

# UNDERSTANDING SQUARE LOSS IN TRAINING OVER-PARAMETRIZED NEURAL NETWORK CLASSIFIERS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Deep learning has achieved many breakthroughs in modern classification tasks. Numerous architectures have been proposed for different data structures but when it comes to the loss function, the cross-entropy loss is the predominant choice. Recently, several alternative losses have seen revived interests for deep classifiers. In particular, empirical evidence seems to promote square loss but a theoretical justification is still lacking. In this work, we contribute to the theoretical understanding of square loss in classification by systematically investigating how it performs for overparametrized neural networks in the neural tangent kernel (NTK) regime. Interesting properties regarding the generalization error, robustness, and calibration error are revealed. We consider two cases, according to whether classes are separable or not. In the general non-separable case, fast convergence rate is established for both misclassification rate and calibration error. When classes are separable, the misclassification rate improves to be exponentially fast. Further, the resulting margin is proven to be lower bounded away from zero, providing theoretical guarantees for robustness. We expect our findings to hold beyond the NTK regime and translate to practical settings. To this end, we conduct extensive empirical studies on practical neural networks, demonstrating the effectiveness of square loss in both synthetic low-dimensional data and real image data. Comparing to cross-entropy, square loss has comparable generalization error but noticeable advantages in robustness and model calibration.

## 1 INTRODUCTION

The pursuit of better classifiers has fueled the progress of machine learning and deep learning research. The abundance of benchmark image datasets, e.g., MNIST, CIFAR, ImageNet, etc., provides test fields for all kinds of new classification models, especially those based on deep neural networks (DNN). With the introduction of CNN, ResNets, and transformers, DNN classifiers are constantly improving and catching up to the human-level performance. In contrast to the active innovations in model architecture, the training objective remains largely stagnant, with cross-entropy loss being the default choice. Despite its popularity, cross-entropy has been shown to be problematic in some applications. Among others, Yu et al. (2020) argued that features learned from cross-entropy lack interpretability and proposed a new loss aiming for maximum coding rate reduction. Pang et al. (2019) linked the use of cross-entropy to adversarial vulnerability and proposed a new classification loss based on latent space matching. Guo et al. (2017) discovered that the confidence of most DNN classifiers trained with cross-entropy is not well-calibrated.

Recently, several alternative losses have seen revived interests for deep classifiers. In particular, many existing works have presented empirical evidence promoting the use of square loss over cross-entropy. Hui & Belkin (2020) conducted large-scale experiments comparing the two and found that square loss tends to perform better in natural language processing related tasks while cross-entropy usually yields slightly better accuracy in image classification. Similar comparisons are also made in Demirkaya et al. (2020). Kornblith et al. (2020) compared a variety of loss functions and output layer regularization strategies on the accuracy and out-of-distribution robustness, and found that square loss has greater class separation and better out-of-distribution robustness.

In comparison to the empirical investigation, theoretical understanding of square loss in training deep learning classifiers is still lacking. Through our lens, square loss has its uniqueness among

classic classification losses, and we argue that it has great potentials for modern classification tasks. Below we list our motivations and reasons why.

**Explicit feature modeling** Deep learning’s success can be largely attributed to its superior ability as feature extractors. For classification, the ideal features should be separated between classes and concentrated within classes. However, when optimizing cross-entropy loss, it’s not clear what the learned features should look like (Yu et al., 2020). In comparison, square loss uses the label codings (one-hot, simplex etc.) as features, which can be modeled explicitly to control class separations.

**Model Calibration** An ideal classifier should not only give the correct class prediction, but also with the correct confidence. Calibration error measures the closeness of the predicted confidence to the underlying conditional probability  $\eta$ . Using square loss in classification can be essentially viewed as regression where it treats discrete labels as continuous code vectors. It can be shown that the optimal classifier under square loss is  $2\eta - 1$ , linear with the ground truth. This distinguishing property allows it to easily recover  $\eta$ . In comparison, the optimal classifiers under the hinge loss and cross-entropy are  $\text{sign}(2\eta - 1)$  and  $\log(\frac{\eta}{1-\eta})$ , respectively. Therefore, hinge loss doesn’t provide reliable information on the prediction confidence, and cross-entropy can be problematic when  $\eta$  is close to 0 or 1 (Zhang, 2004). Hence, in terms of model calibration, square loss is a natural choice.

**Connections to popular approaches** Mixup (Zhang et al., 2017) is a popular data augmentation technique where augmented data are constructed via convex combinations of inputs and their labels. Like in square loss, mixup treats labels as continuous and is shown to improve the generalization of DNN classifiers. In knowledge distillation (Hinton et al., 2015), where a student classifier is trying to learn from a trained teacher, Menon et al. (2021) proved that the “optimal” teacher with the ground truth conditional probabilities provides the lowest variance in student learning. Since classifiers trained using square loss is a natural consistent estimator of  $\eta$ , one can argue that it is a better teacher. In supervised contrastive learning (Khosla et al., 2020), the optimal features are the same as those from square loss with simplex label coding (Graf et al., 2021) (details in Section 4).

Despite its lack of popularity in practice, square loss has many advantages that can be easily overlooked. In this work, we systematically investigate from a statistical estimation perspective, the properties of deep learning classifiers trained using square loss. The neural networks in our analysis are required to be sufficiently overparametrized in the neural tangent kernel (NTK) regime. Even though this restricts the implication of our results, it is a necessary first step towards a deeper understanding. In summary, our main contributions are:

- **Generalization error bound:** We consider two cases, according to whether classes are separable or not. In the general non-separable case, we adopt the classical binary classification setting with smooth conditional probability. Fast rate of convergence is established for overparametrized neural network classifiers with Tsybakov’s noise condition. If two classes are separable with positive margins, we show that overparametrized neural network classifiers can provably reach zero misclassification error with probability *exponentially* tending to one. To the best of our knowledge, this is the *first* such result for separable but not linear separable classes. Furthermore, we bridge these two cases and offer a *unified* view by considering auxiliary random noise injection.
- **Robustness (margin property):** When two classes are separable, the decision boundary is not unique and large-margin classifiers are preferred. In the separable case, we further show that the decision boundary of overparametrized neural network classifiers trained by square loss cannot be too close to the data support and the resulting margin is lower bounded away from zero, providing theoretical guarantees for robustness.
- **Calibration error:** We show that classifiers trained using square loss are inherently well-calibrated, i.e., the trained classifier provides consistent estimation of the ground-truth conditional probability in  $L_\infty$  norm. Such property doesn’t hold for cross-entropy.
- **Empirical evaluation:** We corroborate our theoretical findings with empirical experiments in both synthetic low-dimensional data and real image data. Comparing to cross-entropy, square loss has comparable generalization error but noticeable advantages in robustness and model calibration.

This work contributes towards the theoretical understanding of deep classifiers, from an estimation point of view, which has been a classic topic in statistics literature. Among others, Mammen & Tsybakov (1999) established the optimal convergence rate for 0-1 loss excess risk when the decision boundary is smooth. Zhang (2004); Bartlett et al. (2006) extended the analysis to various surrogate

losses. Audibert & Tsybakov (2007); Kohler & Krzyzak (2007) studied the convergence rates for plug-in classifiers from local averaging estimators. Steinwart et al. (2007) investigated the convergence rate for support vector machine using Gaussian kernels. We build on and extend classic results to neural networks in the NTK regime. Comparing to existing works on deep learning classification, e.g., Kim et al. (2018) derived fast convergence rates of ReLU DNN classifiers that minimize the empirical hinge loss, our results incorporate the training algorithm and apply to trained classifiers.

We require the neural network to be overparametrized, which has been extensively studied recently, under the umbrella term NTK. Most such results are in the regression setting with a handful of exceptions. Ji & Telgarsky (2019) showed that only polylogarithmic width is sufficient for gradient descent to overfit the training data using logistic loss. Hu et al. (2020) proved generalization error bound for regularized NTK in classification. Cao & Gu (2019; 2020) provided optimization and generalization guarantees for overparametrized network trained with cross-entropy. In comparison, our results are sharper in the sense that we take the ground truth data assumptions into consideration. This allows a faster convergence rate, especially when the classes are separable, where the exponential convergence rate is attainable. The NTK framework greatly reduces the technical difficulty for our theoretical analysis. However, our results are mainly due to properties of the square loss itself and we expect them to hold for a wide range of classifiers.

There are other works investigating the use of square loss for training (deep) classifiers. Han et al. (2021) uncovered that the “neural collapse” phenomenon also occurs under square loss where the last-layer features eventually collapse to their simplex-style class-means. Muthukumar et al. (2020) compared classification and regression tasks in the overparameterized linear model with Gaussian features, illustrating different roles and properties of loss functions used at the training and testing phases. Poggio & Liao (2019) made interesting observations on effects of popular regularization techniques such as batch normalization and weight decay on the gradient flow dynamics under square loss. These findings support our theoretical results’ implication, which further strengthens our beliefs that the essence comes from the square loss and our analysis can go beyond NTK regime.

The rest of this paper is arranged as follows. Section 2 presents some preliminaries. Main theoretical results are in Section 3. The simplex label coding is discussed in Section 4 followed by numerical studies in Section 5 and conclusions in Section 6. Technical proofs and details of the numerical studies can be found in the Appendix.

## 2 PRELIMINARIES

**Notation** For a function  $f : \Omega \rightarrow \mathbb{R}$ , let  $\|f\|_\infty = \sup_{\mathbf{x} \in \Omega} |f(\mathbf{x})|$  and  $\|f\|_p = (\int_\Omega |f(\mathbf{x})|^p d\mathbf{x})^{1/p}$ . For a vector  $\mathbf{x}$ ,  $\|\mathbf{x}\|_p$  denotes its  $p$ -norm, for  $1 \leq p \leq \infty$ .  $L_p$  and  $l_p$  are used to distinguish function norms and vector norms. For two positive sequences  $\{a_n\}_{n \in \mathbb{N}}$  and  $\{b_n\}_{n \in \mathbb{N}}$ , we write  $a_n \lesssim b_n$  if there exists a constant  $C > 0$  such that  $a_n \leq Cb_n$  for all sufficiently large  $n$ . We write  $a_n \asymp b_n$  if  $a_n \lesssim b_n$  and  $b_n \lesssim a_n$ . Let  $[N] = \{1, \dots, N\}$  for  $N \in \mathbb{N}$ ,  $\mathbb{I}$  be the indicator function, and  $\mathbf{I}_d$  be the  $d \times d$  identity matrix.  $N(\mu, \Sigma)$  represents Gaussian distribution with mean  $\mu$  and covariance  $\Sigma$ .

**Classification problem settings** Let  $P$  be an underlying probability measure on  $\Omega \times \mathbf{Y}$ , where  $\Omega \subset \mathbb{R}^d$  is compact and  $\mathbf{Y} = \{1, -1\}$ . Let  $(X, Y)$  be a random variable with respect to  $P$ . Suppose we have observations  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \subset (\Omega \times \mathbf{Y})^n$  i.i.d. sampled according to  $P$ . The classification task is to predict the unobserved label  $y$  given a new input  $\mathbf{x} \in \Omega$ . Let  $\eta$  defined on  $\Omega$  denote the conditional probability, i.e.,  $\eta(\mathbf{x}) = \mathbb{P}(y = 1 | \mathbf{x})$ . Let  $P_X$  be the marginal distribution of  $P$  on  $X$ . The key quantity of interest is the misclassification error, i.e., 0-1 loss. In the population level, the 0-1 loss can be written as

$$L(f) = \mathbb{E}_{(X,Y) \sim P} \mathbb{I}\{\text{sign}(f(X)) \neq Y\} = \mathbb{E}_{X \sim P_X} [(1 - \eta(X)) \mathbb{I}\{f(X) \geq 0\} + \eta(X) \mathbb{I}\{f(X) < 0\}], \quad (2.1)$$

where the expectation is taken with respect to the probability measure  $P$ . Clearly, an optimal classifier with the minimal 0-1 loss is  $2\eta - 1$ .

According to whether labels are deterministic, there are two scenarios of interest. If  $\eta$  only takes values from  $\{0, 1\}$ , i.e., labels are deterministic, we call this case the *separable case*<sup>1</sup>. Let  $\Omega_1 =$

<sup>1</sup>In the separable case we consider, the classes are not limited to linearly separable but can be arbitrarily complicated.

$\{x|\eta(x) = 1\}$ ,  $\Omega_2 = \{x|\eta(x) = 0\}$  and  $\Omega = \Omega_1 \cup \Omega_2$ . If the probability measure of  $\{x|\eta(x) \in (0, 1)\}$  is non-zero, i.e., the labels contain randomness, we call this case the *non-separable case*. In the separable case, we further assume that there exists a positive margin, i.e.,  $\text{dist}(\Omega_1, \Omega_2) \geq 2\gamma > 0$ , where  $\gamma$  is a constant, and  $\text{dist}(\Omega_1, \Omega_2) = \inf_{x \in \Omega_1, x' \in \Omega_2} \|x - x'\|_2$ . In the non-separable case, to quantify the difficulty of classification, we adopt the well-established Tsybakov’s noise condition (Audibert & Tsybakov, 2007), which measures how large the “difficult region” is where  $\eta(x) \approx 1/2$ .

**Definition 2.1** (Tsybakov’s noise condition). Let  $\kappa \in [0, \infty]$ . We say  $P$  has Tsybakov noise exponent  $\kappa$  if there exists a constant  $C, T > 0$  such that for all  $0 < t < T$ ,  $P_X(|2\eta(X) - 1| < t) \leq C \cdot t^\kappa$ .

A large value of  $\kappa$  implies the difficult region to be small. It is expected that a larger  $\kappa$  leads to a faster convergence rate of a neural network classifier. This intuition is verified for the overparametrized neural network classifier trained by square loss and  $\ell_2$  regularization. See Section 3 for more details.

**Neural network setup** We mainly focus on the one-hidden-layer ReLU neural network family  $\mathcal{F}$  with  $m$  nodes in the hidden layer, denoted by

$$f_{\mathbf{W}, \mathbf{a}}(\mathbf{x}) = \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \sigma(\mathbf{W}_r^\top \mathbf{x}),$$

where  $\mathbf{x} \in \Omega$ ,  $\mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_m) \in \mathbb{R}^{d \times m}$  is the weight matrix in the hidden layer,  $\mathbf{a} = (a_1, \dots, a_m)^\top \in \mathbb{R}^m$  is the weight vector in the output layer,  $\sigma(z) = \max\{0, z\}$  is the rectified linear unit (ReLU). The initial values of the weights are independently generated from

$$\mathbf{W}_r(0) \sim N(\mathbf{0}, \xi^2 \mathbf{I}_m), \quad a_r \sim \text{unif}\{-1, 1\}, \quad \forall r \in [m].$$

Based on the observations  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , the goal of training a neural network is to find a solution to

$$\min_{\mathbf{W}} \sum_{i=1}^n l(f_{\mathbf{W}, \mathbf{a}}(\mathbf{x}_i), y_i) + \mu \mathcal{R}(\mathbf{W}, \mathbf{a}), \quad (2.2)$$

where  $l$  is the loss function,  $\mathcal{R}$  is the regularization, and  $\mu \geq 0$  is the regularization parameter. Note in Equation 2.2 that we only consider training the weights  $\mathbf{W}$ . This is because  $a \cdot \sigma(z) = \text{sign}(a) \cdot \sigma(|a|z)$ , which allows us to reparametrize the network to have all  $a_i$ ’s to be either 1 or  $-1$ . In this work, we consider square loss associated with  $\ell_2$  regularization, i.e.,  $l(f_{\mathbf{W}, \mathbf{a}}(\mathbf{x}_i), y_i) = (f_{\mathbf{W}, \mathbf{a}}(\mathbf{x}_i) - y_i)^2$  and  $\mathcal{R}(\mathbf{W}, \mathbf{a}) = \|\mathbf{W}\|_2^2$ .

A popular way to train the neural network is via gradient based methods. It has been shown that the training process of DNNs can be characterized by the neural tangent kernel (NTK) (Jacot et al., 2018). As is usually assumed in the NTK literature (Arora et al., 2019; Hu et al., 2020; Bietti & Mairal, 2019; Hu et al., 2021), we consider data on the unit sphere  $\mathbb{S}^{d-1}$ , i.e.,  $\|\mathbf{x}_i\|_2 = 1, \forall i \in [n]$ , and the neural network is highly overparametrized ( $m \gg n$ ) and trained by gradient descent (GD). For details about NTK and GD in one-hidden-layer ReLU neural networks, we refer to Appendix A. In the rest of this work, we use  $f_{\mathbf{W}^{(k)}, \mathbf{a}}$  to denote the GD-trained neural network classifier under square loss associated with  $\ell_2$  regularization, where  $k$  is the iteration number satisfying Assumption D.1 and  $\mathbf{W}^{(k)}$  is the weight matrix after  $k$ -th iteration.

