001

003 004 005

006

007 008 009

010 011

012

013

014

016

017

018

019

021

024

025

026

027

028

029

031

033 034

037

038

040

041

042

043 044

046

047

048

051

052

# WEIGHTED DEEP ENSEMBLE UNDER MISSPECIFICATION

## Anonymous authors Paper under double-blind review

#### **ABSTRACT**

Deep neural networks are supported by the universal approximation theorem, which guarantees that sufficiently large architectures can approximate smooth functions. In practice, however, this guarantee holds only under restrictive conditions, and violations of these conditions give rise to model misspecification. We categorize such misspecification into three sources: variable misspecification, arising from insufficiently informative features; structural misspecification, stemming from the limited width and depth of networks that cannot fully capture the underlying complexity; and inherent misspecification, occurring when the true model possesses properties such as discontinuities that cannot be faithfully represented. To mitigate the impact of these forms of misspecification, ensemble methods have become a common strategy for enhancing predictive performance. However, standard ensembles composed of identically architected and equally weighted models may suffer from "collective blindness", where shared errors are amplified and lead to systematically biased predictions with high confidence. To mitigate this issue, we introduce weighted deep ensemble method that learns the optimal weights. We prove that our method provably attains the convergence rate of the best single model in the ensemble and asymptotically achieves oracle-level predictive risk. To the best of our knowledge, this is the first work to provide rigorous theoretical guarantees for weighted deep ensemble under both well-specified and misspecified settings.

#### 1 Introduction

Model misspecification in statistics arises from the omission of relevant variables, inclusion of irrelevant variables, incorrect functional forms and incorrect distributional assumptions (Maasoumi, 1990; White, 1982). In such cases, the best possible approximation  $f^* \in \mathcal{F}$ , with  $\mathcal{F}$  denoting the function class used for estimation, still maintains a significant approximation error from the true function  $f_0$ . In deep learning, the neural networks are always assumed to be well-specified. As shown in the universal approximation theorem, sufficiently large neural networks have the ability to approximate any continuous function, which in principle allows the approximation error  $\|f_0 - f^*\|$  to approach zero (Hornik et al., 1989; Park et al., 2020; Lu et al., 2021). Therefore, existing studies always focus on overcoming challenges in optimization and estimation errors (Barron, 1994; Soltanolkotabi et al., 2018; Adcock & Dexter, 2021).

However, the assumption that neural networks are well-specified is frequently violated in practice, as model misspecification is common. Unlike in traditional statistics, misspecification in deep learning manifests in several distinct forms. First, it may arise from an information deficit, where the input features and their latent representations lack the necessary information to capture the true data-generating process. Second, practical constraints on network depth and width impose finite capacity, leading to non-negligible approximation error when the true function is highly complex. Finally, misspecification can occur when the true function has properties such as discontinuities, which cannot be exactly represented by neural networks and can only be approximated with non-vanishing error at the discontinuity points.

"All models are wrong, but some are useful (Box, 1976)." Ensemble methods is the most intuitive method to leverage the useful parts of multiple wrong models (Fort et al., 2019; Huang et al., 2024).

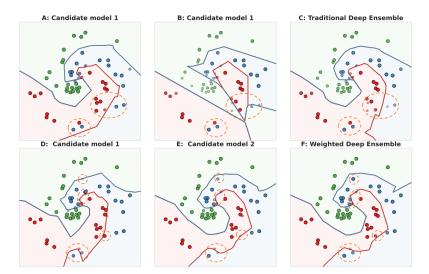


Figure 1: Comparison of the decision boundaries and confidence levels of different ensemble methods, with darker shading indicating higher confidence. The Traditional Deep Ensemble (C) shows "collective blindness" by having lower confidence in correctly classified areas but higher confidence in the misclassified areas, e.g., blue points within the red region. Weighted Deep Ensemble (F) corrects these errors by maintaining high confidence in correct areas while showing low confidence at uncertain boundaries. The key difference is highlighted in the circled area.

However, it still raises a critical question: if a neural network model is misspecified, can the estimators from an ensemble of such models still be trusted? We find that traditional deep ensembles, which consist of models with identical architectures (Zhang et al., 2020; Schweighofer et al., 2024), tend to learn in highly correlated ways when faced with the same misspecification. By deviating in the same incorrect direction, they produce an adverse effect we term "collective blindness". This phenomenon is caused by the ensemble reinforcing, rather than correcting, the biases shared by all members. As illustrated in Figure 1 (C), in the orange circle, the traditional deep ensemble produces misclassified predictions with high confidence. The root of the problem is that traditional ensemble methods not only employ similar model architectures but, more critically, typically aggregate predictions using equal weights. When all models are plagued by the same misspecification, this simplistic averaging only serves to amplify their shared error. To address this, we propose a novel and effective solution: an optimally weighted deep ensemble built upon architectural diversity. By ensuring that exploitable differences exist among the models, we can theoretically derive data-driven weights that minimize the ensemble's prediction error on a held-out validation set. As illustrated in Figure 1 (F), the use of optimal weights helps the ensemble aggregate complementary strengths of its constituent models and produce more accurate predictions than relying on a single misspecified model. Our primary contributions are as follows:

- (1) We are the first to systematically define and categorize misspecification in deep learning into variable, structural, and inherent forms. We propose weighted deep ensemble to mitigate the "collective blindness" effect seen in traditional ensembles and provide a theoretical guarantee for our weighted deep ensemble for well-specified and misspecified models.
- (2) We establish an asymptotic error bound for the weighted deep ensemble estimator and show that the bound converges at the same rate as the smallest error bound among all candidate networks. This guarantees that the ensemble inherits the speed of the best individual model, so including slower or misspecified models can never slow it down, while any rapidly converging model immediately improves the overall rate. Furthermore, we provide detailed analyses for common networks such as MLP, CNN, and RNN.
- (3) We prove that the weight vector yields a prediction risk that converges to the oracle minimum, even though the oracle weight itself depends on unknown population quantities and cannot be computed. Thus the proposed weighted choice method recovers the infeasible optimal weight asymptotically, giving the first rigorous guarantee that weighted deep ensemble can attain oracle-level accuracy using only observable data. To the best of our knowledge, this is the first study to offer a theoretical guarantee for weighted deep ensemble.

#### 2 RELATED WORK

Misspecification. Model misspecification, which occurs when a chosen model fails to accurately capture the true data-generating process, is a common challenge in statistics (McGuirk et al., 1993; Cerreia-Vioglio et al., 2025). Misspecification is categorized into several types: omitted varaible bias, where the exclusion of a relevant variables leads to biased and inconsistent parameter estimates (Gospodinov & Maasoumi, 2021), incorrect functional form, such as assuming a linear relationship when the true function is nonlinear (Gerds & Schumacher, 2001; Kasparis, 2011), and mismatch distribution when the assumed probability distribution for the error term or the response variable in models is incorrect (Masiha et al., 2021; Kuang et al., 2020). These types of misspecification always degrade model predictive performance (Lanzani, 2025). In deep learning, the universal approximation theorem states that a sufficiently large neural network can approximate any continuous function (Raghu et al., 2017; Kratsios et al., 2021). Therefore, previous research always assumed that deep models are well-specified. However, misspecification is a widespread issue in practice due to limited information and finite model capacity of neural networks with limited width and depth. Our work is the first to provide a clear definition for misspecification in deep learning and theoretical guarantees for deep ensembles under misspecified conditions.

Ensemble Learning. Deep ensembles, typically composed of identical architectures with different random initialization, have been shown to outperform single deep learning models in terms of accuracy (Lakshminarayanan et al., 2017; Mohammed & Kora, 2023). However, relying solely on the same model structure may limit the effectiveness of the ensemble. To address it, several works have introduced greater diversity by varying neural network architectures (Zhang et al., 2020) and training methods (Gontijo-Lopes et al., 2022). A large number of studies have theoretically explained the advantages of deep ensembles from the aspects of diversity (Wood et al., 2023; Jeffares et al., 2023) and generalization (Ortega et al., 2022; Odonnat et al., 2024). Specifically, Zhang et al. (2020) demonstrates that diversity in ensemble components can reduce prediction errors. Lin et al. (2024) show that averaging outputs enhances out-of-distribution generalization. However, these studies are centered around equal-weighted ensemble learning. Recent studies have delved into weighted deep ensemble (Kim et al., 2018; Matena & Raffel, 2022), but they only verified it experimentally. In contrast, we are the first to provide the theoretical guarantee for the weighted deep ensemble.

#### 3 METHODOLOGY

#### 3.1 PROBLEM SETUP

We consider a general model where the input features  $\boldsymbol{X}=(X_1,\ldots,X_d)\in\mathbb{R}^d$  and the output Y can be either real-valued or categorical. In regression tasks,  $Y\in\mathbb{R}$  and follows the model  $Y=f_0(\boldsymbol{X})+\varepsilon$ , where  $f_0$  is an unknown true function and  $\varepsilon$  is a noise term satisfying  $\mathbb{E}(\varepsilon|\boldsymbol{X})=0$ . In classification tasks with C classes,  $\boldsymbol{Y}=(Y_1,\ldots,Y_C)^{\top}\in\{0,1\}^C$  is the one-hot vector, where only one entry is 1 indicating the true class and all others are 0. The conditional probability of  $\boldsymbol{Y}$  is modeled as  $P(\boldsymbol{Y}\mid\boldsymbol{X})=f_0(\boldsymbol{X})$  for  $c=1,\ldots,C$ , where  $f_0(\boldsymbol{X})=(f_{0,1}(\boldsymbol{X}),\ldots,f_{0,C}(\boldsymbol{X}))^{\top}$  and  $f_{0,c}(\boldsymbol{X})=P(Y_c=1\mid\boldsymbol{X})$ . We assume that n independent observable samples  $(\boldsymbol{X}_i,Y_i)$  are drawn from a joint distribution over  $(\boldsymbol{X},Y)$ . The supremum norm is defined as  $\|f\|_{\infty}=\sup_{\boldsymbol{X}}|f(\boldsymbol{X})|$ , while the  $L^2$  norm is  $\|f\|_{L^2}=(\int |f(\boldsymbol{X})|^2\,dP_{\boldsymbol{X}}(\boldsymbol{X}))^{1/2}$ .

**Model training.** The observable data is split into two parts, a training set with size  $n_{\text{train}}$  for training the neural network and the other  $n_{\text{val}} = n - n_{\text{train}}$  for choosing weights. Specifically, a base model  $\hat{f}$  is obtained by empirical risk minimization:  $\hat{f} = \arg\min_{f \in \mathcal{F}} 1/n_{\text{train}} \sum_{i=1}^{n_{\text{train}}} \ell(f(\boldsymbol{X}_i), Y_i)$ , where  $\mathcal{F}$  denotes the model function class and  $\ell$  is the loss function. And  $f^*$  is defined as the minimizer of the true expected risk:  $f^* = \arg\min_{f \in \mathcal{F}} \mathbb{E}\left[\ell(f(\boldsymbol{X}), Y)\right]$ .

#### 3.2 MISSPECIFICATION IN DEEP LEARNING

In this section, we systematically define three types of misspecification in deep learning. Let  $f_0: \mathcal{X}_{\text{true}} \to \mathcal{Y}$  denote the true data-generating function.

**Definition 1** (Variable Misspecification). Let  $\mathcal{X}_{model}$  be the feature space available to a given model. Define the projection map  $\pi: \mathcal{X}_{true} \to \mathcal{X}_{model}$  that restricts each  $x \in \mathcal{X}_{true}$  to its coordinates in

 $\mathcal{X}_{model}$ . We say that the model exhibits variable misspecification if  $f_0$  cannot be expressed as a composition of  $\pi$  with a function on  $\mathcal{X}_{model}$ . Formally,

$$\nexists g: \mathcal{X}_{\text{model}} \to \mathcal{Y} \quad \text{such that} \quad f_0(x) - g(\pi(X)) \quad \text{for almost every } X \in \mathcal{X}_{\text{true}}.$$

**Example.** Consider an image classification task where the true label depends on latent features  $Z=(Z_1,Z_2)$  (e.g., ear shape and nose shape). The true function is  $f_0(Z)=\mathbb{I}_{Z_1+Z_2>0}$  (e.g., predicting cat if positive, dog if negative). If the model's feature space  $\mathcal{X}_{\text{model}}$  corresponds only to  $Z_2$  due to missing feature input or incomplete feature extraction, the model suffers from variable misspecification.

