predict quantities of interest (QoI) and facilitate the quantitative validation of the simulation model. Generally, we use differential equations to model complex systems on a domain Ω , in which the uncertainty sources are represented by Ξ . The solution u is governed by the PDEs

$$\begin{aligned}\mathcal{L}(\mathbf{x}, \Xi; u(\mathbf{x}, \Xi)) &= 0, & \mathbf{x} \in \Omega, \\ \mathcal{B}(\mathbf{x}, \Xi; u(\mathbf{x}, \Xi)) &= 0, & \mathbf{x} \in \partial\Omega,\end{aligned}\tag{31}$$

where \mathcal{L} and \mathcal{B} are differential and boundary operators, respectively. Our goal is to approximate the QoI $u(\mathbf{x}_0, \Xi)$ for some fixed spatial location \mathbf{x}_0 . To reduce the notation, we simply write $u(\Xi)$. In many applications, the dimension of random variable Ξ is very high and can be characterized by a d -dimensional random variable. Hence DNNs are good candidates for such problems.

We consider inclusion of gradient measurements with respect to random variables Ξ , i.e., $\partial u / \partial \Xi_k, k = 1, 2, \dots, d$. The gradient measurements can usually be obtained in a relatively inexpensive way via the adjoint method (Luchini and Bottaro, 2014).

4.2 Numerical Examples

We compare the performance between original neural networks without gradient input and gradient-enhanced neural networks. For simplicity, we say that a method is $X\%$ gradient-enhanced if $X\%$ of samples contain derivative information with respect to all dimensions. In our tests, each neural network contains two layers with 1000 hidden neurons. The hyperparameter β , which is used to balance the two part losses introduced by the function values and derivative information, is set to 10. We initialize all trainable parameters using the Glorot normal scheme. For the training procedure, we use the Adam optimizer. To quantitatively evaluate the accuracy of the numerical solution, we shall consider the relative L^2 error $\|u_\theta - u\|_2 / \|u\|_2$, where u and u_θ denote the ground truth and predicted solution. All numerical tests are implemented in PyTorch.

4.2.1 Function Approximations

Before applying the gradient-enhanced method to uncertainty quantification, we first demonstrate the effectiveness of our approach in approximating high-dimensional functions. More precisely, we consider the Gaussian function,

$$f_1(\mathbf{x}) = \exp\left(-\sum_{i=1}^d x_i^2\right), \quad \mathbf{x} = (x_1, x_2, \dots, x_d) \in [-1, 1]^d,$$

and the polynomial function,

$$f_2(\mathbf{x}) = \sum_{i=1}^{d/2} x_i x_{i+1}, \quad \mathbf{x} = (x_1, x_2, \dots, x_d) \in [-1, 1]^d.$$

For these two test functions, we assume that samples $\{x_i\}_{i=1}^n$ are uniformly distributed in $[-1, 1]^d$, y_i is the observation of target function at x_i , and y'_i is the corresponding derivative. Thus $\{x_i, y_i\}_{i=1}^n$ compose the training data for the original DNN. $\{x_i, y_i, y'_i\}_{i=1}^n$ compose the training data for 100% gradient-enhanced DNN. And $\{x_i, y_i\}_{i=1}^n \cup \{\hat{x}_j, \hat{y}'_j\}_{j=1}^m$ compose the training data for 20% gradient-enhanced DNN where m is the rounding-off of 20% n and $\{\hat{x}_j\}_{j=1}^m$ is randomly chosen from $\{x_i\}_{i=1}^n$. The learning rate for the Adam optimizer is set to 0.01 with 20% decay each 500 steps.

For Gaussian function $f_1(x)$, we consider the cases that $d = 2, 4$, and 8. The relative L^2 errors against the number of samples n are presented in top row of Fig. 1. The use of gradient information can indeed improve the accuracy, and furthermore, the more gradient information is included, the better accuracy is obtained. Moreover, we investigate the loss functions of different models for $d = 2$ with 400 samples, $d = 4$ with 1600 samples, and $d = 8$ with 3200 samples, which are depicted in the bottom row of Fig. 1.

For the polynomial function $f_2(x)$, we set the dimension to 4, 8, and 16. Similar to the Gaussian function, we present the relative L^2 errors for different dimensions in Fig. 2, which again shows that the gradient information regularized term can greatly enhance the approximation accuracy. The loss functions for $d = 4$ with 400 samples and $d = 8$ and 16 with 3200 samples are provided in the bottom row of Fig. 2. It can be observed that the loss function of gradient-enhanced methods may be smaller than the original DNN as the iteration number increases, verifying the strength of the gradient-enhanced methods.