### 3 THEORETICAL RESULTS

In this section, we present our main theoretical results. Throughout the analysis, we assume that the overparametrized neural network  $f_{\mathbf{W}, \mathbf{a}}$  and the training process via GD satisfy Assumption D.1 (see Appendix D), which essentially requires the neural network to be sufficiently overparametrized (with a finite width), and imposes conditions on the learning rate and iteration number. Our theoretical results consist of three parts: generalization error, robustness, and calibration error.

#### 3.1 GENERALIZATION ERROR BOUND

In classification, the generalization error is typically referred to as the misclassification error, which can be quantified by  $L(f)$  defined in Equation 2.1. In the non-separable case, the excess risk, defined by  $L(f) - L^*$ , is used to evaluate the quality of a classifier  $f$ , where  $L^* = L(2\eta - 1)$ , which minimizes the 0-1 loss. The following theorem states that the overparametrized neural network with GD and  $\ell_2$  regularization can achieve a small excess risk in the non-separable case.

**Theorem 3.1** (Excess risk in the non-separable case). Suppose Assumptions D.1, D.2, and D.4 hold. Assume the conditional probability  $\eta(\mathbf{x})$  satisfies Tsybakov’s noise condition with component  $\kappa$ . Let  $\mu \asymp n^{\frac{d-1}{2d-1}}$ . Then

$$L(f_{\mathbf{W}(k), \mathbf{a}}) = L^* + O_{\mathbb{P}}(n^{-\frac{d(\kappa+1)}{(2d-1)(\kappa+2)}}). \quad (3.1)$$

From Theorem 3.1, we can see that as  $\kappa$  becomes larger, the convergence rate becomes faster, which is intuitively true. Generalization error bounds in this setting is scarce. To the best of the authors’ knowledge, Hu et al. (2020) is the closest work (the labels are randomly flipped), where the bound is in the order of  $O_{\mathbb{P}}(1/\sqrt{n})$ . Our bound is faster, especially with larger  $\kappa$ . It is known that the optimal convergence rate under Assumptions D.2 and D.4 is  $O_{\mathbb{P}}(n^{-\frac{d(\kappa+1)}{d\kappa+4d-2}})$  (Audibert & Tsybakov, 2007). The differences between Equation 3.1 and the optimal convergence rate is that there is an extra  $(d-1)\kappa$  in the denominator of the convergence rate in Equation 3.1 (since  $n^{-\frac{d(\kappa+1)}{(2d-1)(\kappa+2)}} = n^{-\frac{d(\kappa+1)}{(d-1)\kappa+d\kappa+4d-2}}$ ). If the conditional probability  $\eta$  has a bounded Lipschitz constant, then Kohler & Krzyzak (2007) showed that the convergence rate based on the plug-in kernel estimate is  $O_{\mathbb{P}}(n^{-\frac{\kappa+1}{\kappa+3+d}})$ , which is slower than the rate in Equation 3.1 if  $d$  is large.

Now we turn to the separable case. Since  $\eta$  only takes value from  $\{0, 1\}$  in the separable case,  $\eta$  is bounded away from  $1/2$ . Therefore, one can trivially take  $\kappa \rightarrow \infty$  in Equation 3.1 and obtain the convergence rate  $O_{\mathbb{P}}(n^{-d/(2d-1)})$ . However, this rate can be significantly improved in the separable case, as stated in the following theorem.

**Theorem 3.2** (Generalization error in the separable case). Suppose Assumptions D.1, D.3, and D.5 hold. Let  $\mu = o(1)$ . There exist positive constants  $C_1, C_2$  such that the misclassification rate is 0% with probability at least  $1 - \delta - C_1 \exp(-C_2 n)$ , and  $\delta$  can be arbitrarily small<sup>2</sup> by enlarging the neural network’s width.

Note that in Theorem 3.2, the regularization parameter can take any rate that converges to zero. In particular,  $\mu$  can be zero, and the corresponding classifier overfits the training data. Theorem 3.2 states that the convergence rate in the separable case is exponential, if a sufficiently wide neural network is applied. This is because the observed labels are not corrupted by noise, i.e.,  $\mathbb{P}(y = 1|\mathbf{x})$  is either one or zero. Therefore, it is easier to classify separable data, which is intuitively true.

### 3.2 ROBUSTNESS AND CALIBRATION ERROR

If two classes are separable with positive margin, the decision boundary is not unique. Practitioners often prefer the decision boundary with large margins, which are robust against possible perturbation on input points (Elsayed et al., 2018; Ding et al., 2018). The following theorem states that the square loss trained margin can be lower bounded by a positive constant. Recall that in the separable case,  $\Omega = \Omega_1 \cup \Omega_2$ , where  $\Omega_1 = \{\mathbf{x}|\eta(\mathbf{x}) = 1\}$  and  $\Omega_2 = \{\mathbf{x}|\eta(\mathbf{x}) = 0\}$ .

**Theorem 3.3** (Robustness in the separable case). Suppose the assumptions of Theorem 3.2 are satisfied. Let  $\mu = o(1)$ . Then there exist positive constants  $C, C_1, C_2$  such that for all  $n$ ,

$$\min_{\mathbf{x} \in \mathcal{D}_T, \mathbf{x}' \in \Omega_1 \cup \Omega_2} \|\mathbf{x} - \mathbf{x}'\|_2 \geq C,$$

and the misclassification rate is 0% with probability at least  $1 - \delta - C_1 \exp(-C_2 n)$ , where  $\mathcal{D}_T$  is the decision boundary, and  $\delta$  is as in Theorem 3.2.

**Remark 1.** Note that  $\|\mathbf{x} - \mathbf{x}'\|_{\infty} \geq \sqrt{d}\|\mathbf{x} - \mathbf{x}'\|_2$ , thus Theorem 3.3 also indicates  $l_{\infty}$  robustness.

In the non-separable case,  $\eta(\mathbf{x})$  varies within (0,1) and practitioners may not only want a classifier with a small excess risk, but also want to recover the underlying conditional probability  $\eta$ . Therefore, square loss is naturally preferred since it treats the classification problem as a regression problem. The following theorem states that, one can recover the conditional probability  $\eta$  by using an overparametrized neural network with  $\ell_2$  regularization and GD training.

**Theorem 3.4** (Calibration error). Suppose the conditions in Theorem 3.1, Assumption D.3 and D.4 are fulfilled. Let  $\mu \asymp n^{\frac{d-1}{2d-1}}$ . Then

$$\|(f_{\mathbf{W}(k), \mathbf{a}} + 1)/2 - \eta\|_{L_{\infty}} = O_{\mathbb{P}}(n^{-\frac{1}{4d-2}}). \quad (3.2)$$

<sup>2</sup>The term  $\delta$  only depends on the width of the neural network. A smaller  $\delta$  requires a wider neural network. If  $\delta = 0$ , then the number of nodes in the hidden layer is infinity.

Theorem 3.4 states that the underlying conditional probability in the non-separable case can be recovered by  $(f_{\mathbf{W}^{(k)},a} + 1)/2$ . The form  $(f_{\mathbf{W}^{(k)},a} + 1)/2$  is to account for the  $\{-1, 1\}$  label coding. Under  $\{0, 1\}$  coding, the estimator would be  $f_{\mathbf{W}^{(k)},a}$  itself. The  $L_\infty$  consistency doesn't hold for cross-entropy trained neural networks, due to the form of the optimal solution  $\log(\frac{\eta}{1-\eta})$ . With limited capacity, the network's confidence prediction is bounded away from 0 and 1 (Zhang, 2004). In practice, we want to control the complexity of the neural network thus it is usually the case that  $\|f_{\mathbf{W}^{(k)},a}\|_\infty < C$  for some constant  $C$ . Hence, it cannot accurately estimate  $\eta(\mathbf{x})$  when  $\eta(\mathbf{x}) > \frac{e^C}{1+e^C}$  or  $\eta(\mathbf{x}) < \frac{1}{1+e^C}$ , which makes the calibration error under the cross-entropy loss always bounded away from zero. However, square loss does not have such a problem.

Notice that the calibration error bound in Theorem 3.4 does not depend on the Tsybakov's noise condition, and is slower than the excess risk. This is because, a small calibration error is much stronger than a small excess risk, since the former requires the conditional probability estimation to be *uniformly* accurate, not just matching the sign of  $\eta(\mathbf{x}) - 1/2$ . To be more specific, a good estimated  $\hat{\eta}$  can always result in a low risk plug-in classifier  $\hat{f}(\mathbf{x}) = 2\hat{\eta}(\mathbf{x}) - 1$ , but not vice versa.

**Remark 2** (Technical challenge). Despite the similar forms of regression and classification using square loss, most of the regression analysis techniques cannot be directly applied to the classification problem, even if the supports of two classes are non-separable. Moreover, it is clear that classification problems in the separable case are completely different with regression problems.

**Remark 3** (Extension on NTK). Although our analysis only concerns overparametrized one-hidden-layer ReLU neural networks, it can potentially apply to other types of neural networks in the NTK regime. Recently, it has been shown that overparametrized multi-layer networks correspond to the Laplace kernel (Geifman et al., 2020; Chen & Xu, 2020). As long as the trained neural networks can approximate the classifier induced by the NTK, our results can be naturally extended.

### 3.3 TRANSITION FROM SEPARABLE TO NON-SEPARABLE

The general non-separable case and the special separable case can be connected via Gaussian noise injection. In practice, data augmentation is an effective way to improve robustness and the simplest way is Gaussian noise injection (He et al., 2019). In this section, we only consider it as an auxiliary tool for theoretical analysis purpose and not for actual robust training. Injecting Gaussian noise amounts to convoluting a Gaussian distribution  $N(0, v^2 \mathbf{I}_d)$  to the marginal distribution  $P_X$ , which enlarges both  $\Omega_1$  and  $\Omega_2$  to  $\mathbb{R}^d$  and a unique decision boundary  $\mathcal{D}_v$  can be induced. Correspondingly, the "noisy" conditional probability, denoted as  $\tilde{\eta}_v$ , is also smoothed to be continuous on  $\mathbb{R}^d$ . As  $v \rightarrow 0$ ,  $\|\tilde{\eta}_v - \eta\|_\infty \rightarrow 0$  on  $\Omega_1$  and  $\Omega_2$  and the limiting  $\tilde{\eta}_0$  is a piecewise constant function with discontinuity at the induced decision boundary.

**Lemma 3.5** (Tsybakov's noise condition under Gaussian noises). Let the margin be  $2\gamma > 0$ , the noise be  $N(0, v^2 \mathbf{I}_d)$ . Then there exist some constants  $T, C > 0$  such that for any  $0 < t < T$ ,

$$P_X(|2\tilde{\eta}_v(X) - 1| < t) \leq \frac{Cv^2}{\gamma} \exp\left(-\frac{\gamma^2}{2v^2}\right) \cdot t.$$

**Theorem 3.6** (Exponential convergence rate). Suppose the classes are separable with margin  $2\gamma > 0$ . No matter how complicated  $\Omega_1 \cup \Omega_2$  are, the excess risk of the over parameterized neural network classifier satisfying Assumptions D.1 and D.4 has the rate  $O_{\mathbb{P}}(e^{-n\gamma/7})$ .

The proof of Theorem 3.6 involves taking the auxiliary noise to zero, e.g.,  $v = v_n \asymp 1/\sqrt{n}$ . The exponential convergence rate is a direct outcome of Lemma 3.5 and Theorem 3.1. Note that our exponential convergence rate is much faster than existing ones under the similar separable setting (Ji & Telgarsky, 2019; Cao & Gu, 2019; 2020), which are all polynomial with  $n$ , e.g.,  $O_{\mathbb{P}}(1/\sqrt{n})$ .

**Remark 4.** Theorems 3.4 and 3.6 share the same gist that the over parameterized neural network classifiers can have exponential convergence rate when data are separable with positive margin, while the result of Theorem 3.6 is weaker than that of Theorem 3.4, but with milder conditions. Nevertheless, Theorem 3.6 bridges the non-separable case and separable case.

## 4 MULTICLASS CLASSIFICATION

In binary classification, the labels are usually encoded as  $-1$  and  $1$ . When there are  $K > 2$  classes, the default label coding is one-hot. However, it is empirically observed that this vanilla square loss struggles when the number of classes are large, for which scaling tricks have been proposed (Hui & Belkin, 2020; Demirkaya et al., 2020). Another popular coding scheme is the simplex coding (Mroueh et al., 2012), which takes maximally separated  $K$  points on the sphere as label features. When  $K = 2$ , this reduces to the typical  $-1, 1$  coding. Many advantages of the simplex coding have been discussed, including its relationship with cross-entropy loss and supervised contrastive learning (Papayan et al., 2020; Han et al., 2021; Graf et al., 2021; Fang et al., 2021).

In this work, we adopt the simplex coding. More discussion and empirical comparison about the coding choices can be found in Appendix G.2. Given the label coding, one can easily generalize the theoretical development in Section 3 by employing the following objective function

$$\min_{\mathbf{W}} \sum_{j=1}^K \sum_{i=1}^n (f_{j,\mathbf{W},\alpha}(\mathbf{x}_i) - y_{i,j})^2 + \mu \|\mathbf{W}\|_2^2,$$

where  $f_{\mathbf{W},\alpha} : \Omega \mapsto \mathbb{R}^K$ , and  $\mathbf{y}_i = (y_{i,1}, \dots, y_{i,K})^\top$  is the label of  $i$ -th observation.

The following proposition states a relationship between the simplex coding scheme and the conditional probability.

**Proposition 4.1** (Conditional probability). Let  $f^* : \Omega \rightarrow \mathbb{R}^K$  minimize the mean square error  $\mathbb{E}_X(f^*(X) - \mathbf{v}_y)^2$ , where  $\mathbf{v}_y$  is the simplex coding vector of label  $y$ . Then we have

$$\eta_k(\mathbf{x}) := \mathbb{P}(y = k|\mathbf{x}) = ((K-1)f^*(\mathbf{x})^\top \mathbf{v}_k + 1) / K. \quad (4.1)$$

Unlike the softmax function when using cross entropy, the estimated conditional probability using square loss is not guaranteed to be within 0 and 1. This will cause issues for adversarial attacks, which will be discussed in detail in Appendix G.2.

## 5 NUMERICAL STUDIES

Although our theoretical results are for overparametrized neural network in the NTK regime, we expect our conclusions to generalize to practical network architectures. The focus of this section is not on improving the state-of-the-art performance for deep classifiers, but to illustrate the difference between cross-entropy and square loss. We provide experiment results on both synthetic and real data, to support our theoretical findings and illustrate the practical benefits of square loss in training overparametrized DNN classifiers. Compared with cross-entropy, the square loss has comparable generalization performance, but with stronger robustness and smaller calibration error.

### 5.1 SYNTHETIC DATA

We consider the square loss based and cross-entropy based overparametrized neural networks (ONN) with  $\ell_2$  regularization, denoted as SL-ONN +  $\ell_2$  and CE-ONN +  $\ell_2$ , respectively. The chosen ONNs are two-hidden-layer ReLU neural networks with 500 neurons for each layer, and the parameter  $\mu$  is selected via a validation set. More implementation details are in Appendix G.1.

**Separable case** We consider two separated classes with spiral curve like supports. We also present the performance of the cross-entropy based ONN without  $\ell_2$  regularization (CE-ONN). Figure 1 shows one instance of the test misclassification rate and decision boundaries attained by SL-ONN +  $\ell_2$  (Left), CE-ONN +  $\ell_2$  (Center), and CE-ONN (Right). From this example and other examples in Appendix G.1, it can be seen that SL-ONN +  $\ell_2$  has a smaller test misclassification rate and a much smoother decision boundary. In particular, in the red region, where the training data are sparse, SL-ONN +  $\ell_2$  fits the correct data distribution best.

**Non-separable case** We consider the conditional probability  $\eta(\mathbf{x}) = \sin(\sqrt{2}\pi\|\mathbf{x}\|_2), \mathbf{x} \in [-1, 1]^2$ , and the calibration performance of SL-ONN +  $\ell_2$  and CE-ONN +  $\ell_2$ , where the classifiers are denoted by  $\hat{f}_{l2}$  and  $\hat{f}_{ce}$ , respectively. The results are presented in Figure G.8 in the Appendix.

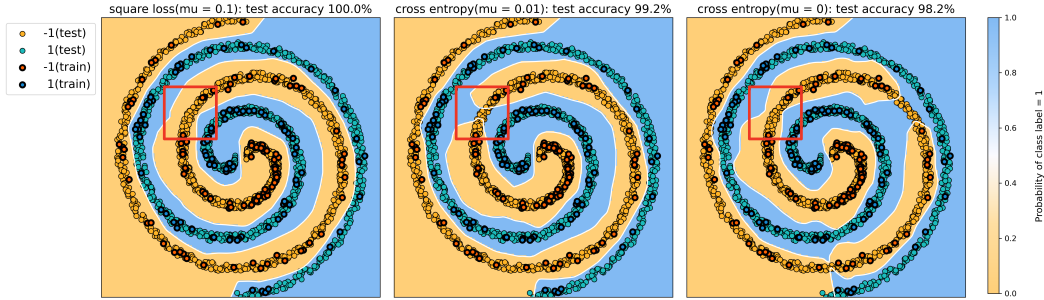


Figure 1: Test misclassification rates and decision boundaries predicted by: SL-ONN +  $\ell_2$  (Left); CE-ONN +  $\ell_2$  (Center); CE-ONN (Right) for the separable case.