**Definition 2** (Structural Misspecification). Let  $\mathcal{H}$  be the function class corresponding to a given neural network architecture. Let  $\mathcal{L}$  be a loss function and define the risk

$$R(h) = \mathbb{E}_{\boldsymbol{X} \sim P_{\boldsymbol{X}}} \left[ \mathcal{L}(h(\boldsymbol{X}), f_0(\boldsymbol{X})) \right].$$

For a tolerance level  $\delta > 0$ , we say that the architecture exhibits **structural misspecification** if the minimal achievable risk within  $\mathcal{H}$  exceeds this tolerance. Formally,

$$\inf_{h \in \mathcal{H}} R(h) > \delta.$$

**Example.** Consider ReLU neural networks with both depth and width restricted to the order of  $\log(n)$ . For  $\beta$ -smooth functions in dimension d, it is known that the approximation error in this class satisfies

$$\inf_{h \in \mathcal{H}} R(h) = c \log(n)^{-4\beta/d},$$

for some constant c>0 depending only on  $\beta$  and d (Jiao et al., 2023). We set the tolerance level to  $\delta=n^{-1/4}$ , which corresponds to the rate condition commonly required in double machine learning (Chernozhukov et al., 2018). For sufficiently large n, we have

$$\inf_{h \in \mathcal{H}} R(h) = c \log(n)^{-4\beta/d} > \delta,$$

so the architecture is structurally misspecified relative to the tolerance  $\delta$ .

**Definition 3** (Inherent Misspecification). Let  $\mathcal{H}$  be the function class corresponding to a given neural network architecture. Approximation results for neural networks typically require  $f_0$  to belong to a smoothness class, such as a Hölder or Sobolev space, in order to guarantee vanishing  $L^2$  approximation error. We say that the architecture exhibits **inherent misspecification** if  $f_0$  does not satisfy the required smoothness conditions, so that the minimal achievable  $L^2$  error is bounded away from zero. Formally, for some tolerance level  $\delta > 0$ ,

$$\inf_{h \in \mathcal{H}} \|h - f_0\|_{L^2} > \delta.$$

**Example.** Consider the Dirichlet function

$$f_0(x) = \begin{cases} 1, & x \in \mathbb{Q} \cap [0, 1], \\ 0, & x \in (\mathbb{R} \setminus \mathbb{Q}) \cap [0, 1]. \end{cases}$$

This function is nowhere continuous and does not belong to any Hölder or Sobolev class. As a result, neural networks cannot approximate  $f_0$  in  $L^2$  with vanishing error, and the approximation gap remains strictly positive. Therefore, the model class  $\mathcal{H}$  suffers from inherent misspecification.

#### 3.3 WEIGHTED DEEP ENSEMBLE

First, we introduce the standard deep ensemble method.

**Deep ensembles.** A standard deep ensemble consists of M candidate models  $\widehat{f}_1(\cdot),\ldots,\widehat{f}_M(\cdot)$ . The deep ensemble prediction  $\bar{f}(\boldsymbol{X})$  is:  $\bar{f}(\boldsymbol{X}) = \sum_{m=1}^M w_m \widehat{f}_m(\boldsymbol{X})$ , where all weights are set equally as  $w_m = 1/M$ . Typically, the candidate models share the same neural network architecture and are trained independently on the same dataset using the same loss function  $\ell(f_m(\boldsymbol{X}),Y)$ , which corresponds to minimizing the empirical risk:  $\mathcal{L}_{\text{avg}} = \mathbb{E}\left[\ell(\bar{f}(\boldsymbol{X}),Y)\right]$ .

As stated before, traditional deep ensembles may suffer from "collective blindness" in the presence of variable, structural, or inherent misspecification. This motivates us to apply a more flexible weighting scheme.

Weighted deep ensemble. In this paper, we propose a method called Weighted Deep Ensemble (WDE), which combines multiple neural networks with different structures into a single predictive model. Let  $\boldsymbol{w} = (w_1, \dots, w_M)^\mathsf{T}$  be the vector of ensemble weights. We restrict the weights to the simplex  $\mathcal{W} = \{\boldsymbol{w} \in [0,1]^M : \sum_{m=1}^M w_m = 1\}$ . The weighted deep ensemble prediction is defined as  $\widehat{f}(\boldsymbol{X}; \boldsymbol{w}) = \sum_{m=1}^M w_m \widehat{f}_m(\boldsymbol{X})$ . In practice, the weight vector  $\boldsymbol{w}$  is unknown and need to be estimated from observable data.

Weight choice criterion. We adopt a validation-risk minimization (VRM) criterion: the estimator  $\hat{w}$  is chosen to minimize the empirical loss on a validation set, under the simplex constraints

$$\widehat{\boldsymbol{w}} = \operatorname*{arg\,min}_{\boldsymbol{w} \in \mathcal{W}} \frac{1}{n_{\mathrm{val}}} \sum_{i=1}^{n_{\mathrm{val}}} \ell(f(\boldsymbol{X}_i; \boldsymbol{w}), Y_i).$$

This VRM strategy links the weighting scheme directly to out-of-sample performance, introducing no extra hyperparameters beyond the standard validation split. Specifically, we consider two general tasks: (i) Regression: with the squared loss  $\ell_{\text{reg}}\big(f(\pmb{X};\pmb{w}),\pmb{Y}\big) = \big(f(\pmb{X};\pmb{w})-\pmb{Y}\big)^2$ , the VRM objective reduces to a strictly convex quadratic program on the simplex  $\mathcal{W}$  (Li et al., 2023; Qu et al., 2025). (ii) Classification: Using cross-entropy loss  $\ell_{\text{class}}\big(f(\pmb{X};\pmb{w}),\pmb{Y}\big) = -\pmb{Y}^{\top}\log(f(\pmb{X};\pmb{w}))$ , where  $\pmb{Y} \in \{0,1\}^C$  denotes the one-hot encoding of the true class, and  $f(\pmb{X};\pmb{w})$  denotes the predicted probability. Since the mapping  $\pmb{w}\mapsto f(\pmb{X};\pmb{w})$  is affine and  $-\log(\cdot)$  is convex, the loss function remains convex in  $\pmb{w}$  and can be optimized using projected-gradient (Boyd & Vandenberghe, 2004) methods. So we construct the weighted deep ensemble estimator  $\widehat{f}(\pmb{X};\widehat{\pmb{w}}) = \sum_{m=1}^M \widehat{w}_m \widehat{f}_m(\pmb{X})$ . Due to the convexity of the loss function with respect to  $\pmb{w}$ , we can obtain the global optimal solution for the ensemble weights.

#### 3.4 ASYMPTOTIC ERROR BOUNDS

Our goal is to establish the asymptotic error bound of the proposed weighted deep ensemble estimator. Specifically, we aim to show that our proposed weighted deep ensemble effectively combines multiple small neural networks to achieve an asymptotic error bound at least as fast as that of the best candidate. By appropriately combining these models, the estimator adapts to various data structures and retains the ability to capture intricate features without resorting to a large, monolithic network. Consequently, the weighted deep ensemble benefits from more flexible modeling choices while still maintaining an asymptotic error bound that is as fast as the best candidate. Importantly, the presence of misspecified or poorly performing candidate models does not slow down the convergence rate of the estimator. We provide theoretical guarantees for both regression (Theorem 1) and classification (Theorem 2) settings. To establish these results, we begin by introducing the following condition:

**Condition 1.** (i). There exists a positive constant C such that  $||f_0(X)||_{\infty} < C$ ,  $||\widehat{f}_m(X)||_{\infty} < C$  for  $m = 1, \ldots, M$ ; (ii).  $\mathbb{E}(\varepsilon | X) = 0$ , and  $\varepsilon$  is sub-Gaussian with parameter  $\sigma$ .

This condition restricts the upper bound of  $f_0$ ,  $f^*$ , and that  $\varepsilon$  has mean zero and sub-Gaussian tails. This condition is also widely used in the literature; see Schmidt-Hieber (2020) and Jiao et al. (2023) for example.

**Theorem 1.** Suppose Condition 1 holds and assume that the candidate model with the fastest convergence rate has an asymptotic error bound of order S, i.e., the candidate model with the fastest asymptotic error bound (denoted as  $\tilde{f}$  without loss of generality) satisfies  $||f_0(X) - \tilde{f}(X)||_{L^2} = O_p(S)$ , then our weighted deep ensemble estimator can also achieve this rate asymptotically:

$$||f_0(\mathbf{X}) - \widehat{f}(\mathbf{X}; \widehat{\mathbf{w}})||_{L^2} = O_p(S + \sqrt{\frac{1}{n}}).$$

Besides the asymptotic error bound S, the result also includes an additional term of order  $1/\sqrt{n}$ . In practice, the minimax rate for nonparametric methods such as neural networks is typically of order  $n^{-\beta/(2\beta+d)}$  for some smoothness  $\beta$  and input dimension d (Stone, 1982), which is slower than

 $1/\sqrt{n}$ ; hence the overall asymptotic error bound is dominated by the nonparametric estimation error. Therefore, Theorem 1 establishes that the asymptotic error bound of the weighted deep ensemble estimator is no worse than that of the best individual candidate model. In other words, regardless of which candidate model achieves the smallest asymptotic error bound, the ensemble procedure guarantees at least comparable asymptotic performance and will not converge more slowly than that benchmark.

Theorem 1 presents the asymptotic error bound of the deep ensemble estimator under regression tasks. In fact, we can establish similar properties for classification tasks. Before presenting the theorem, we define  $R_{0/1}(f) = \mathbb{E}\{I(\arg\max_c Y_c \neq \arg\max_c f_c(\boldsymbol{X}))\} - \mathbb{E}\{I(\arg\max_c Y_c \neq \arg\max_c f_{0,c}(\boldsymbol{X}))\}$  to be the excess misclassification rate.

**Theorem 2.** Suppose  $f_0(X)$  is uniformly bounded away from 0 and 1 and assume that the best candidate model has a misclassification rate of S, i.e., the candidate model with the smallest misclassification rate (denoted as  $\tilde{f}$  without loss of generality) satisfies  $R_{0/1}(\tilde{f}) = O_p(S)$ , then our weighted deep ensemble estimator can also achieve this misclassification rate asymptotically:

$$R_{0/1}(\widehat{f}(\boldsymbol{X}; \widehat{\boldsymbol{w}})) = O_p(S + \sqrt{\frac{1}{n}}).$$

Theorems 1 and 2 show that, in both regression and classification tasks, the asymptotic error bound of our weighted deep ensemble estimator matches the asymptotic error bound of the best single candidate model. Detailed proofs are provided in the Appendix.

In addition, when the pool of candidate models is altered, the estimator's attainable asymptotic error bound necessarily shifts in response to the new composition. For better understanding, the next three corollaries provide the asymptotic error bound of the weighted deep ensemble estimator when the candidate models are chosen from MLP-based networks, CNN-based networks, and RNN-based networks. It is worth noting that Theorem 1 holds under the listed moment conditions, but additional assumptions are implicitly embedded in the asymptotic error bound  $O_p(S)$ . Since different neural network architectures require different conditions to guarantee convergence, we do not enumerate them explicitly here and instead present the result in terms of the general order  $O_p(S)$ . Accordingly, the subsequent corollaries impose further conditions on the candidate models, ensuring that they can indeed attain the corresponding asymptotic error bound in those specific settings. Define  $\widetilde{O}_p(\cdot)$  as the rate by ignoring logarithmic factors.

**Corollary 1** (MLP case). If all candidate models are MLP-based models,  $f_0$  is  $\beta$ -Hölder smooth with  $\beta > 1$ , and Conditions of Theorem 4.2 in Jiao et al. (2023) holds, then with some specifically designed candidate models, the weighted deep ensemble estimators can achieve asymptotic error bound of

$$||f_0(\mathbf{X}) - \widehat{f}(\mathbf{X}; \widehat{\mathbf{w}})||_{L^2}^2 = \widetilde{O}_p(n^{-2\beta/(d+2\beta)}).$$

**Corollary 2** (CNN case). If all candidate models are CNN-based models,  $f_0$  is  $\beta$ -Hölder smooth with  $\beta > 1$ , and Conditions of Theorem 4.6 in Shen et al. (2022) holds, then with some specifically designed candidate models, the weighted deep ensemble estimators can achieve asymptotic error bound of

$$||f_0(\mathbf{X}) - \widehat{f}(\mathbf{X}; \widehat{\mathbf{w}})||_{L^2}^2 = \widetilde{O}_p(n^{-2\beta/(d+2\beta)}).$$

**Corollary 3** (RNN case). If all candidate models are RNN-based models,  $f_0$  is  $\beta$ -Hölder smooth with  $\beta > 1$  and Conditions of Theorem 3 and Theorem 5 in Jiao et al. (2024) holds, then with some specifically designed candidate models, the weighted deep ensemble estimator can achieve asymptotic error bound of

$$||f_0(\mathbf{X}) - \widehat{f}(\mathbf{X}; \widehat{\mathbf{w}})||_{L^2}^2 = \widetilde{O}_p(n^{-2\beta/(dl+2\beta)}),$$

where l is the length of the input sequence, i.e., the number of time steps processed by the RNN.

These results can also be extended to ResNet, Transformer, and other network structures. As long as we have the asymptotic error bound of a single network, Theorems 1 and 2 ensure that the weighted deep ensemble estimator can achieve the same asymptotic error bound as the fastest candidate model.