4.2.2 Elliptic PDE with Random Inputs

We now consider the following stochastic elliptic PDE problem,

$$\begin{cases} -\nabla \cdot (a(x, \omega) \nabla u(x, \omega)) = f(x, \omega) & \text{in } \mathcal{D} \times \Omega, \\ u(x, \omega) = 0 & \text{on } \partial \mathcal{D} \times \Omega, \end{cases} \quad (32)$$

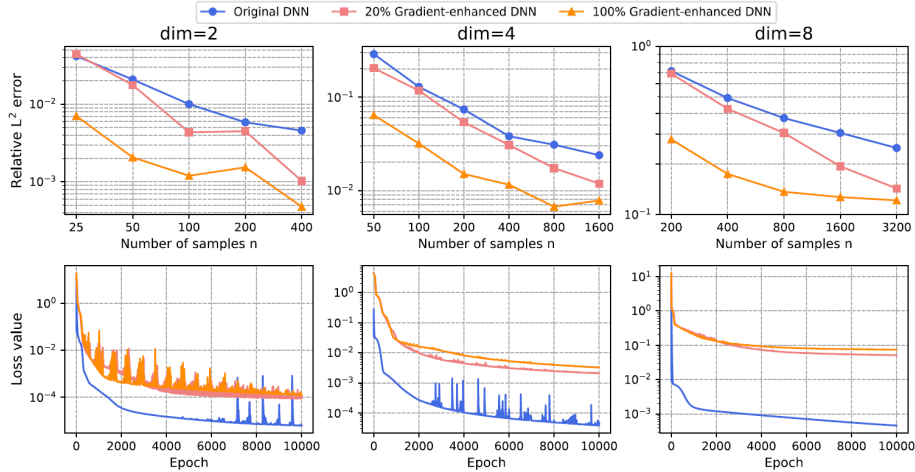


FIG. 1: Example 1. Approximation of $f_1(x)$. Top: The relative L^2 errors against the number of samples. Bottom: The loss functions against increasing epochs with the number of samples 400 for $d = 2$, 1600 for $d = 4$, and 3200 for $d = 8$.

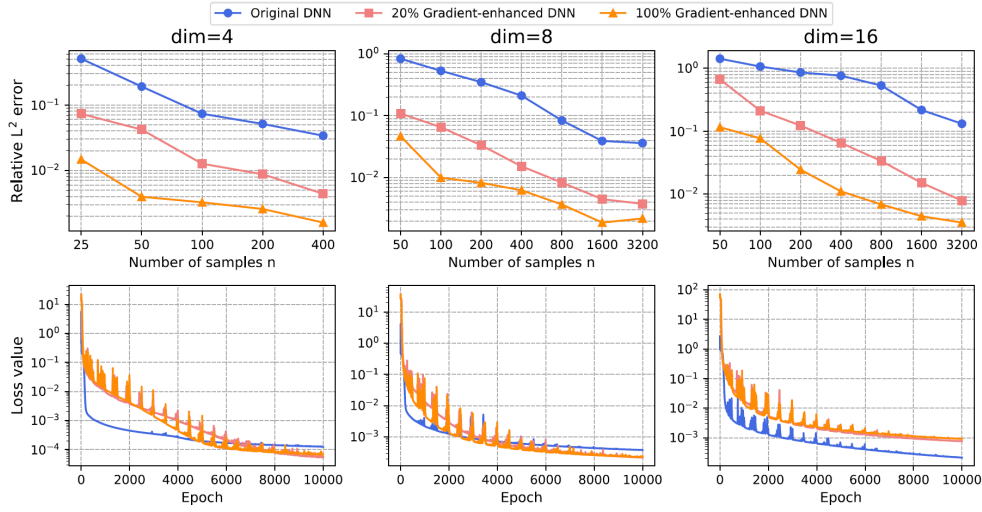


FIG. 2: Example 2. Approximation of $f_2(\mathbf{x})$. Top: The relative L^2 errors against the number of samples. Bottom: The loss functions against increasing epochs with the number of samples 400 for $\text{dim} = 4$ and 3200 for $\text{dim} = 8$ and 16.