The error bar plot of the test calibration error shows that  $\hat{f}_{l_2}$  has the smaller mean and standard deviation than  $\hat{f}_{ce}$ . It suggests that square loss generally outperforms cross entropy in calibration. The histogram and kernel density estimation of the test calibration errors for one case show that the pointwise calibration errors on the test points of  $\hat{f}_{l_2}$  are more concentrated around zero than those of  $\hat{f}_{ce}$ . Moreover, despite a comparable misclassification rate with  $\hat{f}_{ce}$ ,  $\hat{f}_{l_2}$  has a smaller calibration error. Figure G.8 demonstrates that SL-ONN +  $\ell_2$  recovers  $\eta$  much better than CE-ONN +  $\ell_2$ .

## 5.2 REAL DATA

To make a fair comparison, we adopt popular architectures, ResNet (He et al., 2016) and Wide ResNet (Zagoruyko & Komodakis, 2016) and evaluate them on the CIFAR image classification datasets, with only the training loss function changed, from cross-entropy (CE) to square loss with simplex coding (SL). Further, we don’t employ any large scale hyper-parameter tuning and all the parameters are kept as default except for the learning rate (lr) and batch size (bs), where we are choosing from the better of (lr=0.01, bs=32) and (lr=0.1, bs=128). Each experiment setting is replicated 5 times and we report the average performance followed by its standard deviation in the parenthesis. (lr=0.01, bs=32) works better for the most cases except for square loss trained WRN-16-10 on CIFAR-100. More experiment details and additional results can be found in Appendix G.2.

**Generalization** In both CIFAR-10 and CIFAR-100, the performance of cross-entropy and square loss with simplex coding are quite comparable, as observed in Hui & Belkin (2020). Cross-entropy tends to perform slightly better for ResNet, especially on CIFAR-100 with an advantage of less than 1%. There is a more significant gap with Wide ResNet where square loss outperforms cross-entropy by more than 1% on both CIFAR-10 and CIFAR-100. The details can be found in Table 1.

Table 1: Test accuracy on CIFAR datasets. Average accuracy larger than 0 but less than 0.1 is denoted as 0\* without standard deviation.

| Dataset   | Network   | Loss | Clean acc %         | PGD-100 ( $l_\infty$ -strength) |                     |                    | AutoAttack ( $l_\infty$ -strength) |                    |       |
|-----------|-----------|------|---------------------|---------------------------------|---------------------|--------------------|------------------------------------|--------------------|-------|
|           |           |      |                     | 2/255                           | 4/255               | 8/255              | 2/255                              | 4/255              | 8/255 |
| CIFAR-10  | ResNet-18 | CE   | <b>95.15 (0.11)</b> | 8.81 (1.61)                     | 0.65 (0.24)         | 0                  | 2.74 (0.09)                        | 0                  | 0     |
|           |           | SL   | 95.04 (0.07)        | <b>30.53 (0.92)</b>             | <b>6.64 (0.67)</b>  | <b>0.86 (0.24)</b> | <b>4.10 (0.50)</b>                 | <b>0*</b>          | 0     |
|           | WRN-16-10 | CE   | 93.94 (0.16)        | 1.04 (0.10)                     | 0                   | 0                  | 0.33 (0.06)                        | 0                  | 0     |
|           |           | SL   | <b>95.02 (0.11)</b> | <b>37.47 (0.61)</b>             | <b>23.16 (1.28)</b> | <b>7.88 (0.72)</b> | <b>5.37 (0.50)</b>                 | <b>0*</b>          | 0     |
| CIFAR-100 | ResNet-50 | CE   | <b>79.82 (0.14)</b> | 2.31 (0.07)                     | 0*                  | 0                  | 0.99 (0.10)                        | 0*                 | 0     |
|           |           | SL   | 78.91 (0.14)        | <b>13.76 (1.30)</b>             | <b>4.63 (1.20)</b>  | <b>1.21 (0.80)</b> | <b>3.67 (0.60)</b>                 | <b>0.16 (0.05)</b> | 0     |
|           | WRN-16-10 | CE   | 77.89 (0.21)        | 0.83 (0.07)                     | 0*                  | 0                  | 0.42 (0.07)                        | 0                  | 0     |
|           |           | SL   | <b>79.65 (0.15)</b> | <b>6.48 (0.40)</b>              | <b>0.42 (0.04)</b>  | <b>0*</b>          | <b>2.73 (0.20)</b>                 | <b>0*</b>          | 0     |

**Adversarial robustness** Normally trained deep classifiers are found to be adversarially vulnerable and adversarial attacks provide a powerful tool to evaluate classification robustness. For our experiment, we consider the black-box Gaussian noise attack, the classic white-box PGD attack (Madry et al., 2017) and the state-of-the-art AutoAttack (Croce & Hein, 2020), with attack strength level 2/255, 4/255, 8/255 in  $l_\infty$  norm. AutoAttack contains both white-box and black-box attacks and offers a more comprehensive evaluation of adversarial robustness. The Gaussian noises results are



Table 2: Performance on CIFAR-10 dataset for ResNet-18 under standard PGD adversarial training.

| CIFAR10   | Loss | Acc (%)      | PGD steps | Strength( $l_\infty$ ) | Autoattack   |
|-----------|------|--------------|-----------|------------------------|--------------|
| ResNet-18 | CE   | 86.87        | 3         | 8/255                  | 37.08        |
|           |      | 84.50        | 7         | 8/255                  | 41.88        |
| ResNet-18 | SL   | <b>87.31</b> | 3         | 8/255                  | <b>40.46</b> |
|           |      | <b>84.52</b> | 7         | 8/255                  | <b>44.76</b> |

presented in Table G.3 in the Appendix. At different noise levels, square loss consistently outperforms cross-entropy, especially for WRN-16-10, with around 2-4% accuracy improvement. More details can be found in Appendix G.2. The PGD and AutoAttack results are reported in Table 1. Even though classifiers trained with square loss is far away from adversarially robust, it consistently gives significantly higher adversarial accuracy. The same margin can be carried over to standard adversarial training as well. Table 2 lists results from standard PGD adversarial training with CE and SL. By substituting cross-entropy loss to square loss, the robust accuracy increased around 3% while maintaining higher clean accuracy.

One thing to notice is that when constructing white-box attacks, square loss will not work well since it doesn't directly reflect the classification accuracy. More specifically, for a correctly classified image  $(x, y)$ , maximizing the square loss may result in linear scaling of the classifier  $f(x)$ , which doesn't change the predicted class (see Appendix G.2 for more discussion). To this end, we consider a special attack for classifiers trained by square loss by maximizing the cosine similarity between  $f(x)$  and  $v_y$ . We call this angle attack and also utilize it for the PGD adversarial training paired with square loss in Table 2. In our experiments, this special attack rarely outperforms the standard PGD with cross-entropy and the reported PGD accuracy are from the latter settings. This property of square loss may be an advantage in defending adversarial attacks.

**Model calibration** The predicted class probabilities for square loss can be obtained from Equation 4.1. Expected calibration error (ECE) measures the absolute difference between predicted confidence and the actual accuracy. Deep classifiers are usually found to be over-confident (Vaicenavicius et al., 2019). Using ResNet as an example, we report the typical reliability diagram in Figure 2. On CIFAR-10 with ResNet-18, the average ECE for cross-entropy is 0.028 (0.002) while that for square loss is 0.0097 (0.001). On CIFAR-100 with ResNet-50, the average ECE for cross-entropy is 0.094 (0.005) while that for square loss is 0.068 (0.005). Square loss results are much more calibrated with significantly smaller ECE.

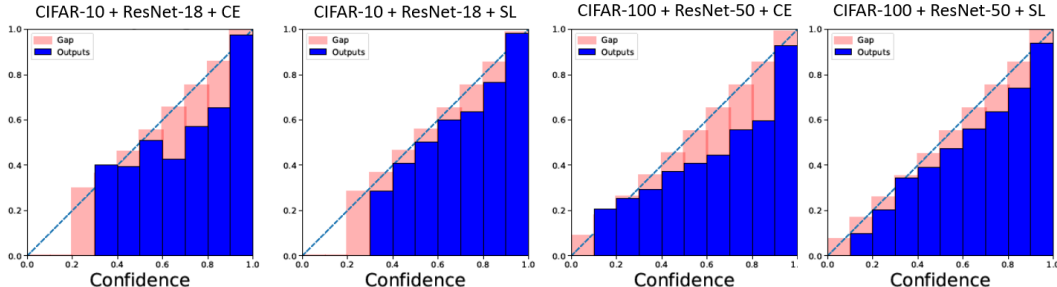


Figure 2: Reliability diagrams of ResNet-18 on CIFAR-10 and ResNet-50 on CIFAR-100. Square loss trained models behave more well-calibrated while cross-entropy trained ones tend to be visibly more over-confident.

## 6 CONCLUSIONS

Classification problems are ubiquitous in deep learning. As a fundamental problem, any progress in classification can potentially benefit numerous relevant tasks. Despite its lack of popularity in practice, square loss has many advantages that can be easily overlooked. Through both theoretical analysis and empirical studies, we identify several ideal properties of using square loss in training neural network classifiers, including provable fast convergence rates, strong robustness, and small calibration error. We encourage readers to try square loss in your own application scenarios.

**Ethnics Statement** We acknowledge the ICLR Code of Ethics. This submission is mostly theoretical and the authors could not think of any potential violations of them in this submission.

**Reproducibility Statement** For our theoretical results, explanations of assumptions can be found in Appendix D and a complete proof of the claims can be found in the Appendix E and Appendix F. Our experiment details can be found in Appendix G, with both data and training descriptions.

## REFERENCES

- N. Aronszajn. Theory of reproducing kernels. *Trans. Amer. Math. Soc.*, 68(3):337–404, 1950.
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. *arXiv preprint arXiv:1901.08584*, 2019.
- Jean-Yves Audibert and Alexandre B Tsybakov. Fast learning rates for plug-in classifiers. *The Annals of statistics*, 35(2):608–633, 2007.
- Peter L Bartlett, Michael I Jordan, and Jon D McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006.
- Colin Bennett and Robert C Sharpley. *Interpolation of Operators*. Academic press, 1988.
- Alberto Bietti and Julien Mairal. On the inductive bias of neural tangent kernels. *arXiv preprint arXiv:1905.12173*, 2019.
- Yuan Cao and Quanquan Gu. Generalization bounds of stochastic gradient descent for wide and deep neural networks. *Advances in Neural Information Processing Systems*, 32:10836–10846, 2019.
- Yuan Cao and Quanquan Gu. Generalization error bounds of gradient descent for learning over-parameterized deep relu networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 3349–3356, 2020.
- Lin Chen and Sheng Xu. Deep neural tangent kernel and Laplace kernel have the same RKHS. *arXiv preprint arXiv:2009.10683*, 2020.
- Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *International Conference on Machine Learning*, pp. 2206–2216. PMLR, 2020.
- Ahmet Demirkaya, Jiasi Chen, and Samet Oymak. Exploring the role of loss functions in multiclass classification. In *2020 54th Annual Conference on Information Sciences and Systems (CISS)*, pp. 1–5. IEEE, 2020.
- Gavin Weiguang Ding, Yash Sharma, Kry Yik Chau Lui, and Ruitong Huang. Mma training: Direct input space margin maximization through adversarial training. *arXiv preprint arXiv:1812.02637*, 2018.
- Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- David Eric Edmunds and Hans Triebel. *Function Spaces, Entropy Numbers, Differential Operators*, volume 120. Cambridge University Press, 2008.
- Gamaleldin F Elsayed, Dilip Krishnan, Hossein Mobahi, Kevin Regan, and Samy Bengio. Large margin deep networks for classification. *arXiv preprint arXiv:1803.05598*, 2018.
- Cong Fang, Hangfeng He, Qi Long, and Weijie J Su. Exploring deep neural networks via layer-peeled model: Minority collapse in imbalanced training. *Proceedings of the National Academy of Sciences*, 118(43), 2021.
- Max H Farrell, Tengyuan Liang, and Sanjog Misra. Deep neural networks for estimation and inference. *Econometrica*, 89(1):181–213, 2021.

- Amnon Geifman, Abhay Yadav, Yoni Kasten, Meirav Galun, David Jacobs, and Ronen Basri. On the similarity between the Laplace and neural tangent kernels. *arXiv preprint arXiv:2007.01580*, 2020.
- Florian Graf, Christoph Hofer, Marc Niethammer, and Roland Kwitt. Dissecting supervised contrastive learning. In *International Conference on Machine Learning*, pp. 3821–3830. PMLR, 2021.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International Conference on Machine Learning*, pp. 1321–1330. PMLR, 2017.
- XY Han, Vardan Papayan, and David L Donoho. Neural collapse under mse loss: Proximity to and dynamics on the central path. *arXiv preprint arXiv:2106.02073*, 2021.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- Zhezhi He, Adnan Siraj Rakin, and Deliang Fan. Parametric noise injection: Trainable randomness to improve deep neural network robustness against adversarial attack. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 588–597, 2019.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Tianyang Hu, Wenjia Wang, Cong Lin, and Guang Cheng. Regularization matters: A nonparametric perspective on overparametrized neural network. In *International Conference on Artificial Intelligence and Statistics*, pp. 829–837. PMLR, 2021.
- Wei Hu, Zhiyuan Li, and Dingli Yu. Simple and effective regularization methods for training on noisily labeled data with generalization guarantee. In *International Conference on Learning Representations*, 2020.
- Like Hui and Mikhail Belkin. Evaluation of neural architectures trained with square loss vs cross-entropy in classification tasks. *arXiv preprint arXiv:2006.07322*, 2020.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, pp. 8571–8580, 2018.
- Ziwei Ji and Matus Telgarsky. Polylogarithmic width suffices for gradient descent to achieve arbitrarily small test error with shallow ReLU networks. *arXiv preprint arXiv:1909.12292*, 2019.
- Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *arXiv preprint arXiv:2004.11362*, 2020.
- Yongdai Kim, Ilsang Ohn, and Dongha Kim. Fast convergence rates of deep neural networks for classification. *arXiv preprint arXiv:1812.03599*, 2018.
- Michael Kohler and Adam Krzyzak. On the rate of convergence of local averaging plug-in classification rules under a margin condition. *IEEE Transactions on Information Theory*, 53(5):1735–1742, 2007.
- Simon Kornblith, Honglak Lee, Ting Chen, and Mohammad Norouzi. Demystifying loss functions for classification. 2020.
- Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In *Advances in Neural Information Processing Systems*, pp. 8157–8166, 2018.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.

- Enno Mammen and Alexandre B Tsybakov. Smooth discrimination analysis. *The Annals of Statistics*, 27(6):1808–1829, 1999.
- Aditya K Menon, Ankit Singh Rawat, Sashank Reddi, Seungyeon Kim, and Sanjiv Kumar. A statistical perspective on distillation. In *International Conference on Machine Learning*, pp. 7632–7642. PMLR, 2021.
- Youssef Mroueh, Tomaso Poggio, Lorenzo Rosasco, and Jean-Jacques Slotine. Multiclass learning with simplex coding. *arXiv preprint arXiv:1209.1360*, 2012.
- Vidya Muthukumar, Adhyayan Narang, Vignesh Subramanian, Mikhail Belkin, Daniel Hsu, and Anant Sahai. Classification vs regression in overparameterized regimes: Does the loss function matter? *arXiv preprint arXiv:2005.08054*, 2020.
- Atsushi Nitanda and Taiji Suzuki. Optimal rates for averaged stochastic gradient descent under neural tangent kernel regime. *arXiv preprint arXiv:2006.12297*, 2020.
- Atsushi Nitanda, Geoffrey Chinot, and Taiji Suzuki. Gradient descent can learn less over-parameterized two-layer neural networks on classification problems. *arXiv preprint arXiv:1905.09870*, 2019.
- Tianyu Pang, Kun Xu, Yinpeng Dong, Chao Du, Ning Chen, and Jun Zhu. Rethinking softmax cross-entropy loss for adversarial robustness. *arXiv preprint arXiv:1905.10626*, 2019.
- Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):24652–24663, 2020.
- Tomaso Poggio and Qianli Liao. Generalization in deep network classifiers trained with the square loss. Technical report, CBMM Memo No, 2019.
- Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with relu activation function. *The Annals of Statistics*, 48(4):1875–1897, 2020.
- Ingo Steinwart, Clint Scovel, et al. Fast rates for support vector machines using Gaussian kernels. *The Annals of Statistics*, 35(2):575–607, 2007.
- Juozas Vaicenavicius, David Widmann, Carl Andersson, Fredrik Lindsten, Jacob Roll, and Thomas Schön. Evaluating model calibration in classification. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 3459–3467. PMLR, 2019.
- Holger Wendland. *Scattered Data Approximation*. Cambridge University Press, 2004.
- Yaodong Yu, Kwan Ho Ryan Chan, Chong You, Chaobing Song, and Yi Ma. Learning diverse and discriminative representations via the principle of maximal coding rate reduction. *Advances in Neural Information Processing Systems*, 33, 2020.
- Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- Tong Zhang. Statistical behavior and consistency of classification methods based on convex risk minimization. *The Annals of Statistics*, pp. 56–85, 2004.