#### 3.5 ASYMPTOTIC OPTIMALITY UNDER MODEL MISSPECIFICATION

Unlike traditional ensembles with uniform weights, our method learns data-dependent weights, ensuring that the ensemble performs at least as well as the best candidate asymptotically. Since

equal weights lie within our admissible weight space, our approach is guaranteed to match or exceed the performance of equal-weighted averaging asymptotically. In practice, the oracle weight vector is unknown and cannot be directly computed from observable data. By instead using the VRM criterion, our method significantly reduces computational cost while maintaining competitive efficiency, making it particularly suitable for large-scale scenarios. In the following, we present theoretical results showing that the proposed weighted deep ensemble estimator achieves asymptotic optimality and attains oracle-level accuracy using only observable data.

Let  $R(\boldsymbol{w}) = \|f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w})\|_{L^2}^2$ ,  $R^*(\boldsymbol{w}) = \|f_0(\boldsymbol{X}) - f^*(\boldsymbol{X}; \boldsymbol{w})\|_{L^2}^2$ ,  $\xi_n = \inf_{\boldsymbol{w} \in \mathcal{W}} R^*(\boldsymbol{w})$ , and  $\phi_n = \sup_{\boldsymbol{w} \in \mathcal{W}} \|\widehat{f}(\boldsymbol{X}; \boldsymbol{w}) - f^*(\boldsymbol{X}; \boldsymbol{w})\|_{L^2}$ . To establish the asymptotic optimality of  $\widehat{\boldsymbol{w}}$ , we require the following condition.

Condition 2. (i). 
$$\xi_n^{-1} n^{-1/2} = o(1)$$
; (ii).  $\xi_n^{-1} \phi_n = o_p(1)$ .

This condition regulates the divergence speed of  $\xi_n$ , and it is frequently used in FMA research, such as Ando & Li (2014); Zhang et al. (2016). Condition 2 requires  $\xi_n$  to grow faster than  $\sqrt{1/n}$  and  $\phi_n$ . Importantly, this condition should be interpreted as a *misspecification condition*, ensuring that the oracle risk  $\xi_n$  does not vanish too quickly relative to the estimation error. It naturally aligns with the three forms of misspecification introduced above:

- (1) **Variable misspecification.** When crucial variables are missing, or when the essential features of  $f_0$  cannot be generated from the available input space,  $f^*$  cannot converge to  $f_0$ . In this case, the approximation error remains bounded away from zero, which implies that the oracle risk  $\xi_n$  does not vanish and Condition 2 (i) is satisfied. This requirement should be viewed as a stronger form of variable misspecification, as it excludes cases where the omitted variables have only negligible influence on  $f_0$ , for instance when their contribution diminishes asymptotically.
- (2) **Structural misspecification.** When networks have limited depth or width, their approximation error remains bounded away from zero. In this case  $\xi_n$  decreases at a much slower rate than  $n^{-1/2}$ , so Condition 2 (i) is satisfied.
- (3) **Inherent misspecification.** When the true function lies outside the smoothness classes typically required for neural network approximation, the minimal  $L^2$  error remains strictly positive. Consequently,  $\xi_n$  is bounded away from zero, and Condition 2 (i) holds.

Unlike  $\xi_n$ , the term  $\phi_n$  measures the estimation error between  $\widehat{f}$  and  $f^*$ , which depends on sample size rather than model specification. Condition 2 (ii) therefore requires that  $\phi_n$  converges to zero at a faster rate than  $\xi_n$ , ensuring that estimation error does not dominate the asymptotic behavior.

**Theorem 3.** Suppose Conditions 1 and 2 hold,

$$\frac{R(\widehat{\boldsymbol{w}})}{\inf_{\boldsymbol{w}\in\mathcal{W}}R(\boldsymbol{w})}\to 1$$

in probability as  $n \to \infty$ .

This theorem states that under Conditions 1 and 2, as the sample size n goes to infinity, the ratio of the risk of  $\widehat{w}$  to the infimum of the risk over all possible weights converges to 1. In other words, although the weight that minimizes the risk is infeasible, the proposed weight choice criterion identifies a weight  $\widehat{w}$  whose risk becomes asymptotically equal to the minimal risk. This means our procedure performs asymptotically as well as this ideal benchmark. This provides strong theoretical support for our weight selection method, assuring that no asymptotic loss is incurred compared to the unattainable optimum. All theoretical proofs are provided in the Appendix B.

#### 4 Numerical Results

In this section, we investigate the models that are well-specified or suffer from variable misspecification, structural misspecification, and inherent misspecification.

**Baselines.** We compare several ensemble strategies: (1) Deep Ensemble (DE): homogeneous MLPs with different initializations and equal weights; (2) Equal-Weight Heterogeneous (EW): heterogeneous MLPs with equal weights; (3) Our Method, WDE: heterogeneous MLPs with optimal weights.

**Experimental details.** Datasets are split 6:2:2 for training, validation, and testing. We use Adam with early stopping (patience=20 epochs). Hyperparameters: learning rate searched in [0.001, 0.1], batch size 128, max 5000 epochs. Our method uses 4 heterogeneous networks with total parameters matching a 4×MLP (2 hidden layers of 30 nodes) DE for fair comparison.

Table 1: Performance comparison across different complexity types with varying missing variables, with sample size 5000.

Task	# Missing	DE	$\mathbf{EW}$	WDE	$\Delta_{ m DE}$	$\Delta_{\mathrm{EW}}$
	0	$7.452 \pm 0.319$	$7.711 \pm 0.527$	$7.025 \pm 0.247$	6.07%	9.76%
	1	$8.175 \pm 0.472$	$8.686 \pm 0.957$	$\boldsymbol{8.028 \pm 1.236}$	1.84%	8.19%
Nested	3	$8.390 \pm 0.681$	$8.825 \pm 0.814$	$8.077\pm1.211$	3.88%	9.27%
	5	$8.440 \pm 0.788$	$8.493 \pm 0.482$	$7.792 \pm 0.287$	8.31%	9.00%
	7	$9.122 \pm 0.498$	$8.687 \pm 0.510$	$7.967 \pm 0.431$	14.50%	9.04%
	0	$1.975 \pm 0.330$	$1.854 \pm 0.187$	$1.773\pm0.131$	11.35%	4.57%
	1	$3.251 \pm 0.740$	$3.640 \pm 0.782$	$2.949 \pm 0.778$	10.22%	23.43%
Interaction	3	$3.949 \pm 0.991$	$4.434 \pm 1.046$	$3.890 \pm 0.984$	20.76%	13.98%
	5	$5.239 \pm 1.149$	$5.223 \pm 1.176$	$5.077 \pm 1.099$	22.36%	2.87%
	7	$5.747 \pm 1.206$	$5.617 \pm 1.266$	$5.548 \pm 1.168$	13.27%	1.24%
	0	$2.154 \pm 0.101$	$2.197 \pm 0.121$	$2.023 \pm 0.116$	6.47%	8.59%
	1	$2.578 \pm 0.175$	$2.844 \pm 0.419$	$2.474 \pm 0.193$	4.17%	14.93%
Periodic	3	$2.994 \pm 0.213$	$3.148 \pm 0.439$	$2.901 \pm 0.251$	3.19%	8.48%
	5	$3.631 \pm 0.254$	$3.666 \pm 0.410$	$3.362 \pm 0.273$	7.98%	9.04%
	7	$3.841 \pm 0.316$	$3.932 \pm 0.481$	$3.714 \pm 0.331$	12.52%	5.87%

Well-specified Models and Variable Misspecification. We consider some simple data-generating processes (DGPs) that can be well-approximated by simple MLPs, i.e., the approximation error is small. Let  $\boldsymbol{X} \in \mathbb{R}^p$  be a random vector where each feature  $X_j$  is independently sampled from  $\mathcal{N}(0,1)$ , with the number of features p=10. We define the different DGPs: (i) Nested:  $f_0(\boldsymbol{x}) = \sin(\sum_{j=1}^{10} x_j^2) + \sum_{j=1}^{10} (\cos x_j)^2$ .; (ii) Interaction:  $f_0(\boldsymbol{X}) = \frac{1}{2} \sum_{j=1}^{5} \sum_{k=6}^{10} X_j X_k$ ; (iii) Periodic:  $f_0(\boldsymbol{X}) = \sum_{j=1}^{5} \sin(X_j) + \sum_{j=6}^{10} \cos(X_j)$ .

When all relevant features are included in the model, the setting is considered well-specified. We introduce variable misspecification by randomly dropping a subset of features during training, with the number of dropped features controlling the degree of misspecification. The results are summarized in Table 2. Our method WDE consistently achieves the lowest MSE.

Table 2: Performance comparison across different parameter discrepancies with sample size 5000

Δ	α	DE	EW	WDE	$\Delta_{ ext{DE}}$	$\Delta_{\mathrm{EW}}$
	0.2	$0.598 \pm 0.222$	$0.402 \pm 0.085$	$0.341 \pm 0.034$	75.47%	17.83%
30000	0.5	$0.630 \pm 0.259$	$0.438 \pm 0.179$	$0.397 \pm 0.146$	58.59%	10.37%
30000	0.7	$0.634 \pm 0.230$	$0.473 \pm 0.128$	$0.398 \pm 0.083$	59.53%	18.87%
	0.9	$0.654 \pm 0.157$	$0.524 \pm 0.185$	$0.417 \pm 0.089$	56.83%	25.76%
	0.2	$0.587 \pm 0.159$	$0.374 \pm 0.075$	$0.339 \pm 0.032$	73.10%	10.22%
50000	0.5	$0.632 \pm 0.193$	$0.425 \pm 0.098$	$0.399 \pm 0.091$	58.31%	6.49%
30000	0.7	$0.701 \pm 0.201$	$0.457 \pm 0.099$	$0.419 \pm 0.067$	67.21%	9.17%
	0.9	$0.790 \pm 0.171$	$0.488 \pm 0.074$	$0.437 \pm 0.061$	80.65%	11.64%
	0.2	$0.576 \pm 0.206$	$0.335 \pm 0.031$	$0.331 \pm 0.026$	74.31%	1.51%
100000	0.5	$0.582 \pm 0.175$	$0.350 \pm 0.050$	$0.340 \pm 0.039$	71.42%	3.15%
100000	0.7	$0.638 \pm 0.139$	$0.353 \pm 0.052$	$0.346 \pm 0.043$	84.26%	1.84%
	0.9	$0.677 \pm 0.226$	$0.350 \pm 0.044$	$0.348 \pm 0.043$	94.50%	0.50%

Structural Misspecification. To investigate structural misspecification, we define the DGP as a convex combination of a simple function and a more complex function:  $f_0(X) = \alpha f_{\text{simple}}(X) + (1-\alpha) f_{\text{complex}}(X)$ , where  $f_{\text{simple}} = \text{MLP}_{2\times30}$  and  $f_{\text{complex}} = \text{MLP}_{3\times100}$ . The function used for fitting is  $f_{\text{simple}}(X)$ . The parameter  $\alpha \in [0,1]$  controls the degree of misspecification, from well-specified ( $\alpha=0$ ) to totally misspecified ( $\alpha=1$ ). We report the relative MSE under varying degrees of misspecification with a sample size of N=5000 in Table 2. The results show that as the misspecification degree increases, the MSE of DE grows rapidly. In contrast, our method (WDE) significantly mitigates this performance degradation and maintains robust accuracy even under high misspecification.

Table 3: Comparison with relative improvement percentages ( $\Delta$ ) across DE, EW and WDE.

Complexity	N	DE	EW	WDE	$\Delta_{ m DE}$	$\Delta_{\mathrm{EW}}$
Square Wave	5000 10000 15000	$0.96068 \pm 0.02707$ $0.81132 \pm 0.02665$ $0.74766 \pm 0.03754$	$0.94702 \pm 0.01385$ $0.80467 \pm 0.02623$ $0.71625 \pm 0.03140$	$\begin{array}{c} 0.93171 \pm 0.02422 \\ 0.76520 \pm 0.03430 \\ 0.64963 \pm 0.06092 \end{array}$	3.11% 6.03% 15.09%	1.64% 5.16% 10.25%
Infinite	5000 10000 15000	$0.00137 \pm 0.00021$ $0.00085 \pm 0.00015$ $0.00057 \pm 0.00013$	$\begin{array}{c} 0.00141 \pm 0.00010 \\ 0.00086 \pm 0.00013 \\ 0.00065 \pm 0.00010 \end{array}$	$\begin{array}{c} 0.00117 \pm 0.00012 \\ 0.00067 \pm 0.00014 \\ 0.00044 \pm 0.00012 \end{array}$	17.34% 27.38% 30.13%	20.40% 28.86% 49.55%

Inherent Misspecification. We consider two discontinuous DGPs to simulate inherent misspecification: (i) Infinite Discontinuity:  $f_0(\boldsymbol{X}) = \sum_{j=1}^d \frac{1}{j} \left( \sum_{k=1}^K \frac{1}{k^2} \cdot \mathbb{I}\left(X_j > \frac{1}{k}\right) \right)$  is the indicator function and K = 1000; (ii) Square Wave:  $f_0(\boldsymbol{X}) = \sum_{j=1}^d \frac{1}{j} \left( \sum_{k=1}^K \frac{1}{k^2} \cdot \mathbb{I}\left(X_j > \frac{1}{k}\right) \right)$ , where  $\operatorname{sgn}(\cdot)$  is the signum function. The results in Table 3 demonstrate that WDE achieves substantial improvement.