where $\mathcal{D} = [0, 1]^2$, $\mathbf{x} = (x_1, x_2)$ is a spatial coordinate and $f(\mathbf{x}, \omega)$ is a deterministic force term $f(\mathbf{x}, \omega) = \cos(x_1) \sin(x_2)$. The random diffusion coefficient $a(\mathbf{x}, \omega) = a_d(\mathbf{x}, \omega)$ with one-dimensional spatial dependence takes the form (Babuka et al., 2010),

$$\log(a_d(\mathbf{x}, \omega) - 0.5) = 1 + Y_1(\omega) \left(\frac{\sqrt{\pi}L}{2} \right)^{1/2} + \sum_{k=2}^d \zeta_k \phi_k(\mathbf{x}) Y_k(\omega), \quad (33)$$

where

$$\zeta_k := (\sqrt{\pi}L)^{1/2} \exp\left(\frac{-(\lfloor \frac{k}{2} \rfloor \pi L)^2}{8}\right) \quad \text{if } k > 1 \text{ and } L = \frac{1}{12},$$

and $\phi_k(\mathbf{x})$ only depends on the first component of \mathbf{x} ,

$$\phi_k(\mathbf{x}) := \begin{cases} \sin(\lfloor \frac{k}{2} \rfloor \pi x_1) & \text{if } k \text{ even,} \\ \cos(\lfloor \frac{k}{2} \rfloor \pi x_1) & \text{if } k \text{ odd.} \end{cases} \quad (34)$$

Here $\{Y_k(\omega)\}_{k=1}^d$ are independent random variables uniformly distributed in the interval $[-1, 1]$. In the following we approximate the QoI q defined by $q(\omega) = u((0.5, 0.5), \omega)$, which is the solution of Eq. (32) at location $\mathbf{x} = (0.5, 0.5)$. Denote $\Psi(\omega) = (Y_1(\omega), \dots, Y_d(\omega))$. The derivatives $dq/d\Psi = \partial u(\mathbf{x}, \omega)/\partial \Psi$ are computed by the adjoint sensitivity method. Both forward and adjoint solvers are implemented in the finite element method (FEM) project FEniCS (Logg et al., 2012). In numerical tests, $\{\Psi(\omega_i)\}_{i=1}^n$ are generated from a uniform distribution in $[-1, 1]^d$, and we solve the forward PDE 320 times for $d = 5$ and 1600 times for $d = 10$. Notice that each partial derivative leads to an adjoint equation; then the number of adjoint equations needed to solve is d times that of the forward equations. It is worth mentioning that the cost of generating derivatives of $q(\omega)$ in elliptic PDE (32) is negligible since they share the same stiff matrix with q .

After obtaining the function values as well as the corresponding gradient information, we apply the original DNN, 20% gradient-enhanced DNN, and 100% gradient-enhanced DNN to approximate the QoI $q(\omega)$. The learning rate for the Adam optimizer is 0.001 with half decay each 1000 steps. The relative L^2 errors for $d = 5$ and 10 are presented at the top row of Fig. 3. We also provide the loss functions of different models for $d = 5$ with the number of samples 320 and $d = 10$ with the number of samples 1600 in the bottom row of Fig. 3. All cases verify that gradient-enhanced methods significantly outperform the original approach. We can achieve the same accuracy using much fewer training samples.

5. CONCLUSION

We have proposed gradient-enhanced deep neural network (DNN) approximations for function approximations and uncertainty quantification. In our approach, the gradient information is included as a regularization term. For this approach, we present posterior estimates (by the two-layer neural networks) similar to those in the path-norm regularized DNN approximations. We also discuss the application of this approach to gradient-enhanced uncertainty quantification, and numerical experiments show that the proposed approach can outperform the traditional DNN approach in many cases of interest. The discussion in this work is limited to supervised learning where labeled data are available, and in our future work, we will consider applying this gradient-enhanced idea to unsupervised learning where the physical equation is considered to yield the loss function.

ACKNOWLEDGMENT

We would like to thank Professor Tao Zhou of the Chinese Academy of Sciences for bringing this topic to our attention and for his encouragement and helpful discussion.