## A GRADIENT DESCENT AND NEURAL TANGENT KERNEL

**Gradient Descent** Since we consider the square loss and  $\ell_2$  regularization, the optimization problem in Equation 2.2 becomes

$$\min_{\mathbf{W}} \sum_{i=1}^n (f_{\mathbf{W},\mathbf{a}}(\mathbf{x}_i) - y_i)^2 + \mu \|\mathbf{W}\|_2^2. \quad (\text{A.1})$$

We consider the GD training of Equation A.1. Let

$$\Phi(\mathbf{W}) = \sum_{i=1}^n (f_{\mathbf{W},\mathbf{a}}(\mathbf{x}_i) - y_i)^2 + \mu \|\mathbf{W}\|_2^2$$

be the objective function in Equation A.1. The gradient of  $\Phi$  with respect to  $\mathbf{w}_r$  can be written as (Arora et al., 2019)

$$\frac{\partial \Phi(\mathbf{W})}{\partial \mathbf{w}_r} = \frac{2}{\sqrt{m}} a_r \sum_{i=1}^n (u_i - y_i) \mathbb{I}_{r,i} \mathbf{x}_i + 2\mu \mathbf{w}_r, \quad r \in [m],$$

where  $u_i = f_{\mathbf{W},\mathbf{a}}(\mathbf{x}_i)$  and  $\mathbb{I}_{r,i} = \mathbb{I}\{\mathbf{w}_r^\top \mathbf{x}_i \geq 0\}$ . Then the GD update rule is

$$\mathbf{w}_r(k+1) = \mathbf{w}_r(k) - \zeta \frac{\partial \Phi(\mathbf{W})}{\partial \mathbf{w}_r} \Big|_{\mathbf{W}=\mathbf{W}(k)},$$

where  $\mathbf{W}(k)$  is the weight matrix at iteration  $k$ , and  $\zeta$  is the learning rate. Define  $\mathbb{I}_{r,i}(k) = \mathbb{I}\{\mathbf{w}_r(k)^\top \mathbf{x}_i \geq 0\}$ ,  $\mathbf{Z}(k) \in \mathbb{R}^{md \times n}$  as

$$\mathbf{Z}(k) = \frac{1}{\sqrt{m}} \begin{pmatrix} a_1 \mathbb{I}_{1,1}(k) \mathbf{x}_1 & \dots & a_1 \mathbb{I}_{1,n}(k) \mathbf{x}_n \\ \vdots & \dots & \vdots \\ a_m \mathbb{I}_{m,1}(k) \mathbf{x}_1 & \dots & a_m \mathbb{I}_{m,n}(k) \mathbf{x}_n \end{pmatrix},$$

$\mathbf{H}(k) = \mathbf{Z}(k)^\top \mathbf{Z}(k)$ , and  $\mathbf{u}(k) = (\mathbf{u}_1(k), \dots, \mathbf{u}_n(k))^\top$  with  $\mathbf{u}_i(k) = f_{\mathbf{W}(k),\mathbf{a}}(\mathbf{x}_i)$ . Then the GD update rule with respect to  $\mathbf{W}$  can be written as

$$\text{vec}(\mathbf{W}(k+1)) = \text{vec}(\mathbf{W}(k)) - 2\zeta (\mathbf{Z}(k)(\mathbf{u}(k) - \mathbf{y}) + \mu \text{vec}(\mathbf{W}(k))), \quad (\text{A.2})$$

where  $\text{vec}(\mathbf{W}) = (\mathbf{w}_1^\top, \dots, \mathbf{w}_m^\top)^\top \in \mathbb{R}^{md \times 1}$  is the vectorized weight matrix and  $\mathbf{y} = (y_1, \dots, y_n)^\top$ .

**Neural Tangent Kernel (NTK)** It has been shown that the following NTK

$$h(\mathbf{s}, \mathbf{t}) = \mathbb{E}_{\mathbf{w} \sim N(0, \mathbf{I}_d)} (\mathbf{s}^\top \mathbf{t} \mathbb{I}\{\mathbf{w}^\top \mathbf{s} \geq 0, \mathbf{w}^\top \mathbf{t} \geq 0\}) = \frac{\mathbf{s}^\top \mathbf{t} (\pi - \arccos(\mathbf{s}^\top \mathbf{t}))}{2\pi} \quad (\text{A.3})$$

plays an important role in the study of one-hidden-layer ReLU neural networks, where  $\mathbf{s}, \mathbf{t}$  are  $d$ -dimensional vectors (Du et al., 2018; Hu et al., 2021). Since  $h$  is positive definite on the unit sphere  $\mathbb{S}^{d-1}$  (Bietti & Mairal, 2019), by Mercer's theorem, it possesses a Mercer decomposition as  $h(\mathbf{s}, \mathbf{t}) = \sum_{j=0}^{\infty} \lambda_j \varphi_j(\mathbf{s}) \varphi_j(\mathbf{t})$ , where  $\lambda_1 \geq \lambda_2 \geq \dots \geq 0$  are the eigenvalues, and  $\{\varphi_j\}_{j=1}^{\infty}$  is an orthonormal basis. The asymptotic behavior of the eigenvalues is described in the following lemma.

**Lemma A.1** (Lemma 3.1 of Hu et al. (2021)). Let  $\lambda_j$  be the eigenvalues of NTK  $h$  defined above. Then we have  $\lambda_j \asymp j^{-\frac{d}{d-1}}$ .

Let  $\mathcal{N}$  denote the reproducing kernel Hilbert space (RKHS) generated by  $h$  on  $\mathbb{S}^{d-1}$ , equipped with norm  $\|\cdot\|_{\mathcal{N}}$ . As a corollary of Lemma A.1, it can be shown that the  $(L_\infty)$  entropy number of a unit ball in  $\mathcal{N}$ , denoted by  $\mathcal{N}(1)$ , can be controlled. The relationship between the eigenvalues and entropy numbers has been well studied; see Edmunds & Triebel (2008).

**Lemma A.2.** The entropy number of  $\mathcal{N}(1)$ , denoted by  $H(\mathcal{N}(1), \delta, \|\cdot\|_{L_\infty})$ , is bounded by  $H(\mathcal{N}(1), \delta, \|\cdot\|_{L_\infty}) \lesssim \delta^{-2(d-1)/d}$ .

There are extensive works studying the generalization error bounds under NTK regime. For regression, Nitanda & Suzuki (2020); Hu et al. (2021) show the optimal convergence rates when using overparametrized one-hidden-layer neural networks, where the square loss is used. Arora et al. (2019) provides generalization error bounds and provable learning scenarios for noiseless data. In the NTK regime, the neural network as a regressor is linked with the nonparametric regression via NTK. There are also other works studying the generalization performance of the neural network as a nonparametric regressor, out of the NTK regime; see Schmidt-Hieber (2020); Farrell et al. (2021).

For classification, most of the existing results are established based on the separable data; see Ji & Telgarsky (2019); Cao & Gu (2019); Nitanda et al. (2019) and references therein. In particular, Hu et al. (2020) consider classification with noisy labels (labels are randomly flipped) and propose to use the square loss with  $\ell_2$  regularization.

Besides the generalization error bounds, another important research direction is to bridge the gap between NTK and finite-width overparametrized neural networks via GD training; see Du et al. (2018); Arora et al. (2019); Li & Liang (2018); Hu et al. (2021), among others.

## B OVERVIEW OF REPRODUCING KERNEL HILBERT SPACE

We provide here a brief overview of reproducing kernel Hilbert space (RKHS).

**Definition B.1** (Positive Definite Kernel). A function  $k : \Omega \times \Omega \mapsto \mathbb{R}$  is said to be a *positive definite kernel*, if  $k(\mathbf{x}, \tilde{\mathbf{x}}) = k(\tilde{\mathbf{x}}, \mathbf{x})$  for all  $\mathbf{x}, \tilde{\mathbf{x}} \in \Omega$ , and

$$\sum_{i=1}^n \sum_{j=1}^n \beta_i \beta_j k(\mathbf{x}_i, \mathbf{x}_j) > 0,$$

for all  $n \in \mathbb{N}$ ,  $\beta_1, \dots, \beta_n \in \mathbb{R}$  such that at least one  $\beta_j \neq 0$ , and  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \Omega$ .

For a positive definite kernel  $k$ , define a linear space

$$\mathcal{N}_k^0 := \left\{ \sum_{i=1}^n \beta_i k(\cdot, \mathbf{x}_i) : n \in \mathbb{N}, \beta_1, \dots, \beta_n \in \mathbb{R}, \mathbf{x}_1, \dots, \mathbf{x}_n \in \Omega \right\},$$

and equip this space with an inner product  $\langle \cdot, \cdot \rangle_{\mathcal{N}_k^0}$  by

$$\left\langle \sum_{i=1}^n \beta_i k(\cdot, \mathbf{x}_i), \sum_{j=1}^{\tilde{n}} \tilde{\beta}_j k(\cdot, \tilde{\mathbf{x}}_j) \right\rangle_{\mathcal{N}_k^0} := \sum_{i=1}^n \sum_{j=1}^{\tilde{n}} \beta_i \tilde{\beta}_j k(\mathbf{x}_i, \tilde{\mathbf{x}}_j).$$

The norm of  $g \in \mathcal{N}_k^0$  is defined by  $\|g\|_{\mathcal{N}_k^0}^2 := \langle g, g \rangle_{\mathcal{N}_k^0}$ . Then the RKHS induced by  $k$ , denoted by  $\mathcal{N}_k(\Omega)$ , is defined as the closure of  $\mathcal{N}_k^0(\Omega)$  with respect to the norm  $\|\cdot\|_{\mathcal{N}_k(\Omega)}$ .

For a subset  $\Omega_0 \subset \Omega$ , define the *restriction* of  $\mathcal{N}_k$  on  $\Omega_0$  as

$$\mathcal{N}_k(\Omega_0) := \{g : \Omega_0 \mapsto \mathbb{R} : g = h|_{\Omega_0} \text{ for some } h \in \mathcal{N}_k\},$$

where  $g = h|_{\Omega_0}$  means  $g(\mathbf{x}) = h(\mathbf{x})$  for all  $\mathbf{x} \in \Omega_0$ . We equip  $\mathcal{N}_k(\Omega_0)$  with norm

$$\|g\|_{\mathcal{N}_k(\Omega_0)} := \inf_{\{h \in \mathcal{N}_k : h|_{\Omega_0} = g\}} \|h\|_{\mathcal{N}_k}.$$

Then,  $\mathcal{N}_k(\Omega_0)$  is a RKHS with norm  $\|\cdot\|_{\mathcal{N}_k(\Omega_0)}$  (see Aronszajn 1950, page 351).

## C SIMPLEX COORDINATES

In simplex label coding, the one-hot labels are replaced by the simplex vertices of a  $(K-1)$ -simplex. The vertices of a regular  $(K-1)$ -simplex centered on the origin can be written as:

$$\mathbf{v}_0 = \frac{1}{\sqrt{2K}} \cdot (1, \dots, 1)$$

and for  $1 \leq i \leq K - 1$ ,

$$\mathbf{v}_i = \frac{1}{\sqrt{2}} \mathbf{e}_i - \frac{1}{(K-1)\sqrt{2}} \left(1 + \frac{1}{\sqrt{K}}\right) \cdot (1, \dots, 1).$$

The pairwise angle between vertices is  $\arccos(-1/(K-1))$  and as  $K \rightarrow \infty$ , the angle converges to  $90^\circ$ .

The vertices of a  $(K-1)$ -simplex can be viewed as maximally separated  $K$  points on a sphere. In theory, the radius of the sphere doesn't matter but in practice, we recommend scaling it for larger number of classes, e.g., radius =  $K$  for  $K$ -class classification. We find that such scaling empirically outperforms the default radius 1 in our experiments. More details can be found in Appendix G.2.

## D ASSUMPTIONS

In this work, we impose the following assumptions. In the rest of the Appendix, we use  $\text{poly}(t_1, t_2, \dots)$  to denote some polynomial function with arguments  $t_1, t_2, \dots$ .

**Assumption D.1.** Let  $\lambda_{\min}(\mathbf{H}^\infty)$  be the minimum eigenvalue of the symmetric matrix  $\mathbf{H}^\infty$ , where  $\mathbf{H}^\infty = (h(\mathbf{x}_i, \mathbf{x}_j))_{n \times n}$  ( $\mathbf{H}^\infty$  is usually called the NTK matrix). Let  $\lambda_0$  be the largest number such that with probability at least  $1 - \delta_n$ ,  $\lambda_{\min}(\mathbf{H}^\infty) \geq \lambda_0$ , and  $\delta_n \rightarrow 0$  as  $n$  goes to infinity<sup>3</sup>. For sufficiently large  $n$ , the regularization parameter  $\mu \asymp n^{\frac{d-1}{2d-1}}$ , the learning rate  $\zeta = o(n^{-\frac{3d-1}{2d-1}})$ , the variance of initialization  $\xi^2 = O(1)$ , the number of nodes in the hidden layer  $m \geq \xi^{-2} \text{poly}(n, \lambda_0^{-1})$ , and the iteration number  $k$  satisfies  $\log(\text{poly}_1(n, \xi, 1/\lambda_0)) \lesssim \zeta \mu k \lesssim \log(\text{poly}_2(\xi, 1/n, \sqrt{m}))$ .

**Assumption D.2.** The conditional probability in the non-separable case satisfies  $\eta \in \mathcal{N}$ .

**Assumption D.3.** The solution to Equation 2.2 satisfies  $\|f_{\mathbf{W}^{(k)}, \mathbf{a}}\|_{\mathcal{N}} \leq C$ , where  $C$  is a constant not depending on  $n$ .

**Remark 5.** Assumption D.3 can be replaced by a stronger assumption, that is,  $f_{\mathbf{W}^{(k)}, \mathbf{a}}$  has a bounded Lipschitz constant, and the constant does not depend on  $n$ .

**Assumption D.4.** The probability density function of the marginal distribution  $P_X$ , denoted by  $p(\mathbf{x})$ , is continuous on  $\Omega$ , and there exists a positive constant  $c_0$  such that

$$p(\mathbf{x}) \leq c_0, \forall \mathbf{x} \in \Omega.$$

**Assumption D.5.** The probability density function of the marginal distribution  $p(\mathbf{x})$  is continuous on  $\Omega$ , and there exist positive constants  $c_1 \leq c_2$  such that

$$c_1 \leq p(\mathbf{x}) \leq c_2, \forall \mathbf{x} \in \Omega.$$

Assumption D.1 is related to the neural network and GD training, where similar settings have been adopted by Arora et al. (2019); Hu et al. (2021). From the results in Arora et al. (2019); Hu et al. (2021), the width of the neural network depends on the minimum eigenvalue of the NTK matrix  $\lambda_{\min}(\mathbf{H}^\infty)$ , where a smaller  $\lambda_{\min}(\mathbf{H}^\infty)$  leads to a wider neural network. Therefore, it is desired that  $\lambda_0$  is as large as possible. However, the consistency requires that the probability is tending to one; thus, we require  $\delta_n \rightarrow 0$  as  $n \rightarrow \infty$ . As  $n$  becomes larger, with probability tending to one, the distance of the two nearest points in  $n$  input points converges to zero, thus making  $\mathbf{H}^\infty$  close to a degenerate matrix, and the minimum eigenvalue of  $\mathbf{H}^\infty$  converges to zero. Therefore, inevitably,  $\lambda_0 \rightarrow 0$  (but  $\lambda_{\min}(\mathbf{H}^\infty)$  is strictly larger than 0 for all  $n$  with probability one). The requirements of the regularization parameter, the learning rate, the variance of initialization, the number of nodes in the hidden layer and the iteration number are all the same as those in Hu et al. (2021).

Assumption D.2 imposes conditions on the underlying true conditional probability in the *non-separable* case. This assumption basically requires that the conditional probability is within the function class generated by the GD-trained neural networks we consider (thus can be calibrated). Given that the neural networks are highly flexible, we believe that most of the functions are within the function class generated by the neural networks. As a simple example, any Lipschitz functions are within this function class.

<sup>3</sup>Potential dependency of  $\lambda_0$  on  $n$  is suppressed for notational simplicity.

Assumption D.3 requires that the solution to Equation 2.2 is well-behaved, i.e., the solution is within a ball in  $\mathcal{N}$  with a certain radius. Roughly speaking, Assumption D.3 requires that the *complexity* of the neural network estimator generated by the GD training is controlled. Since the step size is relatively small and the iteration number is not large (only  $\log(\text{poly}(n, \xi, \lambda_0^{-1}))$ ), we believe it is a mild assumption.