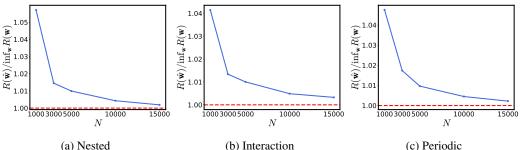


Figure 2: The ratio of MSE of WDE to the oracle weight as sample size increases.

**Theoretical results.** To further validate Theorem 3, we plot the ratio of the MSE of our proposed weighted deep ensemble estimator to that of the oracle weight as the sample size increases, as shown in Figure 2. We can observe that as the sample size grows, the weighted deep ensemble estimator asymptotically approaches the optimal oracle weight using only the observed data.

**Comparing with other weighted methods.** We compare our method against two alternatives in Table 4: (i) Greedy Ensemble: Sequentially adds models to minimize validation loss, then uses equal weighting. (ii) In-sample Ensemble: Weights are optimized directly on the training MSE.

More simulation results can be found in Section D in the Appendix.

Table 4: Performance comparison of weighted methods in variable misspecification

Complexity	# Missing	WDE	In-sample Ensemble	<b>Greedy Ensemble</b>
	0	$7.025 \pm 0.247$	$7.258 \pm 0.329$	$7.136 \pm 0.232$
	1	$8.028 \pm 1.236$	$8.149 \pm 1.150$	$8.108 \pm 1.187$
Nested	3	$8.077\pm1.211$	$8.186 \pm 1.208$	$8.177 \pm 1.217$
	5	$7.792\pm0.287$	$7.963 \pm 0.338$	$7.904 \pm 0.346$
	7	$7.967 \pm 0.431$	$8.076 \pm 0.393$	$8.045 \pm 0.400$
	0	$1.773 \pm 0.131$	$1.869 \pm 0.161$	$1.794 \pm 0.116$
	1	$2.949 \pm 0.778$	$2.990 \pm 0.813$	$2.948 \pm 0.778$
Interaction	3	$3.890 \pm 0.984$	$3.914 \pm 0.995$	$3.908 \pm 0.981$
	5	$5.077\pm1.099$	$5.114 \pm 1.038$	$5.125 \pm 1.084$
	7	$5.548 \pm 1.168$	$5.651 \pm 1.124$	$5.547 \pm 1.170$
	0	$2.023 \pm 0.116$	$2.086 \pm 0.137$	$2.051 \pm 0.126$
	1	$\boldsymbol{2.474 \pm 0.193}$	$2.535 \pm 0.178$	$2.518 \pm 0.226$
Periodic	3	$2.901 \pm 0.251$	$2.961 \pm 0.261$	$2.929 \pm 0.260$
	5	$3.362 \pm 0.273$	$3.400 \pm 0.270$	$3.378 \pm 0.279$
	7	$3.714 \pm 0.331$	$3.751 \pm 0.310$	$3.730 \pm 0.336$

#### 5 CONCLUSION

Our paper formally defines the misspecification problem in deep learning and establishes a theoretical foundation for the weighted deep ensemble, including its error bound and asymptotic optimality. Extensive experiments demonstrate that our method consistently outperforms traditional deep ensembles across various misspecification scenarios and significantly mitigates the adverse effects of model misspecification.

#### ETHICS STATEMENT

Our work is committed to the highest standards of scientific excellence, grounded in the principles of honesty, reliability, and transparency. The core technical contribution of this paper is to address the challenge of model misspecification in deep learning. This is not merely a technical problem but an ethical imperative. A misspecified model can produce unreliable predictions, perpetuate and amplify societal biases, and ultimately cause harm if deployed in critical real-world applications such as healthcare, finance, or autonomous systems. By developing methods to better understand, identify, and correct for misspecification, our research aims to contribute to the creation of more robust, fair, and trustworthy AI systems. We believe that this work is a necessary step toward the responsible development of artificial intelligence, ensuring that its benefits can be realized while minimizing potential negative societal consequences.

#### 7 REPRODUCIBILITY STATEMENT

We are committed to the full reproducibility of our work. To this end, we have included the complete and detailed derivations of our theoretical proofs in the Appendix. Furthermore, all experimental results presented in this paper can be fully reproduced according to experimental details and we will upload the source code.

#### REFERENCES

- Ben Adcock and Nick Dexter. The gap between theory and practice in function approximation with deep neural networks. *SIAM Journal on Mathematics of Data Science*, 3(2):624–655, 2021.
- Tomohiro Ando and Ker-Chau Li. A model-averaging approach for high-dimensional regression. *Journal of the American Statistical Association.*, 109(505):254–265, 2014.
- Andrew R Barron. Approximation and estimation bounds for artificial neural networks. *Machine learning*, 14(1):115–133, 1994.
- Peter L Bartlett, Olivier Bousquet, and Shahar Mendelson. Local rademacher complexities. *The Annals of Statistics*, 33(4):1497–1537, 2005.
- George EP Box. Science and statistics. *Journal of the American Statistical Association*, 71(356): 791–799, 1976.
- Stephen P Boyd and Lieven Vandenberghe. Convex optimization. Cambridge university press., 2004.
- Simone Cerreia-Vioglio, Lars Peter Hansen, Fabio Maccheroni, and Massimo Marinacci. Making decisions under model misspecification. *Review of Economic Studies*, pp. rdaf046, 2025.
- Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, pp. C1–C68, 2018.
- Stanislav Fort, Huiyi Hu, and Balaji Lakshminarayanan. Deep ensembles: A loss landscape perspective. *arXiv preprint arXiv:1912.02757.*, 2019.
- TA Gerds and Martin Schumacher. On functional misspecification of covariates in the cox regression model. *Biometrika*, 88(2):572–580, 2001.
- Raphael Gontijo-Lopes, Yann Dauphin, and Ekin Dogus Cubuk. No one representation to rule them all: Overlapping features of training methods. In *International Conference on Learning Representations.*, 2022.
- Nikolay Gospodinov and Esfandiar Maasoumi. Generalized aggregation of misspecified models: With an application to asset pricing. *Journal of Econometrics*, 222(1):451–467, 2021.
- Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.

- Yichong Huang, Xiaocheng Feng, Baohang Li, Yang Xiang, Hui Wang, Ting Liu, and Bing Qin.
   Ensemble learning for heterogeneous large language models with deep parallel collaboration.
   Advances in Neural Information Processing Systems., 37:119838–119860, 2024.
  - Alan Jeffares, Tennison Liu, Jonathan Crabbé, and Mihaela van der Schaar. Joint training of deep ensembles fails due to learner collusion. *Advances in Neural Information Processing Systems.*, 36: 13559–13589, 2023.
  - Yuling Jiao, Guohao Shen, Yuanyuan Lin, and Jian Huang. Deep nonparametric regression on approximate manifolds: Nonasymptotic error bounds with polynomial prefactors. *The Annals of Statistics*., 51(2):691–716, 2023.
  - Yuling Jiao, Yang Wang, and Bokai Yan. Approximation bounds for recurrent neural networks with application to regression. *arXiv preprint arXiv:2409.05577.*, 2024.
  - Ioannis Kasparis. Functional form misspecification in regressions with a unit root. *Econometric Theory*, 27(2):285–311, 2011.
  - Wonsik Kim, Bhavya Goyal, Kunal Chawla, Jungmin Lee, and Keunjoo Kwon. Attention-based ensemble for deep metric learning. In *European conference on computer vision.*, pp. 736–751, 2018.
  - Anastasis Kratsios, Behnoosh Zamanlooy, Tianlin Liu, and Ivan Dokmanić. Universal approximation under constraints is possible with transformers. *arXiv preprint arXiv:2110.03303*, 2021.
  - Kun Kuang, Ruoxuan Xiong, Peng Cui, Susan Athey, and Bo Li. Stable prediction with model misspecification and agnostic distribution shift. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 4485–4492, 2020.
  - Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in Neural Information Processing Systems.*, 30, 2017.
  - Giacomo Lanzani. Dynamic concern for misspecification. *Econometrica*, 93(4):1333–1370, 2025.
  - Shaoze Li, Zhibin Deng, Cheng Lu, Junhao Wu, Jinyu Dai, and Qiao Wang. An efficient global algorithm for indefinite separable quadratic knapsack problems with box constraints. *Computational Optimization and Applications.*, 86(1):241–273, 2023.
  - Yong Lin, Lu Tan, Yifan Hao, Honam Wong, Hanze Dong, Weizhong Zhang, Yujiu Yang, and Tong Zhang. Spurious feature diversification improves out-of-distribution generalization. In *International Conference on Learning Representations.*, 2024. URL https://openreview.net/forum?id=d6H4RBi7RH.
  - Jianfeng Lu, Zuowei Shen, Haizhao Yang, and Shijun Zhang. Deep network approximation for smooth functions. *SIAM Journal on Mathematical Analysis*, 53(5):5465–5506, 2021.
  - Esfandiar Maasoumi. How to live with misspecification if you must. *Journal of Econometrics*, 44 (1-2):67–86, 1990.
  - Mohammad Saeed Masiha, Amin Gohari, Mohammad Hossein Yassaee, and Mohammad Reza Aref. Learning under distribution mismatch and model misspecification. In *2021 IEEE International Symposium on Information Theory (ISIT)*, pp. 2912–2917. IEEE, 2021.
  - Michael S Matena and Colin A Raffel. Merging models with fisher-weighted averaging. *Advances in Neural Information Processing Systems.*, 35:17703–17716, 2022.
  - Anya M McGuirk, Paul Driscoll, and Jeffrey Alwang. Misspecification testing: a comprehensive approach. *American Journal of Agricultural Economics*, 75(4):1044–1055, 1993.
  - Ammar Mohammed and Rania Kora. A comprehensive review on ensemble deep learning: Opportunities and challenges. *Journal of King Saud University-Computer and Information Sciences.*, 35 (2):757–774, 2023.

- Ambroise Odonnat, Vasilii Feofanov, and Ievgen Redko. Leveraging ensemble diversity for robust self-training in the presence of sample selection bias. In *International Conference on Artificial Intelligence and Statistics.*, pp. 595–603. PMLR, 2024.
  - Luis A Ortega, Rafael Cabañas, and Andres Masegosa. Diversity and generalization in neural network ensembles. In *International Conference on Artificial Intelligence and Statistics.*, pp. 11720–11743. PMLR, 2022.
  - Sejun Park, Chulhee Yun, Jaeho Lee, and Jinwoo Shin. Minimum width for universal approximation. *arXiv preprint arXiv:2006.08859*, 2020.
  - Guangtai Qu, Shaoze Li, Zhibin Deng, and Cheng Lu. A fast global algorithm for multi-linearly constrained separable binary quadratic program. *Journal of Industrial and Management Optimization.*, 21(2):1456–1473, 2025.
  - Maithra Raghu, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the expressive power of deep neural networks. In *international conference on machine learning*, pp. 2847–2854. PMLR, 2017.
  - Anselm Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with ReLU activation function. *The Annals of statistics*, 48(4):1875–1897, 2020.
  - Kajetan Schweighofer, Adrian Arnaiz-Rodriguez, Sepp Hochreiter, and Nuria Oliver. The disparate benefits of deep ensembles. *arXiv preprint arXiv:2410.13831.*, 2024.
  - Guohao Shen, Yuling Jiao, Yuanyuan Lin, and Jian Huang. Approximation with cnns in sobolev space: with applications to classification. *Advances in Neural Information Processing Systems*., 35:2876–2888, 2022.
  - Mahdi Soltanolkotabi, Adel Javanmard, and Jason D Lee. Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *IEEE Transactions on Information Theory*, 65(2):742–769, 2018.
  - Charles J Stone. Optimal global rates of convergence for nonparametric regression. *The Annals of Statistics*, pp. 1040–1053, 1982.
  - Ambuj Tewari and Peter L Bartlett. On the consistency of multiclass classification methods. *Journal of Machine Learning Research*, 8:1007–1025, 2007.
  - Aad W Van Der Vaart and Jon A Wellner. Weak convergence. In *Weak convergence and empirical processes: with applications to statistics*, pp. 16–28. Springer, 1996.
  - Halbert White. Maximum likelihood estimation of misspecified models. *Econometrica: Journal of the econometric society*, pp. 1–25, 1982.
  - Danny Wood, Tingting Mu, Andrew M Webb, Henry WJ Reeve, Mikel Lujan, and Gavin Brown. A unified theory of diversity in ensemble learning. *Journal of Machine Learning Research.*, 24(359): 1–49, 2023.
  - Shaofeng Zhang, Meng Liu, and Junchi Yan. The diversified ensemble neural network. *Advances in Neural Information Processing Systems.*, 33:16001–16011, 2020.
  - Xinyu Zhang. *Model averaging and its applications*. PhD thesis, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, 2010.
  - Xinyu Zhang, Dalei Yu, Guohua Zou, and Hua Liang. Optimal model averaging estimation for generalized linear models and generalized linear mixed-effects models. *Journal of the American Statistical Association.*, 111(516):1775–1790, 2016.