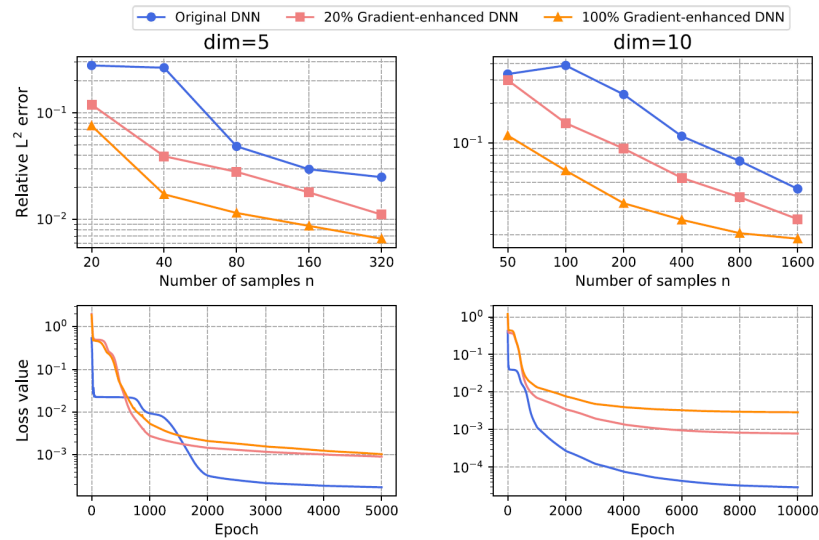


FIG. 3: Example 3. Top: The relative L^2 errors against number of samples for $N = 5, 10$. Bottom: The loss functions against increasing epochs for $d = 5$ with 320 samples and $d = 10$ with 1600 samples.

REFERENCES

- Barron, A., Universal Approximation Bounds for Superpositions of a Sigmoidal Function, *IEEE Trans. Inf. Theor.*, vol. **39**, no. 3, pp. 930–945, 1993.
- Babuka, I., Nobile, F., and Tempone, R., A Stochastic Collocation Method for Elliptic Partial Differential Equations with Random Input Data, *SIAM J. Numer. Anal.*, vol. **52**, no. 2, pp. 317–355, 2010.
- DeVore, R.A. and Lorentz, G.G., *Constructive Approximation*, vol. **303**, Berlin: Springer Science & Business Media, 1993.
- E, W., Ma, C., and Wu, L., *A Priori* Estimates of the Population Risk for Two-Layer Neural Networks, *Commun. Math. Sci.*, vol. **17**, no. 5, pp. 1407–1425, 2019.
- E, W., Ma, C., and Wu, L., The Barron Space and the Flow-Induced Function Spaces for Neural Network Models, *Constr. Approx.*, vol. **55**, no. 1, pp. 369–406, 2022.
- E, W. and Stephan, W., Representation Formulas and Pointwise Properties for Barron Functions, arXiv: 2006.05982, 2020.
- E, W. and Yu, B., The Deep Ritz Method: A Deep Learning-Based Numerical Algorithm for Solving Variational Problems, *Commun. Math. Stat.*, vol. **1**, no. 6, pp. 1–12, 2018.
- Guo, L., Narayan, A., and Zhou, T., A Gradient Enhanced ℓ_1 -Minimization for Sparse Approximation of Polynomial Chaos Expansions, *J. Comput. Phys.*, vol. **367**, pp. 49–64, 2018.
- He, K., Zhang, X., Ren, S., and Sun, J., Deep Residual Learning for Image Recognition, in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, Las Vegas, NV, pp. 770–778, 2016.
- Jakeman, J.D., Eldred, M.S., and Sargsyan, K., Enhancing ℓ_1 -Minimization Estimates of Polynomial Chaos Expansions Using Basis Selection, *J. Comput. Phys.*, vol. **289**, pp. 18–34, 2015.
- Li, Y., Anitescu, M., Roderick, O., and Hickernell, F., Orthogonal Bases for Polynomial Regression with Derivative Information in Uncertainty Quantification, *Vis. Mech. Processes: Int. Online J.*, vol. **1**, no. 4, pp. 297–320, 2011.
- Litjens, G., Kooi, T., Bejnordi, B.E., Setio, A.A.A., Ciompi, F., Ghafoorian, M., Van Der Laak, J.A., Van Ginneken, B., and Sánchez, C.I., A Survey on Deep Learning in Medical Image Analysis, *Med. Image Anal.*, vol. **42**, pp. 60–88, 2017.
- Liu, F., Huang, X., Chen, Y., and Suykens, J.A., Random Features for Kernel Approximation: A Survey on Algorithms, Theory, and Beyond, arXiv: 2004.11154, 2020.
- Lockwood, B. and Mavriplis, D., Gradient-Based Methods for Uncertainty Quantification in Hypersonic Flows, *Comput. Fluids*, vol. **85**, pp. 27–38, 2013.
- Logg, A., Mardal, K.A., and Wells, G., *Automated Solution of Differential Equations by the Finite Element Method: The FEniCS Book*, Berlin: Springer Science & Business Media, vol. **84**, 2012.
- Luchini, P. and Bottaro, A., Adjoint Equations in Stability Analysis, *Annu. Rev. Fluid Mech.*, vol. **46**, pp. 493–517, 2014.
- Majdisova, Z. and Skala, V., Radial Basis Function Approximations: Comparison and Applications, *Appl. Math. Model.*, vol. **51**, pp. 728–774, 2017.
- Maurer, A., A Vector-Contraction Inequality for Rademacher Complexities, in *Int. Conf. on Algorithmic Learning Theory*, Bari, Italy, pp. 3–17, 2016.
- Meng, X. and Karniadakis, G.E., A Composite Neural Network That Learns from Multi-Fidelity Data: Application to Function Approximation and Inverse PDE Problems, *J. Comput. Phys.*, vol. **401**, p. 109020, 2020.
- Peng, J., Hampton, J., and Doostan, A., On Polynomial Chaos Expansion via Gradient-Enhanced ℓ_1 -Minimization, *J. Comput. Phys.*, vol. **310**, pp. 440–458, 2016.
- Qin, T., Chen, Z., Jakeman, J.D., and Xiu, D., Deep Learning of Parameterized Equations with Applications to Uncertainty Quantification, *Int. J. Uncertain. Quantif.*, vol. **11**, no. 2, pp. 63–82, 2021.