Assumption D.4 only requires the probability to be upper bounded from infinity, while Assumption D.5 requires the probability to be upper bounded from infinity and lower bounded away from zero on the support  $\Omega$ . They are standard assumptions used in the classical analysis of classification in statistics; see Audibert & Tsybakov (2007); Kohler & Krzyzak (2007) for example. Clearly, uniform distribution satisfies Assumptions D.4 and D.5. In Audibert & Tsybakov (2007), Assumption D.4 is called mild density assumption and Assumption D.5 is called strong density assumption.

## E PROOFS OF MAIN RESULTS

This section includes the proofs of main results in the paper.

### E.1 PROOF OF THEOREM 3.1

We first introduce some lemmas that are used in the proof of Theorem 3.1.

Let  $l_1(y_i, f(\mathbf{x}_i)) = (1 - y_i f(\mathbf{x}_i))^2 = (y_i - f(\mathbf{x}_i))^2$  be the square loss on a training point  $(\mathbf{x}_i, y_i)$ , the  $l_1$ -risk of  $f$  be  $\mathcal{R}_{l_1}(f) = \mathbb{E}_{X, Y \sim P} l_1(Y, f(X))$ , and  $\mathcal{R}_{l_1} = \min_{f \in \mathcal{N}} \mathbb{E}_{X, Y \sim P} l_1(Y, f(X))$ . Let  $L_1(f, \mathbf{x}, y) = \mu \|f\|_{\mathcal{N}}^2 + l_1(y, f(\mathbf{x}))$  and the  $L_1$ -risk of  $f$  be  $\mathcal{R}_{L_1}(f) = \mathbb{E}_{X, Y \sim P} L_1(f, X, Y)$ . Let  $f_n = \arg \min_{f \in \mathcal{N}} \mathcal{R}_{L_1}(f)$ .

Lemma E.1 is (a weaker version of) Theorem 5.6 of Steinwart et al. (2007), which provides a bound on the deviation between the empirical minimizer and true minimizer. Lemma E.2 is used to verify that one of the conditions of Lemma E.1 is fulfilled. Lemma E.3 shows that under certain conditions, the solution to

$$\min_{f \in \mathcal{N}} \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \frac{\mu}{n} \|f\|_{\mathcal{N}}^2 \quad (\text{E.1})$$

is closely related to the estimator given by the overparameterized neural networks  $f_{\mathbf{W}^{(k)}, \mathbf{a}}$ . Lemma E.3 can be obtained by merely repeating the proof of Theorem 5.2 of Hu et al. (2021), since we require that the probability density function  $p(\mathbf{x})$  of  $P_X$  is upper bounded by a positive constant by Assumption D.4. Therefore, the only difference is that we replace  $\|\cdot\|_2$  (which corresponds to the uniform distribution) to the  $L_2$  norm corresponding to the probability measure  $P_X$ ; thus the proof is omitted. Note also that the second statement of Lemma E.3 corresponds to the noiseless case.

**Lemma E.1.** Let  $Z = \Omega \times \{-1, 1\}$ . Let  $\mathcal{F}$  be a convex set of bounded measurable functions from  $Z$  to  $\mathbb{R}$  and let  $L : \mathcal{F} \times Z \rightarrow [0, \infty)$  be a convex and continuous loss function. For a probability measure  $P$  on  $Z$ , define

$$\mathcal{G} := \{L \circ f - L \circ f_{P, \mathcal{F}} : f \in \mathcal{F}\},$$

where  $f_{P, \mathcal{F}}$  is a minimizer of  $\mathbb{E}_{Z \sim P} L(f, Z)$ . Suppose that there are constants  $c \geq 0$ ,  $0 < \alpha \leq 1$ ,  $\delta \leq 0$  and  $B > 0$  such that  $\mathbb{E}_{Z \sim P} g^2 \leq c(\mathbb{E}_{Z \sim P} g)^\alpha + \delta$  and  $\|g\|_\infty \leq B$  for all  $g \in \mathcal{G}$ . Furthermore, assume that  $\mathcal{G}$  is separable with respect to  $\|\cdot\|_\infty$  and that there are constants  $a \geq 1$  and  $0 < \alpha < 2$  with

$$\sup_{T \in Z^n} H(B^{-1} \mathcal{G}, \epsilon, \|\cdot\|_{L_2(T)}) \leq a \epsilon^{-\beta} \quad (\text{E.2})$$

for all  $\epsilon > 0$ , where  $H(B^{-1} \mathcal{G}, \epsilon, \|\cdot\|_{L_2(T)})$  is the entropy number of the set  $B^{-1} \mathcal{G}$ , and  $\|f\|_{L_2(T)}^2 = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i, y_i)^2$  is the empirical norm. Then there exists a constant  $c_\beta > 0$  depending only on  $\beta$  such that for all  $n \geq 1$  and all  $t \geq 1$  we have

$$\mathbb{P}(T \in Z^n : \mathcal{R}_{L, P}(f_{T, \mathcal{F}}) > \mathcal{R}_{L, P}(f_{P, \mathcal{F}}) + c_\beta \varepsilon(n, a, B, c, \delta, t)) \leq e^{-t},$$



where

$$\begin{aligned} & \varepsilon(n, a, B, c, \delta, t) \\ &= B^{2\beta/(4-2\alpha+\alpha\beta)} c^{(2-\beta)/(4-2\alpha+\alpha\beta)} \left(\frac{a}{n}\right)^{2/(4-2\alpha+\alpha\beta)} + B^{\beta/2} \delta^{(2-\beta)/4} \left(\frac{a}{n}\right)^{1/2} \\ &+ B \left(\frac{a}{n}\right)^{2/(2+\beta)} + \sqrt{\frac{\delta t}{n}} + \left(\frac{ct}{n}\right)^{1/(2-\alpha)} + \frac{Bt}{n}, \end{aligned} \quad (\text{E.3})$$

and  $f_{T,\mathcal{F}}$  is the minimizer with respect to the empirical measure.

**Lemma E.2.** Assume the conditions of Theorem 3.1 hold. Define  $C := 8\|(2\eta - 1)^{-1}\|_{\kappa,\infty} + 32$ , where  $\|\cdot\|_{\kappa,\infty}$  is the norm of Lorentz space  $L_{\kappa,\infty}$  (Bennett & Sharpley, 1988). Let  $\mu > 0$  and  $0 < \gamma \leq n^{1/2}\mu^{-1/2}$ , then for all  $f \in \gamma\mathcal{N}(1)$  we have

$$\mathbb{E}_{X,Y \sim P}(L_1 \circ f - L_1 \circ f_n)^2 \leq C(K\gamma + 1)^2(\mathbb{E}_{X,Y \sim P}(L_1 \circ f - L_1 \circ f_n)) + 2C(K\gamma + 1)^2 a(\mu),$$

where  $a(\mu)$  is the approximation error function given by

$$a(\mu) = \inf_{f \in \mathcal{N}} (n^{-1}\mu \|f\|_{\mathcal{N}}^2 + \mathcal{R}_{l_1}(f) - \mathcal{R}_{l_1}).$$

**Lemma E.3.** Suppose Assumptions D.1 and D.4 hold. Then we have

$$\mathbb{E}_{X \sim P_X} (f_{\mathbf{W}^{(k)},\alpha}(X) - \hat{f}(X))^2 = O_{\mathbb{P}}(n^{-\frac{d}{2d-1}}),$$

where  $\hat{f}$  is the solution to Equation E.1. Furthermore, if there exists a function  $f \in \mathcal{N}$  that does not depend on  $n$  and  $f(\mathbf{x}_i) = y_i$  for all  $i = 1, \dots, n$ , then we can set  $\mu = o(1)$  and obtain

$$\mathbb{E}_{X \sim P_X} (f_{\mathbf{W}^{(k)},\alpha}(X) - \hat{f}(X))^2 = o_{\mathbb{P}}(1).$$

**Remark 6.** According to the proof in Hu et al. (2021), the probability in  $o_{\mathbb{P}}(1)$  of Lemma E.3 only relates to the width of the one-hidden-layer neural network, which can be arbitrarily small by enlarging the neural network's width.

Now we are ready to prove Theorem 3.1. Let  $L = L_1$  in Lemma E.1, which is clearly continuous. Let  $\hat{f}$  be the solution to Equation E.1. The key idea in this proof is using  $\hat{f}$  to bridge two functions  $2\eta - 1$  and  $f_{\mathbf{W}^{(k)},\alpha}$ .

Since  $\hat{f}$  is the solution to Equation E.1, it can be seen that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(\mathbf{x}_i))^2 + \frac{\mu}{n} \|\hat{f}\|_{\mathcal{N}}^2 &\leq \frac{1}{n} \sum_{i=1}^n (y_i - (2\eta(\mathbf{x}_i) - 1))^2 + \frac{\mu}{n} \|2\eta(\mathbf{x}_i) - 1\|_{\mathcal{N}}^2 \\ &\leq \frac{2}{n} \sum_{i=1}^n (y_i^2 + (2\eta(\mathbf{x}_i) - 1)^2) + \frac{\mu}{n} \|2\eta(\mathbf{x}_i) - 1\|_{\mathcal{N}}^2 \\ &\leq C_1, \end{aligned} \quad (\text{E.4})$$

where the second inequality is by the Cauchy-Schwarz inequality, and the third inequality is because  $y_i^2 = 1$  and  $\eta(\mathbf{x})$  is bounded.

The reproducing property implies that

$$\hat{f}(\mathbf{x}) = \langle \hat{f}, h(\mathbf{x}, \cdot) \rangle_{\mathcal{N}} \leq \|\hat{f}\|_{\mathcal{N}} \|h(\mathbf{x}, \cdot)\|_{\mathcal{N}} = \|\hat{f}\|_{\mathcal{N}} \sqrt{h(\mathbf{x}, \mathbf{x})}, \forall \mathbf{x} \in \Omega,$$

which yields

$$\|\hat{f}\|_{L_{\infty}} \leq C_2 \|\hat{f}\|_{\mathcal{N}}.$$

Together with Equation E.4, we obtain

$$\|\hat{f}\|_{L_{\infty}} \leq C_3 \|\hat{f}\|_{\mathcal{N}} \leq C_4 (\mu/n)^{-1/2}. \quad (\text{E.5})$$

Thus, we can take  $B = C_4(\mu/n)^{-1/2}$  in Lemma E.1. The entropy condition can be verified via Lemma A.2, which allows us to take  $\beta = 2(d-1)/d$ . Equation E.5, together with Lemma E.2, also suggests that we can take  $c = C(KB + 1)^2$ ,  $\alpha = 1$ , and  $\delta = 2C(KB + 1)^2 a(\mu)$ .

Next, we provide an upper bound on  $a(\mu)$ . The definition of  $a(\mu)$  implies

$$\begin{aligned} a(\mu) &= n^{-1} \mu \|f_n\|_{\mathcal{N}}^2 + R_{l_1}(f_n) - R_{l_1} \\ &= n^{-1} \mu \|f_n\|_{\mathcal{N}}^2 + \mathbb{E}_{X \sim P_X} (2\eta(X) - 1 - f_n(X))^2 \\ &\leq n^{-1} \mu \|2\eta - 1\|_{\mathcal{N}}^2, \end{aligned} \quad (\text{E.6})$$

where we use the relationship  $\mathcal{R}_{l_1, P}(f) - \mathcal{R}_{l_1, P} = \mathbb{E}_{X \sim P_X} (2\eta(X) - 1 - f(X))^2$ .

Plugging all the terms into Equation E.3, together with Lemma E.1, yields that

$$\mathcal{R}_{L_1, P}(\hat{f}) = \mathcal{R}_{L_1, P}(f_n) + O_{\mathbb{P}}(\varepsilon(n, a, B, c, \delta)), \quad (\text{E.7})$$

where

$$\begin{aligned} \varepsilon(n, a, B, c, \delta) &= B^{\frac{4}{2+\beta}} n^{-\frac{2}{2+\beta}} + B(\mu/n)^{\frac{2-\beta}{4}} n^{-\frac{1}{2}} \|2\eta - 1\|_{\mathcal{N}}^{\frac{2-\beta}{2}} + B^2 n^{-1} \\ &= B^{\frac{4d}{4d-2}} n^{-\frac{2d}{4d-2}} + B\mu^{\frac{1}{2d}} n^{-\frac{1}{2} - \frac{1}{2d}} \|2\eta - 1\|_{\mathcal{N}}^{\frac{1}{d}} + B^2 n^{-1}. \end{aligned} \quad (\text{E.8})$$

Since  $\mathcal{R}_{l_1, P}(f) - \mathcal{R}_{l_1, P} = \mathbb{E}_{X \sim P_X} (2\eta(X) - 1 - f(X))^2$ , we subtract  $\mathcal{R}_{l_1, P}$  on both sides of Equation E.7 and get

$$\begin{aligned} &\mathbb{E}_{X \sim P_X} (2\eta(X) - 1 - \hat{f}(X))^2 + n^{-1} \mu \|\hat{f}\|_{\mathcal{N}}^2 \\ &= \mathbb{E}_{X \sim P_X} (2\eta(X) - 1 - f_n(X))^2 + n^{-1} \mu \|f_n\|_{\mathcal{N}}^2 + O_{\mathbb{P}}(\varepsilon(n, a, B, c, \delta)) \\ &= O_{\mathbb{P}}(n^{-1} \mu \|2\eta - 1\|_{\mathcal{N}}^2 + \varepsilon(n, a, B, c, \delta)), \end{aligned} \quad (\text{E.9})$$

where the last equality (with big  $O$  notation) is by Equation E.6. Combining Equation E.9 and Equation E.5 implies

$$\|\hat{f}\|_{L_{\infty}}^2 \leq C_3^2 \|\hat{f}\|_{\mathcal{N}}^2 = O_{\mathbb{P}}(1 + n\mu^{-1} \varepsilon(n, a, B, c, \delta)).$$

In the following, we will show that by taking  $\mu \asymp n^{\frac{d-1}{2d-1}}$ ,

$$\mathbb{E}_{X \sim P_X} (2\eta(X) - 1 - \hat{f}(X))^2 + n^{-1} \mu \|\hat{f}\|_{\mathcal{N}}^2 = O_{\mathbb{P}}(n^{-\frac{d}{2d-1}} \max(1, \|2\eta - 1\|_{\mathcal{N}}^{\frac{2}{d}})) = O_{\mathbb{P}}(n^{-\frac{d}{2d-1}}). \quad (\text{E.10})$$

If  $n\mu^{-1} \varepsilon(n, a, B, c, \delta) \lesssim 1$ , then  $\varepsilon(n, a, B, c, \delta) \lesssim \mu/n$ , and Equation E.10 holds. Otherwise, we can replace  $B^2$  by its upper bound  $O_{\mathbb{P}}(n\mu^{-1} \varepsilon(n, a, B, c, \delta))$  in Equation E.8 and obtain that

$$\varepsilon = O_{\mathbb{P}}(\varepsilon^{\frac{d}{2d-1}} \mu^{-\frac{d}{2d-1}} + \varepsilon^{\frac{1}{2}} \mu^{-\frac{d-1}{2d}} n^{-\frac{1}{2d}} \|2\eta - 1\|_{\mathcal{N}}^{\frac{1}{d}}),$$

where we set  $\varepsilon = \varepsilon(n, a, B, c, \delta)$  for notational simplicity. Let us hereby denote  $I_1 = \varepsilon^{\frac{d}{2d-1}} \mu^{-\frac{d}{2d-1}}$  and  $I_2 = \varepsilon^{\frac{1}{2}} \mu^{-\frac{d-1}{2d}} n^{-\frac{1}{2d}} \|2\eta - 1\|_{\mathcal{N}}^{\frac{1}{d}}$ , and consider the following two cases.

**Case 1:**  $I_1 \geq I_2$ , then we have

$$\varepsilon = O_{\mathbb{P}}(\varepsilon^{\frac{d}{2d-1}} \mu^{-\frac{d}{2d-1}}).$$

Solving this equality leads to

$$\varepsilon = O_{\mathbb{P}}(\mu^{-\frac{d}{d-1}}). \quad (\text{E.11})$$

Plugging Equation E.11 into Equation E.9 and minimize the right-hand side of Equation E.9 with respect to  $\mu$  gives us  $\mu \asymp n^{\frac{d-1}{2d-1}}$ ; thus Equation E.10 holds.

**Case 2:**  $I_1 < I_2$ , then we have

$$\varepsilon = O_{\mathbb{P}}(\varepsilon^{\frac{1}{2}} \mu^{-\frac{d-1}{2d}} n^{-\frac{1}{2d}} \|2\eta - 1\|_{\mathcal{N}}^{\frac{1}{d}}),$$

which leads to

$$\varepsilon = O_{\mathbb{P}}(\mu^{-\frac{d-1}{d}} n^{-\frac{1}{d}} \|2\eta - 1\|_{\mathcal{N}}^{\frac{2}{d}}). \quad (\text{E.12})$$

Similarly, we plug Equation E.12 into Equation E.9 and minimize the right-hand side of Equation E.9 with respect to  $\mu$  and obtain  $\mu \asymp n^{\frac{d-1}{2d-1}}$ , which also leads to Equation E.10.