### A LARGE LANGUAGE MODEL USAGE DISCLOSURE

In this work, we made limited use of a large language model (LLM) as an auxiliary tool. In particular:

**Language polishing**: We used ChatGPT-5 to improve the readability, grammar, and fluency of the English text. The authors reviewed all edits and manually adjusted phrasing as needed.

**Code assistance**: We asked ChatGPT-5 to assist in generating boilerplate code for data preprocessing, but in a minimal and constrained way; the authors carefully verified, tested, modified, and adapted all generated code to ensure correctness.

We emphasize that all content in the submission is attributed to the authors. We take full responsibility for the correctness of all claims and any content originally generated by the LLM that contained errors or inconsistencies that were revised or removed. We confirm that the LLM was not included as an author, and no portion of the submission is entirely generated without human oversight.

#### B PROOF OF THEOREMS

To facilitate the proof of the theorem, we begin by stating a useful lemma.

#### B.1 LEMMA 1

Lemma 1 (Lemma 1 in Zhang (2010)). Let

$$\widehat{\boldsymbol{w}} = \arg\min_{\boldsymbol{w} \in \mathcal{W}} \{ R(\boldsymbol{w}) + a_n(\boldsymbol{w}) + b_n \}.$$

If

$$\sup_{\boldsymbol{w}\in\mathcal{W}}\frac{|a_n(\boldsymbol{w})|}{R^*(\boldsymbol{w})}=o_p(1)$$

and

$$\sup_{\boldsymbol{w}\in W}\frac{|R(\boldsymbol{w})-R^*(\boldsymbol{w})|}{R^*(\boldsymbol{w})}=o_p(1),$$

and there exists a positive constant c so that  $\lim_{n\to\infty}\inf_{w\in\mathcal{W}}R^*(w)\geq c$  almost surely, then we have

$$\frac{R(\widehat{\boldsymbol{w}})}{\inf_{\boldsymbol{w}\in\mathcal{W}}R(\boldsymbol{w})}\to 1$$

in probability.

#### B.2 Proof of Theorem 1

The weight choice criterion can be decomposed as

$$\begin{split} & \mathcal{L}(\boldsymbol{w}) \\ &= \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{Y_i - \widehat{f}(\boldsymbol{X}_i; \boldsymbol{w})\}^2 \\ &= \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{Y_i - f_0(\boldsymbol{X}_i) + f_0(\boldsymbol{X}_i) - \widehat{f}(\boldsymbol{X}_i; \boldsymbol{w})\}^2 \\ &= \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{Y_i - f_0(\boldsymbol{X}_i)\}^2 + \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{f_0(\boldsymbol{X}_i) - \widehat{f}(\boldsymbol{X}_i; \boldsymbol{w})\}^2 \\ &+ \frac{2}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{Y_i - f_0(\boldsymbol{X}_i)\} \{f_0(\boldsymbol{X}_i) - \widehat{f}(\boldsymbol{X}_i; \boldsymbol{w})\}. \end{split}$$

We first analyze  $1/n_{\text{val}} \sum_{i=1}^{n_{\text{val}}} \{Y_i - f_0(\boldsymbol{X}_i)\} \{f_0(\boldsymbol{X}_i) - \widehat{f}(\boldsymbol{X}_i; \boldsymbol{w})\}$ . Let  $\mathcal{G}_r = \{g_{\boldsymbol{w}}(\boldsymbol{X}) = f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w}) : \boldsymbol{w} \in \mathcal{W}, \|f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w})\|_{L^2} \le r\}$ , and let  $\mathcal{D} = \{\mathcal{D}_{\text{train}}, \mathcal{D}_{\text{validation}}\}$  collect all

observed samples. By the multiplier inequality (Van Der Vaart & Wellner, 1996; Bartlett et al., 2005), we have

$$\sup_{g \in \mathcal{G}_r} \left[ \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} (Y_i - f_0(\boldsymbol{X}_i)) g(\boldsymbol{X}_i) - \mathbb{E}\{(Y - f_0(\boldsymbol{X})) g(\boldsymbol{X}) | \mathcal{D}\} \right] = O_p \left( \sigma \mathbb{E} \mathcal{R}_{n_{\text{val}}} \mathcal{G}_r + \sigma \frac{1}{\sqrt{n_{\text{val}}}} r \right), \quad (1)$$

where  $\mathcal{R}_{n_{\mathrm{val}}}\mathcal{G}_r$  is the Rademacher complexity of  $\mathcal{G}_r$ , and  $\sigma$  is the sub-Gaussian parameter of the noise  $\varepsilon=Y-f_0(X)$ . Given that  $\mathcal{W}=\{\boldsymbol{w}\in[0,1]^M,\sum_{m=1}^Mw_m=1\}$ , we have  $\mathcal{R}_{n_{\mathrm{val}}}\mathcal{G}_r\leq r\sqrt{2\log(M)/n}$ . Since M,  $\sigma$  are finite, and  $n_{\mathrm{val}}$  has the same order as the total sample size n, the first term on the right hand side of (1) is  $O_p(\sigma\mathbb{E}\mathcal{R}_n\mathcal{G}_r)=O_p(r/\sqrt{n})$ , and the second term is  $O_p(\sigma\sqrt{\log(n_{\mathrm{val}})/n_{\mathrm{val}}}r)=O_p(r/\sqrt{n})$ . Then (1) becomes

$$\sup_{g \in \mathcal{G}_r} \left[ \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} (Y_i - f_0(\boldsymbol{X}_i)) g(\boldsymbol{X}_i) - \mathbb{E}\{(Y - f_0(\boldsymbol{X})) g(\boldsymbol{X}) | \mathcal{D}\} \right] = O_p\left(\frac{r}{\sqrt{n}}\right). \tag{2}$$

Let  $w_0$  be the one-hot vector with entry 1 at the position corresponding to the model with the fastest convergence rate and 0 elsewhere. Then it is straightforward to show  $||f_0(\boldsymbol{X}) - \hat{f}(\boldsymbol{X}; \boldsymbol{w}_0)||_{L^2}^2 = ||f_0(\boldsymbol{X}) - \tilde{f}(\boldsymbol{X})||_{L^2}^2 = O_p(S^2)$ . Moreover, since  $\hat{f}$  only depends on  $\mathcal{D}$ , and  $\boldsymbol{X}$  is an independent sample drawn from the same distribution but independent of  $\mathcal{D}$ , we have

$$\mathbb{E}\{(Y - f_0(\boldsymbol{X}))(f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w}_0))|\mathcal{D}\}$$

$$= \mathbb{E}\left[\mathbb{E}\left\{(Y - f_0(\boldsymbol{X}))(f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w}_0))|\boldsymbol{X}, \mathcal{D}\right\}|\mathcal{D}\right]$$

$$= \mathbb{E}\left\{\mathbb{E}(Y - f_0(\boldsymbol{X})|\boldsymbol{X}, \mathcal{D})(f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w}_0))|\mathcal{D}\right\}$$

$$= \mathbb{E}\{\mathbb{E}(\varepsilon|\boldsymbol{X})(f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w}_0))|\mathcal{D}\}$$

$$= 0. \tag{3}$$

Taking r = S, we have  $g_{w_0} \in \mathcal{G}_S$  and thus by (2) and (3), we have

$$\frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{Y_i - f_0(\boldsymbol{X}_i)\} \{f_0(\boldsymbol{X}_i) - \widehat{f}(\boldsymbol{X}_i; \boldsymbol{w}_0)\} 
= \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} (Y_i - f_0(\boldsymbol{X}_i)) (f_0(\boldsymbol{X}_i) - \widehat{f}(\boldsymbol{X}_i; \boldsymbol{w}_0)) - \mathbb{E}\{(Y - f_0(\boldsymbol{X})) (f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w}_0)) | \mathcal{D}\} 
+ \mathbb{E}\{(Y - f_0(\boldsymbol{X})) (f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w}_0)) | \mathcal{D}\} 
\leq \sup_{g \in \mathcal{G}_r} \left[ \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} (Y_i - f_0(\boldsymbol{X}_i)) g(\boldsymbol{X}_i) - \mathbb{E}\{(Y - f_0(\boldsymbol{X})) g(\boldsymbol{X}) | \mathcal{D}\}\right] 
= O_p(\sqrt{\frac{1}{n}} S),$$
(4)

Note that (3) remains valid when  $w_0$  is replaced by  $\widehat{w}$ , because  $\widehat{w}$  is entirely determined by  $\mathcal{D}$ , and hence is independent of the new sample X. Taking  $r = ||f_0(X) - \widehat{f}(X; \widehat{w})||_{L^2}$ , we have

$$\frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{Y_{i} - f_{0}(\boldsymbol{X}_{i})\} \{f_{0}(\boldsymbol{X}_{i}) - \widehat{f}(\boldsymbol{X}_{i}; \widehat{\boldsymbol{w}})\} 
= \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} (Y_{i} - f_{0}(\boldsymbol{X}_{i})) (f_{0}(\boldsymbol{X}_{i}) - \widehat{f}(\boldsymbol{X}_{i}; \widehat{\boldsymbol{w}})) - \mathbb{E}\{(Y - f_{0}(\boldsymbol{X})) (f_{0}(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \widehat{\boldsymbol{w}})) | \mathcal{D}\} 
+ \mathbb{E}\{(Y - f_{0}(\boldsymbol{X})) (f_{0}(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \widehat{\boldsymbol{w}})) | \mathcal{D}\} 
\leq \sup_{g \in \mathcal{G}_{r}} \left[ \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} (Y_{i} - f_{0}(\boldsymbol{X}_{i})) g(\boldsymbol{X}_{i}) - \mathbb{E}\{(Y - f_{0}(\boldsymbol{X})) g(\boldsymbol{X}) | \mathcal{D}\} \right] 
= O_{p}(\sqrt{\frac{1}{n}} \|f_{0}(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \widehat{\boldsymbol{w}})\|_{L^{2}}).$$
(5)

Then we analyze  $1/n_{\text{val}} \sum_{i=1}^{n_{\text{val}}} \{f_0(\boldsymbol{X}_i) - \widehat{f}(\boldsymbol{X}_i; \boldsymbol{w})\}^2$ . Let  $\mathcal{H}_r = \{h_{\boldsymbol{w}} = \{f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w})\}^2 : \boldsymbol{w} \in \mathcal{W}, \text{Var}[\{f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w})\}^2 | \mathcal{D}] \leq r^2\}$ . Similar to (1), we obtain

$$\sup_{h \in \mathcal{H}_r} \left[ \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} h(\boldsymbol{X}_i)^2 - \mathbb{E}\left\{ h(\boldsymbol{X})^2 | \mathcal{D} \right\} \right] = O_p\left(\mathbb{E}\mathcal{R}_n \mathcal{H}_r + \sqrt{\frac{1}{n_{\text{val}}}} r \right) = O_p\left(\sqrt{\frac{1}{n}} r \right).$$
(6)

Since  $f_0$  and  $\widehat{f}_m$  are uniformly bounded, there exists a positive constant C such that

$$\operatorname{Var}[\{f_{0}(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w})\}^{2} | \mathcal{D}]$$

$$\leq \mathbb{E}[\{f_{0}(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w})\}^{4} | \mathcal{D}]$$

$$\leq \|f_{0}(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w})\|_{L^{2}}^{2} \|f_{0}(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w})\|_{\infty}^{2}$$

$$\leq C^{2} \|f_{0}(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w})\|_{L^{2}}^{2}.$$
(7)