- Raissi, M., Perdikaris, P., and Karniadakis, G.E., Physics-Informed Neural Networks: A Deep Learning Framework for Solving Forward and Inverse Problems Involving Nonlinear Partial Differential Equations, *J. Comput. Phys.*, vol. **378**, pp. 686–707, 2019.
- Ross, A. and Doshi-Velez, F., Improving the Adversarial Robustness and Interpretability of Deep Neural Networks by Regularizing Their Input Gradients, *Proc. of the AAAI Conf. on Artificial Intelligence*, New Orleans, LA, vol. **32**, no. 1, 2018.
- Schmidhuber, J. and Hochreiter, S., Long Short-Term Memory, *Neural Comput.*, vol. **9**, no. 8, pp. 1735–1780, 1997.
- Schwab, C. and Zech, J., Deep Learning in High Dimension: Neural Network Expression Rates for Generalized Polynomial Chaos Expansions in UQ, *Anal. Appl.*, vol. **17**, no. 1, pp. 19–55, 2019.
- Shalev-Shwartz, S. and Ben-David, S., *Understanding Machine Learning: From Theory to Algorithms*, Cambridge, UK: Cambridge University Press, 2014.
- Siegel, J.W. and Xu, J., Approximation Rates for Neural Networks with General Activation Functions, *Neural Networks*, vol. **128**, pp. 313–321, 2020.
- Sirignano, J. and Spiliopoulos, K., DGM: A Deep Learning Algorithm for Solving Partial Differential Equations, *J. Comput. Phys.*, vol. **375**, pp. 1339–1364, 2018.
- Spitzbart, A., A Generalization of Hermite’s Interpolation Formula, *Am. Math. Mon.*, vol. **67**, no. 1, pp. 42–46, 1960.
- Wu, Z., Hermite–Birkhoff Interpolation of Scattered Data by Radial Basis Functions, *Approx. Theory Appl.*, vol. **8**, no. 2, pp. 1–10, 1992.
- Yan, M., Yang, J., Chen, C., Zhou, J., Pan, Y., and Zeng, Z., Enhanced Gradient Learning for Deep Neural Networks, *IET Image Process.*, vol. **16**, no. 2, pp. 365–377, 2022.
- Yang, L., Meng, X., and Karniadakis, G.E., B-PINNs: Bayesian Physics-Informed Neural Networks for Forward and Inverse PDE Problems with Noisy Data, *J. Comput. Phys.*, vol. **425**, p. 109913, 2021.
- Yu, J., Lu, L., Meng, X., and Karniadakis, G.E., Gradient-Enhanced Physics-Informed Neural Networks for Forward and Inverse PDE Problems, *Comput. Methods Appl. Mech. Eng.*, vol. **393**, p. 114823, 2022.
- Zhuang, X., Nguyen, L.C., Nguyen-Xuan, H., Alajlan, N., and Rabczuk, T., Efficient Deep Learning for Gradient-Enhanced Stress Dependent Damage Model, *Appl. Sci.*, vol. **10**, no. 7, p. 2556, 2020.