Now we can obtain an upper bound on the excess risk. For the notation simplicity, let  $f = f_{\mathbf{W}(k), \mathbf{a}}$ . The excess risk can be bounded by

$$L(f) - L^* \leq \mathbb{E}_{X \sim P_X} \mathbb{I}\{(2\eta(X) - 1)f(X) \leq 0, |\eta(X) - 0.5| < \delta\} |2\eta(X) - 1| \\ + \mathbb{E}_{X \sim P_X} \mathbb{I}\{(2\eta(X) - 1)f(X) \leq 0, |\eta(X) - 0.5| \geq \delta\} |2\eta(X) - 1|. \quad (\text{E.13})$$

The first term can be bounded via Tsybakov's noise condition as

$$\mathbb{E}_{X \sim P_X} \mathbb{I}\{(2\eta(X) - 1)f(X) \leq 0, |\eta(X) - 0.5| < \delta\} |2\eta(X) - 1| \leq 2\delta \mathbb{E}[\mathbb{I}\{|\eta(X) - 0.5| < \delta\}] \\ = 2\delta \mathbb{P}(|\eta(X) - 0.5| < \delta) \leq 2C\delta^{\kappa+1}. \quad (\text{E.14})$$

It remains to bound the second term in Equation E.13. If  $p(\mathbf{x})$  is continuous, then by the fact that  $|2\eta(X) - 1| \leq |2\eta(X) - 1 - f(X)|$  if  $(2\eta(X) - 1)f(X) \leq 0$ , we have

$$\mathbb{E}_{X \sim P_X} \mathbb{I}\{(2\eta(X) - 1)f(X) \leq 0, |\eta(X) - 0.5| \geq \delta\} |2\eta(X) - 1| \\ \leq 2\delta^{-1} \mathbb{E}_{X \sim P_X} \mathbb{I}\{(2\eta(X) - 1)f(X) \leq 0, |\eta(X) - 0.5| \geq \delta\} |2\eta(X) - 1|^2 \\ \leq 2\delta^{-1} \mathbb{E}_{X \sim P_X} \mathbb{I}\{|\eta(X) - 0.5| \geq \delta\} (2\eta(X) - 1 - f(X))^2 \\ \leq 2\delta^{-1} \mathbb{E}_{X \sim P_X} (2\eta(X) - 1 - f(X))^2 \\ \leq 4\delta^{-1} \mathbb{E}_{X \sim P_X} (2\eta(X) - 1 - \hat{f}(X))^2 + 4\delta^{-1} \mathbb{E}_{X \sim P_X} (f(X) - \hat{f}(X))^2 \\ = O_{\mathbb{P}}(\delta^{-1} n^{-\frac{d}{2d-1}}), \quad (\text{E.15})$$

where the fourth inequality is by the Cauchy-Schwarz inequality, and the last equality (with big  $O$  notation) is by Equation E.10 and Lemma E.3. Taking  $\delta = n^{-\frac{d}{(2d-1)(\kappa+2)}}$ , and plugging Equation E.14 and Equation E.15 into Equation E.13 leads to

$$L(f) = L^* + O_{\mathbb{P}}(n^{-\frac{d(\kappa+1)}{(2d-1)(\kappa+2)}}).$$

This finishes the proof.

## E.2 PROOF OF THEOREM 3.2

We first present a lemma.

**Lemma E.4.** Suppose two sets are separable with a positive margin  $\gamma > 0$ . Then there exists a function  $f_T$  satisfying

$$f_T(\mathbf{x}) = 1, \forall \mathbf{x} \in \Omega_1, \quad f_T(\mathbf{x}) = -1, \forall \mathbf{x} \in \Omega_2.$$

*Proof of Theorem 3.2.* By the equivalence of the RKHS generated by the Laplace kernel and  $\mathcal{N}$  (Geifman et al., 2020; Chen & Xu, 2020), it can be shown that  $\mathcal{N}$  can be embedded into the Sobolev space  $W_2^\nu$  for some  $\nu > d/2$ . Consider the Hölder space  $C_b^{0,\alpha}$  for  $0 < \alpha \leq 1$  equipped with the norm

$$\|f\|_{C_b^{0,\alpha}} := \sup_{\mathbf{x}, \mathbf{x}' \in \Omega, \mathbf{x} \neq \mathbf{x}'} \frac{|f(\mathbf{x}) - f(\mathbf{x}')|}{\|\mathbf{x} - \mathbf{x}'\|_2^\alpha}. \quad (\text{E.16})$$

By the Sobolev embedding theorem, we have the embedding relationship

$$\|f\|_{C_b^{0,\tau}} \leq C_1 \|f\|_{W_2^\nu} \leq C_2 \|f\|_{\mathcal{N}} \quad (\text{E.17})$$

for all  $f \in \mathcal{N}$ , where  $\tau = \min(\nu - d/2, 1)$ .

Without loss of generality, let us consider  $\mathbf{x} \in \Omega_1$ . The case of  $\mathbf{x} \in \Omega_2$  can be proved similarly. For any  $\mathbf{x} \in \Omega_1$ , take  $\mathbf{x}' = \arg \min_{\mathbf{x}_i} \|\mathbf{x}_i - \mathbf{x}\|_2$ . Thus, the definition of the Hölder space and Equation E.17 imply

$$|f_{\mathbf{W}(k), \mathbf{a}}(\mathbf{x}) - f_{\mathbf{W}(k), \mathbf{a}}(\mathbf{x}')| \leq C_2 \|f_{\mathbf{W}(k), \mathbf{a}}\|_{\mathcal{N}} \|\mathbf{x}' - \mathbf{x}\|_2^\tau \leq C_3 \|\mathbf{x}' - \mathbf{x}\|_2^\tau, \quad (\text{E.18})$$

where the last inequality is by Assumption D.3.

Let  $\hat{f}$  be the solution to Equation E.1, and let  $\hat{f}_1$  be the solution to Equation E.1 with  $\mu = 0$ . Let  $f_T$  be as in Lemma E.4. Note that  $\hat{f}_1$  satisfies  $\hat{f}_1(\mathbf{x}_i) = f_T(\mathbf{x}_i)$ . Thus, by the identity  $\|\hat{f}_1\|_{\mathcal{N}}^2 + \|\hat{f}_1 - f_T\|_{\mathcal{N}}^2 = \|\hat{f}_T\|_{\mathcal{N}}^2$  (Wendland, 2004), we have  $\|\hat{f}_1\|_{\mathcal{N}} \leq \|f_T\|_{\mathcal{N}}$ . Since  $\hat{f}$  is the solution to Equation E.1, we have

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(\mathbf{x}_i))^2 + \frac{\mu}{n} \|\hat{f}\|_{\mathcal{N}}^2 \leq \frac{1}{n} \sum_{i=1}^n (y_i - f_T(\mathbf{x}_i))^2 + \frac{\mu}{n} \|f_T\|_{\mathcal{N}}^2 = \frac{\mu}{n} \|f_T\|_{\mathcal{N}}^2, \quad (\text{E.19})$$

which implies  $\|\hat{f}\|_{\mathcal{N}} \leq \|f_T\|_{\mathcal{N}}$ , where we utilize  $y_i = f_T(\mathbf{x}_i)$  in the separable case.

Direct computation shows

$$\begin{aligned} f_{\mathbf{W}(k),a}(\mathbf{x}') &= \hat{f}_1(\mathbf{x}') - (\hat{f}_1(\mathbf{x}') - \hat{f}(\mathbf{x}')) - (\hat{f}(\mathbf{x}') - f_{\mathbf{W}(k),a}(\mathbf{x}')) \\ &= 1 - I_1 - I_2, \end{aligned} \quad (\text{E.20})$$

where we use  $\hat{f}_1(\mathbf{x}_i) = 1$  for any  $\mathbf{x}_i \in \Omega$ ; thus  $\hat{f}_1(\mathbf{x}') = 1$ .

By the representer theorem,  $\hat{f}$  and  $\hat{f}_1$  can be expressed as

$$\hat{f}_1(\mathbf{x}) = h(\mathbf{x}, \mathbf{X})(\mathbf{H}^\infty + \mu \mathbf{I}_n)^{-1} \mathbf{y}, \quad \hat{f}(\mathbf{x}) = h(\mathbf{x}, \mathbf{X})(\mathbf{H}^\infty)^{-1} \mathbf{y},$$

where  $h(\mathbf{x}, \mathbf{X}) = (h(\mathbf{x}, \mathbf{x}_1), \dots, h(\mathbf{x}, \mathbf{x}_n)) \in \mathbb{R}^{1 \times n}$ ,  $\mathbf{H}^\infty = (h(\mathbf{x}_i, \mathbf{x}_j))_{n \times n}$ , and  $\mathbf{y} = (y_1, \dots, y_n)^\top = (f_T(\mathbf{x}_1), \dots, f_T(\mathbf{x}_n))^\top$ . Thus, the first term  $I_1$  in Equation E.20 can be bounded by

$$\begin{aligned} |I_1| &= |\hat{f}_1(\mathbf{x}') - \hat{f}(\mathbf{x}')| = |h(\mathbf{x}', \mathbf{X})(\mathbf{H}^\infty)^{-1} \mathbf{y} - h(\mathbf{x}', \mathbf{X})(\mathbf{H}^\infty + \mu \mathbf{I}_n)^{-1} \mathbf{y}| \\ &= |\mu h(\mathbf{x}', \mathbf{X})(\mathbf{H}^\infty)^{-1} (\mathbf{H}^\infty + \mu \mathbf{I}_n)^{-1} \mathbf{y}| \\ &\leq \mu \sqrt{h(\mathbf{x}', \mathbf{X})(\mathbf{H}^\infty)^{-1} (\mathbf{H}^\infty + \mu \mathbf{I}_n)^{-1} (\mathbf{H}^\infty)^{-1} h(\mathbf{x}', \mathbf{X})^\top \mathbf{y}^\top (\mathbf{H}^\infty + \mu \mathbf{I}_n)^{-1} \mathbf{y}} \\ &= \mu \sqrt{h(\mathbf{x}', \mathbf{X})(\mathbf{H}^\infty)^{-1} (\mathbf{H}^\infty + \mu \mathbf{I}_n)^{-1} (\mathbf{H}^\infty)^{-1} h(\mathbf{x}', \mathbf{X})^\top} \|\hat{f}_1\|_{\mathcal{N}} \\ &\leq \sqrt{\mu} \sqrt{h(\mathbf{x}', \mathbf{X})(\mathbf{H}^\infty)^{-2} h(\mathbf{x}', \mathbf{X})^\top} \|f_T\|_{\mathcal{N}} = \sqrt{\mu} \|f_T\|_{\mathcal{N}}, \end{aligned} \quad (\text{E.21})$$

where the first inequality is by the Cauchy-Schwarz inequality, the second inequality is because  $(\mathbf{H}^\infty + \mu \mathbf{I}_n)^{-1} \preceq \mu^{-1} \mathbf{I}_n$ , and the last equality is because for any  $\mathbf{x}_i$ ,  $(\mathbf{H}^\infty)^{-1} h(\mathbf{x}_i, \mathbf{X})^\top = \mathbf{e}_i$ . Therefore,  $I_1$  converges to zero as  $n \rightarrow \infty$  since  $\mu = o(1)$ . Specifically, there exists an  $n_1$  such that when  $n \geq n_1$ ,  $|I_1| \leq 1/4$ .

The second term  $I_2$  in Equation E.20 can be bounded by

$$\begin{aligned} |I_2| &\leq \|\hat{f} - f_{\mathbf{W}(k),a}\|_\infty \leq C_4 \|\hat{f} - f_{\mathbf{W}(k),a}\|_{\mathcal{N}}^{\frac{d-1}{d}} \|\hat{f} - f_{\mathbf{W}(k),a}\|_2^{\frac{1}{d}} \\ &\leq C_4 (\|\hat{f}\|_{\mathcal{N}} + \|f_{\mathbf{W}(k),a}\|_{\mathcal{N}})^{\frac{d-1}{d}} \|\hat{f} - f_{\mathbf{W}(k),a}\|_2^{\frac{1}{d}} \\ &\leq C_5 \left( \mathbb{E}_{X \sim P_X} (\hat{f}(X) - f_{\mathbf{W}(k),a}(X))^2 \right)^{\frac{1}{2d}}, \end{aligned} \quad (\text{E.22})$$

which converges to zero by Lemma E.3. In Equation E.22, the second inequality is by the interpolation inequality, the third inequality is by the triangle inequality, and the last inequality is because of Assumption D.5. Therefore, there exists an  $n_2$  such that when  $n \geq n_2$ , with probability at least  $1 - \delta$ ,  $|I_2| \leq 1/4$ .

Take  $n_0 = \max(n_1, n_2)$ . For  $n \geq n_0$ , Equation E.20 gives us  $f_{\mathbf{W}(k),a}(\mathbf{x}') \geq 1/2$  with probability at least  $1 - \delta$ . Therefore, by Equation E.18, as long as

$$\|\mathbf{x}' - \mathbf{x}\|_2 = \min_{\mathbf{x}_i} \|\mathbf{x}_i - \mathbf{x}\|_2 \leq (4C_3)^{-1/\tau} := C_6, \forall \mathbf{x} \in \Omega_1, \quad (\text{E.23})$$

we have  $f_{\mathbf{W}(k),a}(\mathbf{x}) \geq 1/4$  for all  $\mathbf{x} \in \Omega_1$ , which implies that the missclassification rate is zero.

Let  $N(\delta, \Omega_1, \|\cdot\|_2)$  be the covering number of  $\Omega_1$  and  $N_0 = N(C_6/2, \Omega_1, \|\cdot\|_2)$ . Since  $\Omega_1$  is compact and  $C_6 > 0$ ,  $N_0$  is finite (and is a constant). Therefore,  $\Omega_1$  can be covered by  $N_0$  balls

with radius  $C_6/2$  (denoted by  $\mathbf{B}$ ), and as long as for each ball, there exists one point  $\mathbf{x}_j$  in this ball, Equation E.23 is satisfied. Since  $\mathbf{x}_k$  has a probability density function with lower bound  $c_1$ , it remains to upper bound the probability that there exists one ball such that there is no point in it. Define this event as  $\mathcal{A}$ . The union bound of probability implies that for  $n > n_0$ ,

$$\mathbb{P}(\mathcal{A}) \leq N_0 \left(1 - \frac{c_1 \text{Vol}(\mathbf{B})}{\text{Vol}(\Omega_1)}\right)^n \leq N_0 \exp(-C_7 n),$$

where  $C_7 = -\log \left(1 - \frac{c_1 \text{Vol}(\mathbf{B})}{\text{Vol}(\Omega_1)}\right)$  is a positive constant. Clearly, we can adjust the constants such that the results in Theorem 3.2 holds for all  $n$ . This finishes the proof.

### E.3 PROOF OF THEOREM 3.3

Note that  $f_{\mathbf{W}(k),\mathbf{a}}$  is a classifier, and the decision boundary is defined by  $\mathcal{D}_T := \{\mathbf{x} | f_{\mathbf{W}(k),\mathbf{a}}(\mathbf{x}) = 0\}$ . Take any point  $\mathbf{x}'$  in  $\mathcal{D}_T$ . The definition of the Hölder space and Equation E.17 imply

$$\frac{|f_{\mathbf{W}(k),\mathbf{a}}(\mathbf{x}) - f_{\mathbf{W}(k),\mathbf{a}}(\mathbf{x}')|}{\|\mathbf{x} - \mathbf{x}'\|_2^\tau} \leq C_2 \|f_{\mathbf{W}(k),\mathbf{a}}\|_{\mathcal{N}} \leq C_3, \forall \mathbf{x} \in \Omega, \quad (\text{E.24})$$

which is the same as

$$\|\mathbf{x} - \mathbf{x}'\|_2^\tau \geq |f_{\mathbf{W}(k),\mathbf{a}}(\mathbf{x})|/C_3, \forall \mathbf{x} \in \Omega, \quad (\text{E.25})$$

where the last inequality in Equation E.24 is because of Assumption D.3. Therefore, it suffices to provide a lower bound of  $|f_{\mathbf{W}(k),\mathbf{a}}(\mathbf{x})|$ . Without loss of generality, let  $\mathbf{x} \in \Omega_1$ , because the case  $\mathbf{x} \in \Omega_2$  can be proved similarly. However, this has already been proved in the proof of Theorem 3.2, where we showed that with probability at least  $1 - \delta - C_4 \exp(-C_5 n)$ ,  $f_{\mathbf{W}(k),\mathbf{a}}(\mathbf{x}) \geq 1/4$  for all  $\mathbf{x} \in \Omega_1$ .