When taking  $\boldsymbol{w} = \boldsymbol{w}_0$  and r = CS in  $\mathcal{H}_r$ , we have  $h_{\boldsymbol{w}_0} = \{f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w}_0)\}^2 \in \mathcal{H}_{CS}$  because  $\operatorname{Var}[\{f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w}_0)\}^2 | \mathcal{D}] \leq C^2 \|f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \boldsymbol{w}_0)\|_{L^2}^2 = C^2 S^2$  from (7). By (6), it follows that

$$\frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{f_0(\boldsymbol{X}_i) - \hat{f}(\boldsymbol{X}_i; \boldsymbol{w}_0)\}^2$$

$$= \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{f_0(\boldsymbol{X}_i) - \hat{f}(\boldsymbol{X}_i; \boldsymbol{w}_0)\}^2 - \mathbb{E}\left[\{f_0(\boldsymbol{X}) - \hat{f}(\boldsymbol{X}; \boldsymbol{w}_0)\}^2 | \mathcal{D}\right] + \mathbb{E}\left[\{f_0(\boldsymbol{X}) - \hat{f}(\boldsymbol{X}; \boldsymbol{w}_0)\}^2 | \mathcal{D}\right]$$

$$\leq \sup_{h \in \mathcal{H}_{CS}} \left[ \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} h(\boldsymbol{X}_i)^2 - \mathbb{E}\left\{h(\boldsymbol{X})^2 | \mathcal{D}\right\}\right] + \|f_0(\boldsymbol{X}) - \hat{f}(\boldsymbol{X}; \boldsymbol{w}_0)\|_{L^2}^2$$

$$= O_p(\sqrt{\frac{1}{n}}S + S^2). \tag{8}$$

Similarly, taking  $\mathbf{w} = \widehat{\mathbf{w}}$  and  $r_{\widehat{\mathbf{w}}} = C \|f_0(\mathbf{X}) - \widehat{f}(\mathbf{X}; \widehat{\mathbf{w}})\|_{L^2}$ , we have  $h_{\widehat{\mathbf{w}}} = \{f_0(\mathbf{X}) - \widehat{f}(\mathbf{X}; \widehat{\mathbf{w}})\}^2 \in \mathcal{H}_{r_{\widehat{\mathbf{w}}}}$  and thus by (6), the following bound holds:

$$\frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{f_0(\boldsymbol{X}_i) - \widehat{f}(\boldsymbol{X}_i; \widehat{\boldsymbol{w}})\}^2$$

$$= \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{f_0(\boldsymbol{X}_i) - \widehat{f}(\boldsymbol{X}_i; \widehat{\boldsymbol{w}})\}^2 - \mathbb{E}\left[\{f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \widehat{\boldsymbol{w}})\}^2 | \mathcal{D}\right] + \mathbb{E}\left[\{f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \widehat{\boldsymbol{w}})\}^2 | \mathcal{D}\right]$$

$$\leq \sup_{h \in \mathcal{H}_{r_{\widehat{\boldsymbol{w}}}}} \left[ \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} h(\boldsymbol{X}_i)^2 - \mathbb{E}\left\{h(\boldsymbol{X})^2 | \mathcal{D}\right\}\right] + \|f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \widehat{\boldsymbol{w}})\|_{L^2}^2$$

$$= O_p(\sqrt{\frac{\log(n)}{n}} \|f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \widehat{\boldsymbol{w}})\|_{L^2}) + \|f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \widehat{\boldsymbol{w}})\|_{L^2}^2. \tag{9}$$

Although  $\widehat{w}$  depends on the validation data, the bound still holds because the supremum is taken over the class  $\mathcal{H}_{r_{\widehat{w}}}$ . Combining (5) and (9),  $\mathcal{L}(\widehat{w})$  can be written as

$$\mathcal{L}(\widehat{\boldsymbol{w}}) = \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{ Y_i - f_0(\boldsymbol{X}_i) \}^2 + \| f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \widehat{\boldsymbol{w}}) \|_{L^2}^2 + O_p(\sqrt{\frac{1}{n}}) \| f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \widehat{\boldsymbol{w}}) \|_{L^2}^2, (10)$$

and according to (4) and (8),  $\mathcal{L}(w_0)$  can be written as

$$\mathcal{L}(\boldsymbol{w}_0) = \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{Y_i - f_0(\boldsymbol{X}_i)\}^2 + O_p(S^2 + \sqrt{\frac{1}{n}}S).$$
 (11)

Using the expansions in (10) and (11), together with the fact that  $\mathcal{L}(\widehat{w})$  minimizes the validation loss, i.e.,  $\mathcal{L}(\widehat{w}) \leq \mathcal{L}(w_0)$ , we have

$$||f_0(\mathbf{X}) - \widehat{f}(\mathbf{X}; \widehat{\mathbf{w}})||_{L^2}^2 + O_p(\sqrt{\frac{1}{n}})||f_0(\mathbf{X}) - \widehat{f}(\mathbf{X}; \widehat{\mathbf{w}})||_{L^2} = O_p(S^2 + \sqrt{\frac{1}{n}}S).$$
(12)

By completing the square, (12) can be written as

$$\left\Vert f_{0}(oldsymbol{X}% )-f_{0}(oldsymbol{X})
ight\Vert d^{2}(oldsymbol{X}) \left\Vert f_{0}(oldsymbol{X})
ight\Vert d^{2}(oldsymbol{X})$$

 $\left[ \|f_0(\boldsymbol{X}) - \widehat{f}(\boldsymbol{X}; \widehat{\boldsymbol{w}})\|_{L^2} + O_p(\sqrt{\frac{1}{n}}) \right]^2 = \{O_p(S + \sqrt{\frac{1}{n}})\}^2,$ 

i.e.,

$$||f_0(\boldsymbol{X})||$$

 $||f_0(\mathbf{X}) - \widehat{f}(\mathbf{X}; \widehat{\mathbf{w}})||_{L^2} = O_p(S + \sqrt{\frac{1}{n}}).$ 

This completes the proof of Theorem 1.

#### B.3 PROOF OF THEOREM 2

The regression and multiclass-classification problems differ only in the choice of loss function. In classification problems, the performance measure of primary interest is the misclassification error (the 0–1 loss). Because this loss is discontinuous, it is typically replaced by a continuous surrogate, most commonly the cross-entropy loss. The surrogate is smooth and differentiable, which facilitates gradient-based optimization. Importantly, since cross-entropy is a calibrated surrogate, its excess risk dominates the squared excess misclassification risk; see Tewari & Bartlett (2007) for example. Let  $w_0$  be the one-hot vector that places 1 on the coordinate corresponding to the model with the fastest convergence rate and 0 elsewhere. Thus, if the excess misclassification risk of  $w_0$  is of order S, the corresponding cross-entropy excess risk satisfies

$$\mathbb{E}\Big\{f_0(\boldsymbol{X})\log\frac{f_0(\boldsymbol{X})}{\widehat{f}(\boldsymbol{X};\boldsymbol{w}_0)}|\mathcal{D}\Big\} = O_p(S^2).$$

Let

$$\mathcal{H}_r = \left\{ h_{\boldsymbol{w}}(\boldsymbol{X}) = \left\{ \log f_0(\boldsymbol{X}) - \log \widehat{f}(\boldsymbol{X}; \boldsymbol{w}) \right\} : \ w \in \mathcal{W}, \ \operatorname{Var}\{h_{\boldsymbol{w}}(\boldsymbol{X}) | \mathcal{D}\} \le r^2 \right\}.$$

By the multiplier inequality,

$$\sup_{h \in \mathcal{H}_r} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{ Y_i - f_0(\boldsymbol{X}_i) \} h(\boldsymbol{X}_i) - \mathbb{E}[\{ Y - f_0(\boldsymbol{X}) \} h(\boldsymbol{X}) | \mathcal{D}] \right| = O_p \left( \gamma \mathbb{E} \mathcal{R}_{n_{\text{val}}} \mathcal{H}_r + \gamma \sqrt{\frac{1}{n}} r \right), (13)$$

where  $\gamma^2 = \text{Var}(Y - f_0(X))$  is finite because  $\text{Var}(Y - f_0(X)) \le 1/4$  in classification tasks. Moreover, we have  $\mathbb{E}\mathcal{R}_{n_{\text{val}}}\mathcal{H}_r \leq Cr\sqrt{\log(M)/n_{\text{val}}}$ , and  $n_{\text{val}}$  has the same order as n, then (13) becomes

$$\sup_{h \in \mathcal{H}_r} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{ Y_i - f_0(\boldsymbol{X}_i) \} h(\boldsymbol{X}_i) - \mathbb{E}[\{ Y - f_0(\boldsymbol{X}) \} h(\boldsymbol{X}) | \mathcal{D}] \right| = O_p\left(\sqrt{\frac{1}{n}} r\right). \tag{14}$$

Similarly,

$$\sup_{h \in \mathcal{H}_r} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} f_0(\boldsymbol{X}_i) h(\boldsymbol{X}_i) - \mathbb{E}[f_0(\boldsymbol{X}) h(\boldsymbol{X}) | \mathcal{D}] \right| = O_p\left(\sqrt{\frac{1}{n}} r\right). \tag{15}$$

Moreover, we have

$$\mathbb{E}[\{Y - f_0(\boldsymbol{X})\}h(\boldsymbol{X})|\mathcal{D}] = \mathbb{E}[\mathbb{E}\{Y - f_0(\boldsymbol{X})|\boldsymbol{X}\}h(\boldsymbol{X})|\mathcal{D}] = 0.$$
(16)

By (14)-(16), we have

$$\frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} Y_i h(\boldsymbol{X}_i) 
= \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{Y_i - f_0(\boldsymbol{X}_i)\} h(\boldsymbol{X}_i) + \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} f_0(\boldsymbol{X}_i) h(\boldsymbol{X}_i) 
\leq \sup_{h \in \mathcal{H}_r} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{Y_i - f_0(\boldsymbol{X}_i)\} h(\boldsymbol{X}_i) - \mathbb{E}[\{Y - f_0(\boldsymbol{X})\} h(\boldsymbol{X})|\mathcal{D}] \right| 
+ \sup_{h \in \mathcal{H}_r} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} f_0(\boldsymbol{X}_i) h(\boldsymbol{X}_i) - \mathbb{E}[f_0(\boldsymbol{X}) h(\boldsymbol{X})|\mathcal{D}] \right| 
+ \mathbb{E}[\{Y - f_0(\boldsymbol{X})\} h(\boldsymbol{X})|\mathcal{D}] + \mathbb{E}[f_0(\boldsymbol{X}) h(\boldsymbol{X})|\mathcal{D}]$$

$$= \mathbb{E}[f_0(\boldsymbol{X}) h(\boldsymbol{X})|\mathcal{D}] + O_p\left(\sqrt{\frac{1}{n}}r\right). \tag{17}$$

Taking r = S and  $h(\mathbf{X}) = \log(f_0(\mathbf{X})/f_0(\mathbf{X}; \mathbf{w}_0))$ , (17) becomes

$$\mathcal{L}(\boldsymbol{w}_{0}) = \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} Y_{i} \log(\widehat{f}(\boldsymbol{X}_{i}; \boldsymbol{w}_{0}))$$

$$= \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} Y_{i} \log \frac{f_{0}(\boldsymbol{X}_{i})}{\widehat{f}(\boldsymbol{X}_{i}; \boldsymbol{w}_{0})} - \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} Y_{i} \log(f_{0}(\boldsymbol{X}_{i}))$$

$$= \left| \mathbb{E}\{f_{0}(\boldsymbol{X}) \log \frac{f_{0}(\boldsymbol{X})}{\widehat{f}(\boldsymbol{X}; \boldsymbol{w}_{0})} | \mathcal{D}\} \right| + O_{p}\left(\sqrt{\frac{1}{n}}S\right) - \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} Y_{i} \log(f_{0}(\boldsymbol{X}_{i}))$$

$$= O_{p}\left(S^{2} + \sqrt{\frac{1}{n}}S\right) - \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} Y_{i} \log(f_{0}(\boldsymbol{X}_{i}))$$

Similarly, we have

$$\mathcal{L}(\widehat{\boldsymbol{w}}) = \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} Y_i \log(\widehat{f}(\boldsymbol{X}_i; \widehat{\boldsymbol{w}}))$$

$$= \mathbb{E}\{f_0(\boldsymbol{X}) \log \frac{f_0(\boldsymbol{X})}{\widehat{f}(\boldsymbol{X}; \widehat{\boldsymbol{w}})} | \mathcal{D}\} + O_p(\sqrt{\frac{\log(n)}{n}}) \sqrt{\mathbb{E}\{f_0(\boldsymbol{X}) \log \frac{f_0(\boldsymbol{X})}{\widehat{f}(\boldsymbol{X}; \widehat{\boldsymbol{w}})} | \mathcal{D}\}}$$

$$-\frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} Y_i \log(f_0(\boldsymbol{X}_i)).$$

Since  $\mathcal{L}(\hat{\boldsymbol{w}}) \leq \mathcal{L}(\boldsymbol{w}_0)$ , This implies

$$\mathbb{E}\{f_0(\boldsymbol{X})\log\frac{f_0(\boldsymbol{X})}{\widehat{f}(\boldsymbol{X};\widehat{\boldsymbol{w}})}|\mathcal{D}\} + O_p(\sqrt{\frac{1}{n}})\sqrt{\mathbb{E}\{f_0(\boldsymbol{X})\log\frac{f_0(\boldsymbol{X})}{\widehat{f}(\boldsymbol{X};\widehat{\boldsymbol{w}})}|\mathcal{D}\}} = O_p(S^2 + \sqrt{\frac{1}{n}}S),$$

and thus  $\sqrt{\mathbb{E}\{f_0(\boldsymbol{X})\log\{f_0(\boldsymbol{X})/\widehat{f}(\boldsymbol{X};\widehat{\boldsymbol{w}})\}|\mathcal{D}\}} = O_p(\sqrt{\frac{1}{n}} + S)$ . This further implies that the misclassification rate  $R_{0/1}(\widehat{\boldsymbol{w}}) = O_p(S + \sqrt{\frac{1}{n}})$ .