### E.4 PROOF OF THEOREM 3.4

By applying the interpolation inequality, the  $L_\infty$  norm of  $2\eta - 1 - f_{\mathbf{W}(k),\mathbf{a}}$  can be bounded by

$$\begin{aligned} \|2\eta - 1 - f_{\mathbf{W}(k),\mathbf{a}}\|_\infty &\leq C_0 \|2\eta - 1 - f_{\mathbf{W}(k),\mathbf{a}}\|_2^{\frac{1}{d}} \|2\eta - 1 - f_{\mathbf{W}(k),\mathbf{a}}\|_{W^\nu}^{\frac{d-1}{d}} \\ &\leq C_1 \|2\eta - 1 - f_{\mathbf{W}(k),\mathbf{a}}\|_2^{\frac{1}{d}} \|2\eta - 1 - f_{\mathbf{W}(k),\mathbf{a}}\|_{\mathcal{N}}^{\frac{d-1}{d}} \\ &\leq C_2 \|2\eta - 1 - f_{\mathbf{W}(k),\mathbf{a}}\|_2^{\frac{1}{d}} (\|2\eta - 1\|_{\mathcal{N}} + \|f_{\mathbf{W}(k),\mathbf{a}}\|_{\mathcal{N}})^{\frac{d-1}{d}} \\ &\leq C_3 \|2\eta - 1 - f_{\mathbf{W}(k),\mathbf{a}}\|_2^{\frac{1}{d}} \\ &\leq C_4 (\mathbb{E}_{X \sim P_X} (2\eta(X) - 1 - f_{\mathbf{W}(k),\mathbf{a}}(X))^2)^{\frac{1}{2d}} = O_{\mathbb{P}}(n^{-\frac{1}{4d-2}}) \end{aligned} \quad (\text{E.26})$$

where the second equality is by the equivariance of the Sobolev space  $W^\nu$  and the RKHS  $\mathcal{N}$ ; the third inequality is by the triangle inequality; the fourth inequality is by Assumptions D.2 and D.3; the fifth inequality is because of Assumption D.5; and the last equality (with big  $O$  notation) is because of Equation E.15. This finishes the proof.

### E.5 PROOF OF LEMMA 3.5

Let's first consider the simplest  $d = 1$  case where  $\Omega_1 = \{\gamma\}$  and  $\Omega_2 = \{-\gamma\}$ . Let  $\phi$  denote the standard normal  $N(0, 1)$  density. By injecting Gaussian noises  $N(0, v^2)$ , the induced conditional probability can be written as

$$\tilde{\eta}_v(x) = \frac{\phi(\frac{x-\gamma}{v})}{\phi(\frac{x-\gamma}{v}) + \phi(\frac{x+\gamma}{v})} = \frac{1}{1 + \exp(-\frac{2\gamma x}{v^2})}.$$

For small enough  $1/2 > t > 0$ , direct calculation yields  $\{x \in \mathbb{R} : |2\tilde{\eta}_v(x) - 1| < t\} = (-x_t, x_t)$  where

$$x_t = \frac{v^2}{2\gamma} \log \left( \frac{1+t}{1-t} \right) \leq \frac{2v^2}{\gamma} t.$$

Hence,

$$\begin{aligned} P_X(|2\tilde{\eta}_v(x) - 1| < t) &= P_X(-x_t < x < x_t) \leq 2x_t\phi((\gamma + x_t)/v) \\ &\leq \frac{Cv^2}{\gamma} \exp\left(-\frac{\gamma^2}{2v^2}\right) \cdot t. \end{aligned}$$

In general cases, notice that Tsybakov's noise condition measures the separation between classes. Therefore, the bottleneck for the inequality is where  $\Omega_1$  and  $\Omega_2$  are the closest, i.e., where margin  $2\gamma$  is attained. Let  $\mathbf{x}_+ \in \Omega_1$  and  $\mathbf{x}_- \in \Omega_2$  satisfy  $\|\mathbf{x}_+ - \mathbf{x}_-\|_2 = 2\gamma$  (which can be attained since  $\Omega$  is compact). Consider the delta distribution at  $\mathbf{x}_+$  and  $\mathbf{x}_-$ , which is less separated than the original distribution. Then, it reduces to the simplest case.

## E.6 PROOF OF THEOREM 3.6

A closer look at the proof of Theorem 3.1 reveals that the convergence rate depends polynomially on the constant in Tsybakov's noise condition. Specifically, it can be checked that  $\|\tilde{\eta}_v\|_{\mathcal{N}}$  converges to infinity and  $\mu$  converges to zero polynomially with  $v \rightarrow 0$ . Under Tsybakov's noise condition, the convergence rate can be obtained via the proof of Theorem 3.1 as

$$L(\hat{f}) - L^* = O_{\mathbb{P}}\left(C^{\frac{1}{\kappa+2}} n^{-\frac{d(\kappa+1)}{(2d-1)(\kappa+2)}}\right) = O_{\mathbb{P}}\left(\text{poly}\left(\frac{1}{v}\right) C^{\frac{1}{\kappa+2}} n^{-\frac{d(\kappa+1)}{(2d-1)(\kappa+2)}}\right).$$

In the Gaussian noise injection case, if we choose  $v = v_n = n^{-1/2}$ , applying Lemma 3.5 yields

$$L(\hat{f}) - L^* = O_{\mathbb{P}}(e^{-n\gamma/6} \text{poly}(n)) = O_{\mathbb{P}}(e^{-n\gamma/7}).$$

## E.7 PROOF OF THEOREM 4.1

Direct computation implies that

$$f_i^*(\mathbf{x}) = \sum_{j=1}^K \eta_j(\mathbf{x}) v_{ji},$$

which implies

$$f^*(\mathbf{x}) = (\mathbf{v}_1, \dots, \mathbf{v}_K) \eta(\mathbf{x}), \quad (\text{E.27})$$

where  $v_{ji}$  is the  $i$ -th element of  $\mathbf{v}_j$ . Let  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_K)$ . Multiplying  $\mathbf{V}^\top$  on both sides of Equation E.27 leads to

$$\begin{aligned} \mathbf{V}^\top f^*(\mathbf{x}) &= \mathbf{V}^\top (\mathbf{v}_1, \dots, \mathbf{v}_K) \eta(\mathbf{x}) = \left( \frac{K}{K-1} \mathbf{I} - \frac{1}{K-1} \mathbf{1}\mathbf{1}^\top \right) \eta(\mathbf{x}) \\ &= \frac{K}{K-1} \eta(\mathbf{x}) - \frac{1}{K-1} \mathbf{1}\mathbf{1}^\top \eta(\mathbf{x}) \\ &= \frac{K}{K-1} \eta(\mathbf{x}) - \frac{1}{K-1} \mathbf{1}, \end{aligned} \quad (\text{E.28})$$

where  $\mathbf{1} = (1, \dots, 1)^\top$ . In Equation E.28, the second equality is because  $\mathbf{v}_i^\top \mathbf{v}_j = -1/(K-1)$  if  $i \neq j$  and  $\mathbf{v}_i^\top \mathbf{v}_i = 1$ ; and the last equality is because  $\sum_{i=1}^n \eta_i(\mathbf{x}) = 1$ . By Equation E.28, it can be seen that

$$\eta_j(\mathbf{x}) = \frac{(K-1)\mathbf{v}_j^\top f^*(\mathbf{x}) + 1}{K},$$

which finishes the proof.

## F PROOF OF LEMMAS IN THE APPENDIX

### F.1 PROOF OF LEMMA E.2

We follow the approach in the proof of Lemma 6.1 and Proposition 6.3 in Steinwart et al. (2007). Note that  $f_{l,P} = 2p - 1$  minimizes  $\mathcal{R}_{l,P}$ . We first show that for all  $f \in \mathcal{F}$  and all  $\alpha \geq 0$ ,

$$\mathbb{E}_{X,Y \sim P}(l_1 \circ f - l_1 \circ f_{l,P})^2 \leq C_{\eta,\kappa} (\|f\|_\infty + 1)^{\frac{2\kappa+4\alpha}{\kappa+\alpha}} \|(2\eta - 1)^{-1}\|_{q,\infty}^{\frac{\alpha}{\kappa+\alpha}} \mathbb{E}_{X,Y \sim P}(l_1 \circ f - l_1 \circ f_{l,P})^{\frac{\kappa}{\kappa+\alpha}}, \quad (\text{F.1})$$

where  $C_{\eta,\kappa} := \|(2\eta - 1)^{-1}\|_{\kappa,\infty} + 4$ . In particular, one can take  $\alpha = 0$  and obtain

$$\mathbb{E}_{X,Y \sim P}(l_1 \circ f - l_1 \circ f_{l,P})^2 \leq C_{\eta,\kappa}(\|f\|_\infty + 1)^2 \mathbb{E}_{X,Y \sim P}(l_1 \circ f - l_1 \circ f_{l,P}). \quad (\text{F.2})$$

Clearly, Tsybakov's noise condition implies that  $\|(2\eta - 1)^{-1}\|_{\kappa,\infty}$  exists. For  $\mathbf{x} \in \Omega$ , let  $p := \mathbb{P}(Y = 1|\mathbf{x})$  and  $t := f(\mathbf{x})$ . Without loss of generality, let  $p > 1/2$ . Additionally, we denote

$$\begin{aligned} v(p, t) &= p(l(1, t) - l(1, f_{l,P}(\mathbf{x})))^2 + (1-p)(l(-1, t) - l(-1, f_{l,P}(\mathbf{x})))^2, \\ m(p, t) &= p(l(1, t) - l(1, f_{l,P}(\mathbf{x}))) + (1-p)(l(-1, t) - l(-1, f_{l,P}(\mathbf{x}))). \end{aligned} \quad (\text{F.3})$$

Note  $f_{l,P} = 2p - 1$  implies  $l(1, f_{l,P}(\mathbf{x})) = 4(p - 1)^2$  and  $l(-1, f_{l,P}(\mathbf{x})) = 4p^2$ . Plugging them into Equation F.3, it can be checked that

$$\begin{aligned} m(p, t) &= 1 + t^2 + 2(1 - 2p)t - 4p(1 - p) = (1 + t - 2p)^2, \\ v(p, t) &= (1 + t - 2p)^2((t + 1)^2 + 12p - 4pt - 12p^2). \end{aligned}$$

By taking

$$\alpha \geq \frac{\log 4 - \log(12p - 12p^2 - 4pt + 2 - (t - 1)^2)}{\log |2p - 1|}, \quad (\text{F.4})$$

it can be shown that

$$v(p, t) \leq \left(2t^2 + \frac{4}{|2p - 1|^\alpha}\right) m(p, t). \quad (\text{F.5})$$

Since

$$\frac{\log 4 - \log(12p - 12p^2 - 4pt + 2 - (t - 1)^2)}{\log |2p - 1|} \leq \frac{\log 4 - \log(-\frac{2}{3}t^2 + 4)}{\log |2p - 1|} \leq 0,$$

it suffices to take  $\alpha \geq 0$ . We further define

$$\begin{aligned} g(y, \mathbf{x}) &:= l(y, f(\mathbf{x})) - l(y, f_{l,P}(\mathbf{x})), \\ h_1(\mathbf{x}) &:= \eta(\mathbf{x})g(1, \mathbf{x}) + (1 - \eta(\mathbf{x}))g(-1, \mathbf{x}), \\ h_2(\mathbf{x}) &:= \eta(\mathbf{x})g^2(1, \mathbf{x}) + (1 - \eta(\mathbf{x}))g^2(-1, \mathbf{x}). \end{aligned}$$

Therefore, Equation F.5 implies  $h_2(\mathbf{x}) \leq (2\|f\|_\infty^2 + \frac{4}{|2\eta(\mathbf{x}) - 1|^\alpha})h_1(\mathbf{x})$  for all  $\mathbf{x}$  with  $\eta(\mathbf{x}) \neq 1/2$ .

Hence, we obtain

$$\begin{aligned} \mathbb{E}_{X,Y \sim P}g^2 &= \int_{\{\mathbf{x} | |2\eta(\mathbf{x}) - 1|^{-1} < t\}} h_2(\mathbf{x})dP_X + \int_{\{\mathbf{x} | |2\eta(\mathbf{x}) - 1|^{-1} \geq t\}} h_2(\mathbf{x})dP_X \\ &\leq (2\|f\|_\infty^2 + 4t^\alpha) \int_{\{\mathbf{x} | |2\eta(\mathbf{x}) - 1|^{-1} < t\}} h_1(\mathbf{x})dP_X + \int_{\{\mathbf{x} | |2\eta(\mathbf{x}) - 1|^{-1} \geq t\}} (\|f\|_\infty + 1)^4 dP_X \\ &\leq 4(\|f\|_\infty^2 + t^\alpha) \mathbb{E}_{X,Y \sim P}g + (\|f\|_\infty + 1)^4 \|(2\eta - 1)^{-1}\|_{q,\infty} t^{-\kappa} \\ &\leq 4t^\alpha (\|f\|_\infty + 1)^2 \mathbb{E}_{X,Y \sim P}g + (\|f\|_\infty + 1)^4 \|(2\eta - 1)^{-1}\|_{q,\infty} t^{-\kappa} \\ &\leq 4t^\alpha (\|f\|_\infty + 1)^2 \mathbb{E}_{X,Y \sim P}g + (\|f\|_\infty + 1)^4 \|(2\eta - 1)^{-1}\|_{\kappa,\infty} t^{-\kappa} \\ &\leq C_{\eta,\kappa} (\|f\|_\infty + 1)^{\frac{2\kappa+4\alpha}{\kappa+\alpha}} \|(2\eta - 1)^{-1}\|_{\kappa,\infty}^{\frac{\alpha}{\kappa+\alpha}} \mathbb{E}_{X,Y \sim P}g^{\frac{\kappa}{\kappa+\alpha}}, \end{aligned}$$

where the last equality is implied by taking  $t^{\kappa+\alpha} := (\|f\|_\infty + 1)^2 (\mathbb{E}_{X,Y \sim P}g)^{-1}$ . This shows Equation F.1 holds.

Based on Equation F.1, we can show that Lemma E.1 holds. To see this, let  $\widehat{C} := (K\gamma + 1)^{(2\kappa+4\alpha)/(\kappa+\alpha)}$  and fix an  $f \in \gamma B_N$ . The term  $\mathbb{E}_{X,Y \sim P}(L_1 \circ f - L_1 \circ f_n)^2$  can be bounded by

$$\begin{aligned} &\mathbb{E}_{X,Y \sim P}(L_1 \circ f - L_1 \circ f_n)^2 \\ &\leq 2\mu^2 n^{-2} \|f\|^4 + 2\mu^2 n^{-2} \|f_n\|^4 + 2\mathbb{E}_{X,Y \sim P}(l_1 \circ f - l_1 \circ f_n)^2 \\ &\leq 4\mathbb{E}_{X,Y \sim P}(l_1 \circ f - l_1 \circ f_{l,P})^2 + 4\mathbb{E}_{X,Y \sim P}(l_1 \circ f_{l,P} - l_1 \circ f_n)^2 + 2\mu^2 n^{-2} \|f\|^4 + 2\mu^2 n^{-2} \|f_n\|^4 \\ &\leq 8C_{\eta,\kappa} \widehat{C} (\mathbb{E}_{X,Y \sim P}(l_1 \circ f - l_1 \circ f_{l,P}) + \mathbb{E}_{X,Y \sim P}(l \circ f_n - l \circ f_{l,P}))^{\kappa/(\kappa+\alpha)} + 2\mu^2 n^{-2} \|f\|^4 + 2\mu^2 \|f_n\|^4 \\ &\leq C\widehat{C} (\mathbb{E}_{X,Y \sim P}(l \circ f - l \circ f_{l,P}) + \mathbb{E}_{X,Y \sim P}(l \circ f_n - l \circ f_{l,P}) + \mu^2 n^{-2} \|f\|^4 + \mu^2 n^{-2} \|f_n\|^4)^{\kappa/(\kappa+\alpha)} \\ &\leq C\widehat{C} (\mathbb{E}_{X,Y \sim P}(L_1 \circ f - L_1 \circ f_n) + 2\mathbb{E}_{X,Y \sim P}(l \circ f_n - l \circ f_{l,P}) + 2\mu n^{-1} \|f_n\|^2)^{\kappa/(\kappa+\alpha)} \\ &\leq C\widehat{C} (\mathbb{E}_{X,Y \sim P}(L_1 \circ f - L_1 \circ f_n))^{\kappa/(\kappa+\alpha)} + 2C\widehat{C} a^{\kappa/(\kappa+\alpha)}(\mu). \end{aligned} \quad (\text{F.6})$$

In Equation F.6 the first and second inequalities is because of the Cauchy-Schwarz inequality; the third inequality is because of Equation F.1 and  $a^p + b^p < 2(a+b)^p$  for all  $a, b \geq 0, 0 < p \leq 1$ ; the fourth inequality follows from  $a^p + b^p < 2(a+b)^p$  for all  $a, b \geq 0, 0 < p \leq 1$ ; the fifth inequality is because  $n^{-1}\mu\|f\|^2 \leq 1$  and  $n^{-1}\mu\|f_n\|^2 \leq 1$ ; the last inequality follows  $(a+b)^p < a^p + b^p$  for all  $a, b \geq 0, 0 < p \leq 1$ . This finishes the proof of Lemma E.1.