#### **B.4** Proof of Corollaries

For the MLP case, suppose the true regression function  $f_0$  lies in the Hölder class  $\mathcal{H}^{\beta}([0,1]^d, B_0)$ , defined as

$$\mathcal{H}^{\beta}([0,1]^d, B_0) = \Big\{ f : [0,1]^d \to \mathbb{R} : \max_{\|\boldsymbol{\alpha}\|_1 \le s} \|\partial^{\boldsymbol{\alpha}} f\|_{\infty} \le B_0, \max_{\|\boldsymbol{\alpha}\|_1 = s} \sup_{x \ne y} \frac{|\partial^{\boldsymbol{\alpha}} f(x) - \partial^{\boldsymbol{\alpha}} f(y)|}{\|x - y\|^r} \le B_0 \Big\},$$

where  $\partial^{\alpha} = \partial^{\alpha_1} \cdots \partial^{\alpha_d}$  with  $\alpha = (\alpha_1, \dots, \alpha_d)^{\top} \in \mathbb{N}_0^d$  and  $\|\alpha\|_1 = \sum_{i=1}^d \alpha_i$ . We consider candidate models from the ReLU MLP class

$$\mathcal{F}_{B_0,\mathcal{W},\mathcal{D},\mathcal{S},d} = \{f : \mathbb{R}^d \to \mathbb{R} : \text{width } \mathcal{W}, \text{ depth } \mathcal{D}, \text{ total parameters } \mathcal{S}, \|f\|_{\infty} \leq B_0 \}.$$
 (18)

Corollary 5.3 in Jiao et al. (2023) guarantees that for the function class of ReLU MLPs with

$$\mathcal{W}_0 = O\left(n^{\frac{d}{4(d+2\beta)}}\log_2(n)\right), \quad \mathcal{D}_0 = O\left(n^{\frac{d}{4(d+2\beta)}}\log_2(n)\right), \quad \mathcal{S}_0 = O\left(n^{\frac{3d}{4(d+2\beta)}}(\log n)^4\right),$$

the empirical risk minimizer  $\widehat{f}=\arg\min_{f\in\mathcal{F}_{B_0,\mathcal{W}_0,\mathcal{D}_0,\mathcal{S}_0,d}}$  satisfies

$$\mathbb{E}\|\widehat{f}(\boldsymbol{X}) - f_0(\boldsymbol{X})\|_{L^2}^2 \le \mathcal{B}_0^5(\lfloor \beta \rfloor + 1)^4 d^{2\lfloor \beta \rfloor + \beta \vee 1} n^{-2\beta/(d+2\beta)} (\log n)^{11} = \widetilde{O}(n^{-2\beta/(d+2\beta)}).$$

Substituting  $S = n^{-\beta/(d+2\beta)}$  into Theorem 1 implies that, when the candidate models are MLPs and include one with width  $W_0$ , depth  $\mathcal{D}_0$ , and size  $\mathcal{S}_0$ , the asymptotic error bound of the model averaging estimator is  $\widetilde{O}_p(n^{-2\beta/(d+2\beta)})$ .

The other two corollaries can be proved similarly, and the proof is omitted here.

#### B.5 PROOF OF THEOREM 3

Let

$$\widetilde{\mathcal{L}}(\boldsymbol{w}) = \mathcal{L}(\boldsymbol{w}) - \frac{1}{n_{\mathrm{val}}} \sum_{i=1}^{n_{\mathrm{val}}} (Y_i^2 - f_0(\boldsymbol{X}_i)^2).$$

Observe that the newly added component is unrelated to w, so we have

$$\widehat{\boldsymbol{w}} = \operatorname*{arg\,min}_{\boldsymbol{w} \in \mathcal{W}} \mathcal{L}(\boldsymbol{w}) = \operatorname*{arg\,min}_{\boldsymbol{w} \in \mathcal{W}} \{ R(\boldsymbol{w}) + \widetilde{\mathcal{L}}(\boldsymbol{w}) - R(\boldsymbol{w}) \}.$$

According to Lemma 1, to prove Theorem 3, it is sufficient to prove

$$\sup_{\boldsymbol{w}\in\mathcal{W}} \frac{|R(\boldsymbol{w}) - R^*(\boldsymbol{w})|}{R^*(\boldsymbol{w})} = o_p(1)$$
(19)

and

$$\sup_{\boldsymbol{w}\in\mathcal{W}} \frac{|\widetilde{\mathcal{L}}(\boldsymbol{w}) - R(\boldsymbol{w})|}{R^*(\boldsymbol{w})} = o_p(1). \tag{20}$$

For (19), we have

$$\sup_{\boldsymbol{w} \in \mathcal{W}} \frac{|R(\boldsymbol{w}) - R^*(\boldsymbol{w})|}{R^*(\boldsymbol{w})} \\
\leq \frac{\sup_{\boldsymbol{w} \in \mathcal{W}} |R(\boldsymbol{w}) - R^*(\boldsymbol{w})|}{\inf_{\boldsymbol{w} \in \mathcal{W}} R^*(\boldsymbol{w})} \\
= \xi_n^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} |\|\widehat{f}(\boldsymbol{X}; \boldsymbol{w}) - f_0(\boldsymbol{X})\|_{L^2}^2 - \|f^*(\boldsymbol{X}; \boldsymbol{w}) - f_0(\boldsymbol{X})\|_{L^2}^2| \\
\leq \xi_n^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} |(\|\widehat{f}(\boldsymbol{X}; \boldsymbol{w}) - f_0(\boldsymbol{X})\|_{L^2} + \|f^*(\boldsymbol{X}; \boldsymbol{w}) - f_0(\boldsymbol{X})\|_{L^2}) \|\widehat{f}(\boldsymbol{X}; \boldsymbol{w}) - f^*(\boldsymbol{X}; \boldsymbol{w})\|_{L^2}| \\
\leq C\xi_n^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} \|\widehat{f}(\boldsymbol{X}; \boldsymbol{w}) - f^*(\boldsymbol{X}; \boldsymbol{w})\|_{L^2} \\
= O_p(\xi_n^{-1} \phi_n) \\
= o_p(1),$$

where the last step comes from Condition 2. Thus (19) is obtained. Next, we will prove (20):

$$\sup_{\boldsymbol{w} \in \mathcal{W}} \frac{|\widetilde{\mathcal{L}}(\boldsymbol{w}) - R(\boldsymbol{w})|}{R^{*}(\boldsymbol{w})}$$

$$\leq \xi_{n}^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} |\widetilde{\mathcal{L}}(\boldsymbol{w}) - R(\boldsymbol{w})|$$

$$= \xi_{n}^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \left\{ y_{i} - \widehat{f}(\boldsymbol{x}_{i}; \boldsymbol{w}) \right\}^{2} - \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \left\{ y_{i}^{2} - f_{0}(\boldsymbol{x}_{i})^{2} \right\} - \|\widehat{f}(\boldsymbol{X}; \boldsymbol{w}) - f_{0}(\boldsymbol{X})\|_{L^{2}}^{2} \right|$$

$$\leq \xi_{n}^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \left\{ \widehat{f}(\boldsymbol{x}_{i}; \boldsymbol{w}) - f_{0}(\boldsymbol{X}_{i}) \right\}^{2} - \|\widehat{f}(\boldsymbol{X}; \boldsymbol{w}) - f_{0}(\boldsymbol{X})\|_{L^{2}}^{2} \right|$$

$$+ \xi_{n}^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \left\{ \widehat{f}(\boldsymbol{X}_{i}; \boldsymbol{w}) - f^{*}(\boldsymbol{X}_{i}; \boldsymbol{w}) \right\} \left\{ Y_{i} - f_{0}(\boldsymbol{X}_{i}) \right\} \right|.$$

$$+ \xi_{n}^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} f^{*}(\boldsymbol{X}_{i}; \boldsymbol{w}) \left\{ Y_{i} - f_{0}(\boldsymbol{X}_{i}) \right\} \right|.$$

$$(21)$$

To prove (20), it is sufficient to prove the next three equations

$$\xi_n^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \left\{ \widehat{f}(\boldsymbol{x}_i; \boldsymbol{w}) - f_0(\boldsymbol{X}_i) \right\}^2 - \|\widehat{f}(\boldsymbol{X}; \boldsymbol{w}) - f_0(\boldsymbol{X})\|_{L^2}^2 \right| = o_p(1), \quad (22)$$

$$\xi_n^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{ \widehat{f}(\boldsymbol{X}_i; \boldsymbol{w}) - f^*(\boldsymbol{X}_i; \boldsymbol{w}) \} \{ Y_i - f_0(\boldsymbol{X}_i) \} \right| = o_p(1), \tag{23}$$

and

$$\xi_n^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} f^*(\boldsymbol{X}_i; \boldsymbol{w}) \{ Y_i - f_0(\boldsymbol{X}_i) \} \right| = o_p(1).$$
 (24)

Let  $\mathcal{H}_r = \{h_{\boldsymbol{w}} = \{\widehat{f}(\boldsymbol{X}; \boldsymbol{w}) - f_0(\boldsymbol{X})\}^2 : \boldsymbol{w} \in \mathcal{W}\}$ . By the multiplier inequality, we have

$$\sup_{\boldsymbol{w} \in \mathcal{W}} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \left\{ \widehat{f}(\boldsymbol{x}_i; \boldsymbol{w}) - f_0(\boldsymbol{X}_i) \right\}^2 - \|\widehat{f}(\boldsymbol{X}; \boldsymbol{w}) - f_0(\boldsymbol{X})\|_{L^2}^2 \right|$$

$$= \sup_{h \in \mathcal{H}_r} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} h(\boldsymbol{X}_i) - \mathbb{E}\{h(\boldsymbol{X})|\mathcal{D}\} \right|$$

$$= O_p\left(\sqrt{\frac{1}{n}}\right), \tag{25}$$

The difference between (6) and (25) is that we set r=1 here. This is because we take the supremum over all  $w \in \mathcal{W}$ , instead of localizing to a particular neighborhood of w. And the bound in (25) is also bigger than that in (6).

Therefore, by Condition 2, (22) can be proved by

$$\xi_n^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \left\{ \widehat{f}(\boldsymbol{x}_i; \boldsymbol{w}) - f_0(\boldsymbol{X}_i) \right\}^2 - \|\widehat{f}(\boldsymbol{X}; \boldsymbol{w}) - f_0(\boldsymbol{X})\|_{L^2}^2 \right| = O_p(\xi_n^{-1} / \sqrt{n}) = o_p(1).$$
 (26)

For (23), we have

$$\xi_{n}^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \{ \widehat{f}(\boldsymbol{X}_{i}; \boldsymbol{w}) - f^{*}(\boldsymbol{X}_{i}; \boldsymbol{w}) \} \{ Y_{i} - f_{0}(\boldsymbol{X}_{i}) \} \right| \\
= \xi_{n}^{-1} (P_{n} - P + P) \{ \widehat{f}(\boldsymbol{X}; \boldsymbol{w}) - f^{*}(\boldsymbol{X}; \boldsymbol{w}) \} \{ Y - f_{0}(\boldsymbol{X}) \} \\
= \xi_{n}^{-1} O_{p}(\sqrt{1/n}) \\
= o_{p}(1). \tag{27}$$

For (24), according to the Chebyshev's inequality, we have

$$\Pr\left\{\xi_{n}^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} f^{*}(\boldsymbol{X}_{i}; \boldsymbol{w}) \varepsilon_{i} \right| > \nu \right\}$$

$$= \Pr\left\{\xi_{n}^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} \sum_{m=1}^{M} w_{m} f_{m}^{*}(\boldsymbol{X}_{i}) \varepsilon_{i} \right| > \nu \right\}$$

$$\leq \sum_{m=1}^{M} \Pr\left\{\xi_{n}^{-1} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} f_{m}^{*}(\boldsymbol{X}_{i}) \varepsilon_{i} \right| > \nu \right\}$$

$$\leq \sum_{m=1}^{M} \xi_{n}^{-2} \nu^{-2} \operatorname{Var} \left\{ \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} f_{m}^{*}(\boldsymbol{X}_{i}) \varepsilon_{i} \right\}$$

$$= \xi_{n}^{-2} \nu^{-2} n^{-1} M \operatorname{Var} \left\{f_{m}^{*}(\boldsymbol{X}) \varepsilon\right\}$$

$$= C\xi_{n}^{-2} \nu^{-2} n^{-1}, \tag{28}$$

which means  $\xi_n^{-1} \sup_{\boldsymbol{w} \in \mathcal{W}} \left| \frac{1}{n_{\text{val}}} \sum_{i=1}^{n_{\text{val}}} f^*(\boldsymbol{X}_i; \boldsymbol{w}) \varepsilon_i \right| = o_p(1)$  by Condition 2.

Equations (26)-(28) imply that (22)-(24) hold, and thus we obtain (20). Since (19) and (20) hold, we complete the proof of Theorem 3.

#### C MORE EXPERIMENTAL DETAILS

**Training resouces.** We use the A800 80G for training the models with PyTorch version 2.5.1.

#### D ADDITIONAL RESULTS

We report the ratio of WDE MSE to oracle MSE across parameter discrepancies and sample sizes in structural misspecification in Table 5.

Table 5: Ratio of WDE MSE to oracle MSE across parameter discrepancies and sample sizes

		$\Delta = 30000$			$\Delta = 50000$			$\Delta = 100000$		
$\alpha$		N			N			N		
a	5000	10000	15000	5000	10000	15000	5000	10000	15000	
0.9	1.0066	1.0038	1.0020	1.0061	1.0033	1.0013	1.0076	1.0032	1.0031	
0.8	1.0062	1.0031	1.0022	1.0046	1.0022	1.0020	1.0054	1.0034	1.0025	
0.7	1.0045	1.0030	1.0023	1.0078	1.0041	1.0021	1.0061	1.0034	1.0021	
0.6	1.0056	1.0020	1.0025	1.0087	1.0034	1.0018	1.0068	1.0038	1.0022	
0.5	1.0059	1.0022	1.0014	1.0055	1.0036	1.0025	1.0040	1.0030	1.0014	
0.4	1.0061	1.0028	1.0034	1.0081	1.0027	1.0025	1.0051	1.0041	1.0010	
0.3	1.0058	1.0031	1.0022	1.0057	1.0027	1.0019	1.0046	1.0038	1.0036	
0.2	1.0091	1.0037	1.0024	1.0078	1.0027	1.0018	1.0066	1.0028	1.0014	
0.1	1.0088	1.0029	1.0017	1.0062	1.0028	1.0015	1.0058	1.0049	1.0021	
0.0	1.0069	1.0020	1.0015	1.0036	1.0034	1.0021	1.0079	1.0045	1.0023	

We compare with other weighted methods in structural misspecification from Table 6 to Table 8.

#### E FUTURE WORK

While the current theoretical results establish asymptotic error bounds for the proposed weighted deep ensemble estimator, extending these results to non-asymptotic settings remains an important direction for future research. Such analysis could provide tighter guarantees in finite-sample regimes, which are often relevant in practical applications. Another promising avenue is to investigate principled approaches for determining the number and composition of candidate models in the ensemble. Understanding how model diversity and ensemble size affect performance could lead to more efficient and adaptive ensemble design strategies.

Table 6: Relative MSE comparison for sample size N=5,000

		WDE	In-sample Ensemble	Greedy Ensemble
$\Delta$	$1 - \alpha$	Mean ± Std	Mean $\pm$ Std	Mean $\pm$ Std
	0.1	$0.417 \pm 0.089$	$0.421 \pm 0.090$	$0.418 \pm 0.091$
	0.2	$0.408 \pm 0.081$	$0.416 \pm 0.081$	$0.413 \pm 0.083$
	0.3	$0.398 \pm 0.083$	$0.402 \pm 0.086$	$0.400 \pm 0.085$
	0.4	$0.387 \pm 0.095$	$0.394 \pm 0.097$	$0.391 \pm 0.098$
30,000	0.5	$0.397 \pm 0.146$	$0.399 \pm 0.146$	$0.399 \pm 0.146$
30,000	0.6	$0.360 \pm 0.048$	$0.366 \pm 0.052$	$0.363 \pm 0.049$
	0.7	$0.347 \pm 0.043$	$0.356 \pm 0.048$	$0.351 \pm 0.044$
	0.8	$0.341 \pm 0.034$	$0.349 \pm 0.037$	$0.343 \pm 0.036$
	0.9	$0.337 \pm 0.031$	$0.343 \pm 0.034$	$0.340 \pm 0.032$
	1.0	$0.334 \pm 0.031$	$0.342 \pm 0.030$	$0.336 \pm 0.031$
	0.1	$0.437 \pm 0.061$	$0.442 \pm 0.063$	$0.439 \pm 0.063$
	0.2	$0.476 \pm 0.192$	$0.483 \pm 0.190$	$0.479 \pm 0.191$
	0.3	$0.419 \pm 0.067$	$0.424 \pm 0.070$	$0.423 \pm 0.070$
	0.4	$0.391 \pm 0.050$	$0.396 \pm 0.051$	$0.394 \pm 0.051$
50.000	0.5	$0.399 \pm 0.091$	$0.382 \pm 0.041$	$0.402 \pm 0.103$
50,000	0.6	$0.354 \pm 0.041$	$0.361 \pm 0.041$	$0.357 \pm 0.041$
	0.7	$0.344 \pm 0.030$	$0.350 \pm 0.032$	$0.345 \pm 0.031$
	0.8	$0.339 \pm 0.032$	$0.345 \pm 0.035$	$0.340 \pm 0.033$
	0.9	$0.334 \pm 0.027$	$0.342 \pm 0.029$	$0.337 \pm 0.028$
	1.0	$0.330 \pm 0.027$	$0.335 \pm 0.027$	$0.334 \pm 0.027$
	0.1	$0.348 \pm 0.043$	$0.357 \pm 0.044$	$0.349 \pm 0.043$
	0.2	$0.349 \pm 0.042$	$0.357 \pm 0.047$	$0.352 \pm 0.044$
	0.3	$0.346 \pm 0.043$	$0.355 \pm 0.043$	$0.348 \pm 0.043$
	0.4	$0.342 \pm 0.042$	$0.348 \pm 0.043$	$0.344 \pm 0.042$
100.000	0.5	$0.340 \pm 0.039$	$0.347 \pm 0.040$	$0.342 \pm 0.040$
100,000	0.6	$0.335 \pm 0.035$	$0.343 \pm 0.038$	$0.338 \pm 0.036$
	0.7	$0.332 \pm 0.034$	$0.339 \pm 0.037$	$0.335 \pm 0.035$
	0.8	$0.331 \pm 0.026$	$0.336 \pm 0.028$	$0.331 \pm 0.026$
	0.9	$0.330 \pm 0.025$	$0.337 \pm 0.026$	$0.334 \pm 0.027$
	1.0	$0.327 \pm 0.024$	$0.336 \pm 0.021$	$0.331 \pm 0.026$

Table 7: Relative MSE comparison for sample size N=10,000

		WDE	In-sample Ensemble	Greedy Ensemble
$\Delta$	$1 - \alpha$	Mean ± Std	Mean $\pm$ Std	Mean $\pm$ Std
	0.1	$0.399 \pm 0.108$	$0.404 \pm 0.108$	$0.404 \pm 0.110$
	0.2	$0.396 \pm 0.117$	$0.400 \pm 0.116$	$0.399 \pm 0.120$
	0.3	$0.377 \pm 0.094$	$0.380 \pm 0.093$	$0.379 \pm 0.094$
	0.4	$0.369 \pm 0.069$	$0.371 \pm 0.068$	$0.371 \pm 0.069$
20,000	0.5	$0.351 \pm 0.053$	$0.355 \pm 0.051$	$0.352 \pm 0.054$
30,000	0.6	$0.339 \pm 0.044$	$0.344 \pm 0.045$	$0.342 \pm 0.047$
	0.7	$0.328 \pm 0.038$	$0.331 \pm 0.038$	$0.330 \pm 0.038$
	0.8	$0.327 \pm 0.034$	$0.330 \pm 0.034$	$0.329 \pm 0.036$
	0.9	$0.318 \pm 0.032$	$0.321 \pm 0.032$	$0.319 \pm 0.032$
	1.0	$0.320 \pm 0.034$	$0.323 \pm 0.035$	$0.323 \pm 0.036$
	0.1	$0.399 \pm 0.052$	$0.405 \pm 0.052$	$0.401 \pm 0.052$
	0.2	$0.398 \pm 0.063$	$0.401 \pm 0.065$	$0.400 \pm 0.064$
	0.3	$0.384 \pm 0.060$	$0.388 \pm 0.059$	$0.387 \pm 0.061$
	0.4	$0.396 \pm 0.105$	$0.398 \pm 0.104$	$0.397 \pm 0.105$
50,000	0.5	$0.349 \pm 0.048$	$0.353 \pm 0.048$	$0.351 \pm 0.049$
50,000	0.6	$0.334 \pm 0.039$	$0.338 \pm 0.040$	$0.336 \pm 0.040$
	0.7	$0.324 \pm 0.036$	$0.327 \pm 0.036$	$0.325 \pm 0.036$
	0.8	$0.318 \pm 0.034$	$0.321 \pm 0.033$	$0.319 \pm 0.034$
	0.9	$0.314 \pm 0.035$	$0.320 \pm 0.039$	$0.316 \pm 0.036$
	1.0	$0.311 \pm 0.029$	$0.314 \pm 0.030$	$0.313 \pm 0.029$

Table 8: Relative MSE comparison for sample size N=15,000

		WDE	In-sample Ensemble	Greedy Ensemble	
$\Delta$	$1 - \alpha$	Mean ± Std	Mean $\pm$ Std	Mean $\pm$ Std	
	0.1	$0.381 \pm 0.086$	$0.384 \pm 0.084$	$0.384 \pm 0.089$	
	0.2	$0.381 \pm 0.106$	$0.384 \pm 0.105$	$0.386 \pm 0.111$	
	0.3	$0.350 \pm 0.052$	$0.354 \pm 0.052$	$0.351 \pm 0.053$	
	0.4	$0.340 \pm 0.047$	$0.343 \pm 0.046$	$0.342 \pm 0.047$	
30,000	0.5	$0.328 \pm 0.040$	$0.333 \pm 0.039$	$0.330 \pm 0.041$	
30,000	0.6	$0.323 \pm 0.032$	$0.327 \pm 0.032$	$0.326 \pm 0.032$	
	0.7	$0.323 \pm 0.062$	$0.326 \pm 0.062$	$0.325 \pm 0.065$	
	0.8	$0.309 \pm 0.026$	$0.314 \pm 0.029$	$0.312 \pm 0.029$	
	0.9	$0.335 \pm 0.106$	$0.338 \pm 0.105$	$0.336 \pm 0.106$	
	1.0	$0.303 \pm 0.023$	$0.307 \pm 0.024$	$0.305 \pm 0.024$	
	0.1	$0.388 \pm 0.049$	$0.390 \pm 0.049$	$0.389 \pm 0.050$	
	0.2	$0.379 \pm 0.052$	$0.382 \pm 0.053$	$0.383 \pm 0.055$	
	0.3	$0.384 \pm 0.071$	$0.387 \pm 0.071$	$0.387 \pm 0.071$	
	0.4	$0.360 \pm 0.049$	$0.363 \pm 0.049$	$0.362 \pm 0.049$	
50,000	0.5	$0.335 \pm 0.040$	$0.339 \pm 0.040$	$0.337 \pm 0.041$	
50,000	0.6	$0.331 \pm 0.040$	$0.334 \pm 0.042$	$0.333 \pm 0.042$	
	0.7	$0.312 \pm 0.025$	$0.314 \pm 0.025$	$0.313 \pm 0.026$	
	0.8	$0.306 \pm 0.024$	$0.308 \pm 0.025$	$0.306 \pm 0.024$	
	0.9	$0.301 \pm 0.023$	$0.303 \pm 0.023$	$0.303 \pm 0.023$	
	1.0	$0.302 \pm 0.023$	$0.304 \pm 0.022$	$0.303 \pm 0.022$	
	0.1	$0.306 \pm 0.027$	$0.308 \pm 0.027$	$0.307 \pm 0.027$	
	0.2	$0.303 \pm 0.025$	$0.306 \pm 0.025$	$0.305 \pm 0.026$	
	0.3	$0.304 \pm 0.027$	$0.306 \pm 0.026$	$0.305 \pm 0.027$	
	0.4	$0.301 \pm 0.023$	$0.304 \pm 0.024$	$0.302 \pm 0.024$	
100.000	0.5	$0.299 \pm 0.024$	$0.302 \pm 0.024$	$0.300 \pm 0.024$	
100,000	0.6	$0.297 \pm 0.022$	$0.299 \pm 0.022$	$0.298 \pm 0.022$	
	0.7	$0.297 \pm 0.021$	$0.299 \pm 0.021$	$0.298 \pm 0.021$	
	0.8	$0.294 \pm 0.019$	$0.296 \pm 0.019$	$0.295 \pm 0.019$	
	0.9	$0.297 \pm 0.018$	$0.299 \pm 0.018$	$0.298 \pm 0.018$	
	1.0	$0.299 \pm 0.020$	$0.301 \pm 0.021$	$0.300 \pm 0.021$	