## F.2 PROOF OF LEMMA E.4

Since there is a positive margin between  $\Omega_1$  and  $\Omega_2$ , we can always find two sets  $\tilde{\Omega}_1$  and  $\tilde{\Omega}_2$  with infinitely smooth boundaries such that  $\Omega_1 \subset \tilde{\Omega}_1$ , and  $\Omega_2 \subset \tilde{\Omega}_2$ . Then the result follows from the Sobolev extension theorem.

## G APPENDIX FOR DETAILED EXPERIMENTS

### G.1 SYNTHETIC DATA

During the neural network training, we use the popular RMSProp optimizer with the default settings, and select the tuning parameter  $\mu$  for SL-ONN +  $\ell_2$  and CE-ONN +  $\ell_2$  by a validation set.

**Separable case** In the separable case, we consider a two-dimension distribution  $P = (\rho \sin \theta + 0.04, \rho \cos \theta)$  where  $\rho = (\theta/4\pi)^{4/5} + \epsilon$  with selected  $\theta$  from  $(0, 4\pi]$  and  $\epsilon \sim \text{unif}([-0.03, 0.03])$ . We draw 100 positive and 100 negative training samples from  $-P$  and  $P$ , respectively. We select the tuning parameter  $\mu$  for SL-ONN +  $\ell_2$  and CE-ONN +  $\ell_2$  by minimizing the validation misclassification rate, where the candidate set of  $\mu$  is  $\{0, 0.001, 0.01, 0.1, 1\}$ . For each  $\mu$ , we generate 40 replications to estimate the mean and standard deviation of validation misclassification rate. We observe that SL-ONN +  $\ell_2$  and CE-ONN +  $\ell_2$  have the least mean and least standard deviation for the validation misclassification rate at  $\mu = 0.1$  and  $\mu = 0.01$ , respectively. The errorbar plot for each  $\mu$  is shown in Figure G.3.

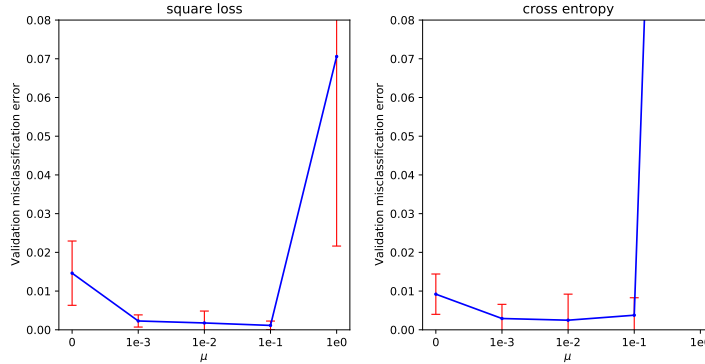


Figure G.3: The errorbar plot of validation misclassification rate with respect to different  $\mu$  in the separable case.

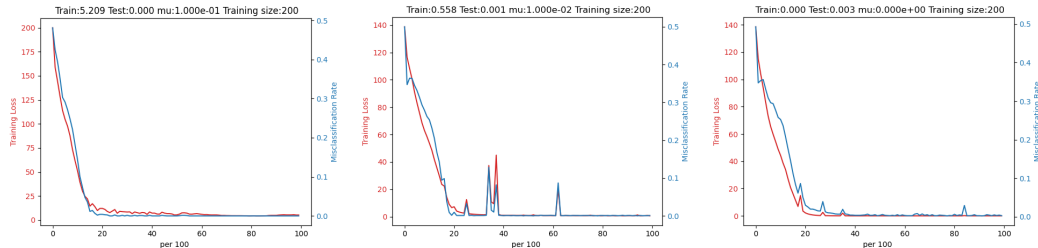


Figure G.4: An instance about the training process of SL-ONN +  $\ell_2$  (Left), CE-ONN +  $\ell_2$  (Center) and CE-ONN (Right).



As mentioned in Section 5, we also consider the cross-entropy loss based ONN without  $\ell_2$  regularization (CE-ONN). All three models are trained for 10000 iterations and achieve 0% training misclassification rate. In Figure G.5, we present five more examples about the decision boundary prediction and test accuracy of SL-ONN +  $\ell_2$ , CE-ONN +  $\ell_2$  and CE-ONN. We can find that SL-ONN +  $\ell_2$  still beats CE-ONN +  $\ell_2$  and CE-ONN in almost all the cases. SL-ONN +  $\ell_2$  attains the smallest misclassification rate and depicts the largest margin decision boundary which separates the positive and negative samples best. In addition, we can observe that CE-ONN +  $\ell_2$  outperforms CE-ONN in all cases, although the  $\ell_2$  regularization term bring some oscillation to the training of CE-ONN +  $\ell_2$ , as shown in Figure G.4.

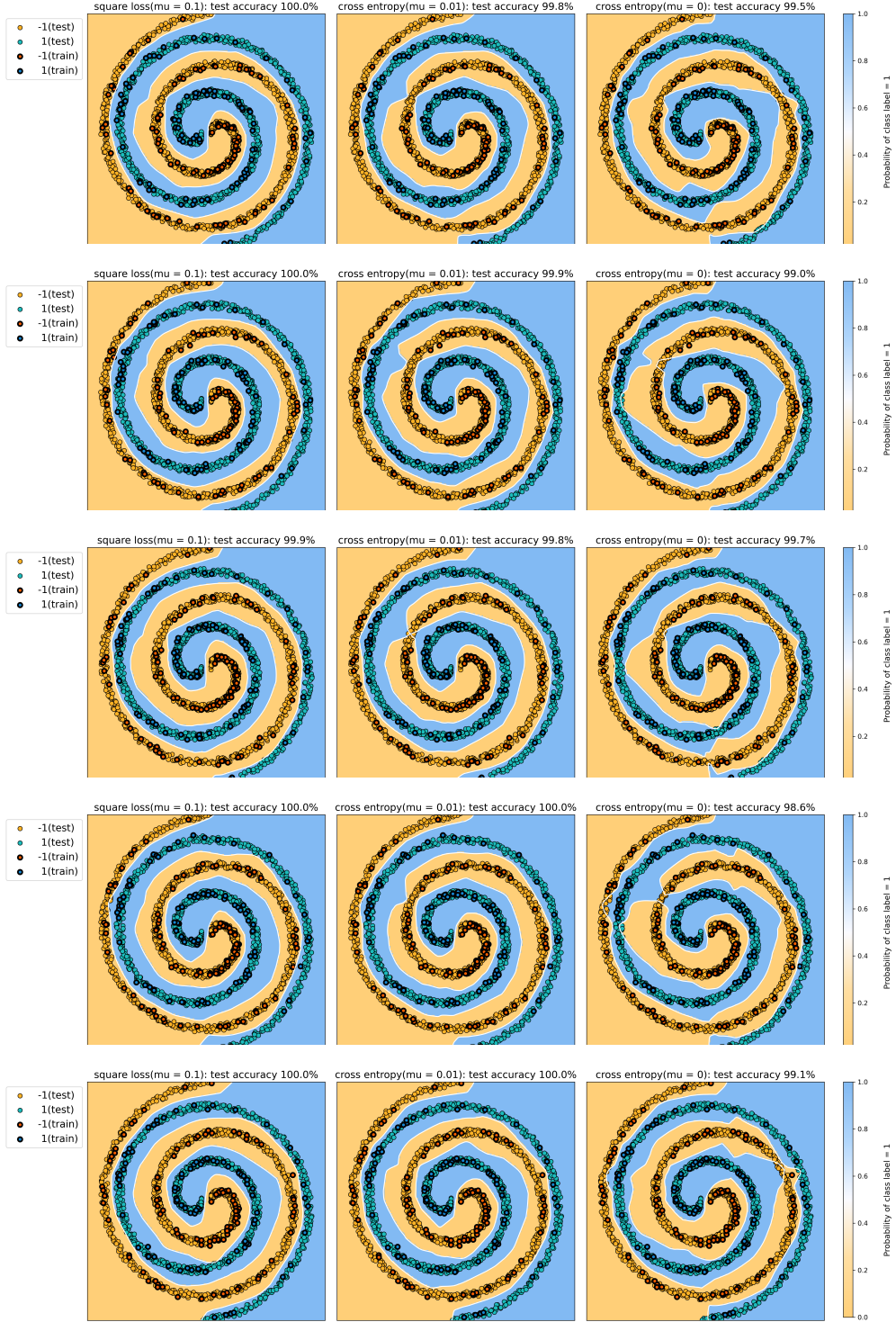


Figure G.5: Five examples of the separable case.

**Non-separable case** In the non-separable case, the training data points  $x_1, \dots, x_n$  are i.i.d. sampled from  $\text{unif}([-1, 1]^2)$  and the training labels  $y_1, \dots, y_n$  are sampled according to

$Bernoulli(\eta(\mathbf{x}_i))$ , where  $\eta(\mathbf{x}) = \sin(\sqrt{2\pi}\|\mathbf{x}\|_2)$ , and  $n = 8000$ . The 3-dimensional plot of  $\eta(\mathbf{x})$  is presented in Figure G.6.

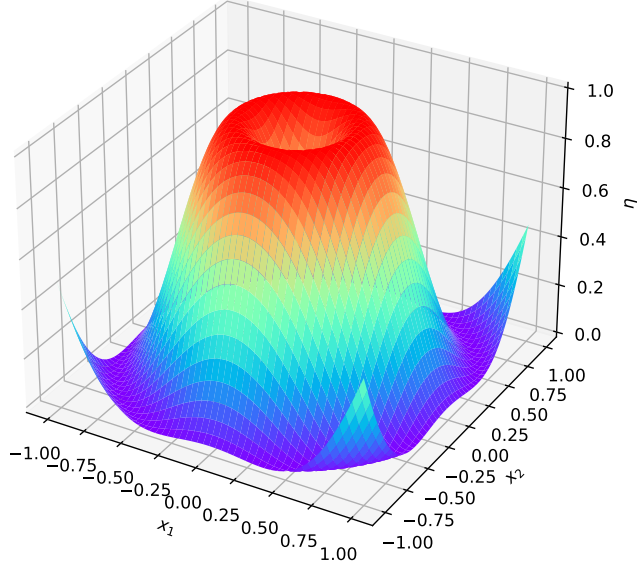


Figure G.6: The 3-dimensional plot of  $\eta(\mathbf{x})$  in the non-separable case.

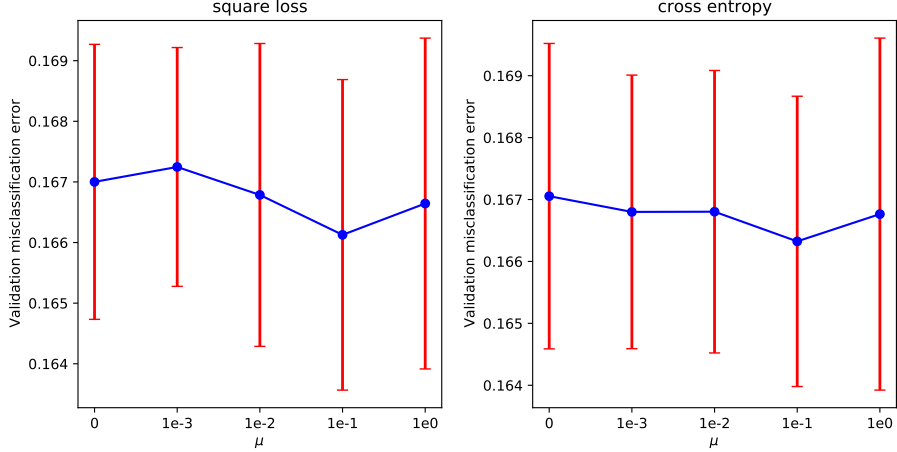


Figure G.7: The errorbar plot of validation misclassification rate with respect to different  $\mu$  in the non-separable case.

We select the tuning parameter  $\mu$  for SL-ONN +  $\ell_2$  and CE-ONN +  $\ell_2$  via a validation set, and the candidate set of  $\mu$  is  $\{0, 0.001, 0.01, 0.1, 1\}$ . For each  $\mu$ , we run 40 replications to estimate the mean and standard deviation of validation misclassification rate. The iteration number of training is 2000. We find SL-ONN +  $\ell_2$  and CE-ONN +  $\ell_2$  have the smallest mean and standard deviation for the validation misclassification rate at  $\mu = 0.1$  and  $\mu = 0.1$ , respectively. The error bar plot <sup>4</sup> for  $\mu$  equaling to 0, 0.001, 0.01, 0.1 and 1 is shown in Figure G.7.

The calibration error results are presented in Figure G.8. The error bar plot of the test calibration error shows that  $\hat{f}_{l2}$  has the smaller mean and standard deviation than  $\hat{f}_{ce}$ .

<sup>4</sup>In an error bar plot, the center of each plot is the mean, and the upper and lower red dashes denote (mean+one standard deviation) and (mean − one standard deviation), respectively.

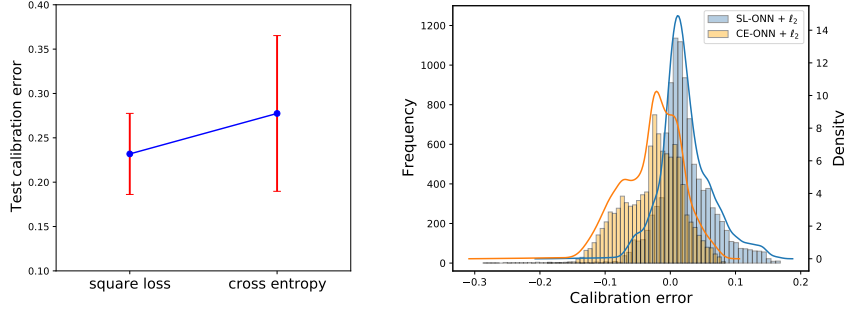


Figure G.8: (Left) The error bar plot of test calibration errors for the 40 replicated runs. (Right) The histogram and kernel density estimation of test calibration errors in one instance from 40 replications. In this instance,  $\|(\hat{f}_{l_2} + 1)/2 - \eta\|_{L_\infty} = 0.188$ ,  $\|(\hat{f}_{ce} + 1)/2 - \eta\|_{L_\infty} = 0.287$ , and the test misclassification rate for  $\hat{f}_{l_2}$  and  $\hat{f}_{ce}$  are 0.167 and 0.164, respectively.

## G.2 REAL DATA

**Data and network architecture** We use the popular CIFAR-10 and CIFAR-100 datasets, with training and testing split of 50000 and 10000. The data loader is from `torch.utils.data`. As typically employed in practice, our training includes data augmentations, a composition of random crop and horizontal flip. We trained two types of neural networks, ResNet (He et al., 2016) and Wide ResNet (Zagoruyko & Komodakis, 2016). To be more specific, we used the default ResNet-18 and ResNet-50 for CIFAR-10 and CIFAR-100 respectively, and the default WRN-16-10 for both CIFAR-10 and CIFAR-100. All experiments are run in PyTorch version 1.9.0 and cuda 10.2.

**Training details** The training algorithm is the default SGD with momentum (0.9) and weight decay (0.0005). The learning rate scheduler is the `StepLR()` from `torch.optim.lr_scheduler` with step size 50. In our experiment, the only parameters that we tuned are the learning rate (lr) and batch size (bs), with only two options, (lr=0.01, bs=32) and (lr=0.1, bs=128). We find that (lr=0.01, bs=32) performs better for most cases except for square loss trained WRN-16-10 on CIFAR-100, where the average accuracy for (lr=0.01, bs=32) is 77.96%, around 1.5% less than that for (lr=0.1, bs=128). Meanwhile, for cross-entropy trained WRN-16-10, (lr=0.1, bs=128) yields an average accuracy of 76.83%, around 1% less than that for (lr=0.01, bs=32). The two training settings perform quite comparable for WRN-16-10 on CIFAR-10. For consistency, we stick with (lr=0.01, bs=32) in this case.

**Adversarial robustness** For square loss, training deep classifiers is the same as regression. When attacking classifiers trained with square loss, the default way of constructing adversarial examples doesn’t work well. To be more specific, for a correctly classified training image  $(\mathbf{x}, y)$ , the adversarial examples are typically generated by

$$\max_{\|\delta\|_\infty=\alpha} L(f(\mathbf{x} + \delta), y).$$

Such an attacking scheme works fine for cross-entropy, where

$$L(f(\mathbf{x}), y) = -\log(\text{softmax}(f(\mathbf{x}))) = -\log\left(\frac{\exp(f_y(\mathbf{x}))}{\sum_{k \neq y} \exp(f_k(\mathbf{x}))}\right),$$

but is problematic for regression losses such as square loss. The fundamental reason lies in Proposition 4.1 and its proof. Recall that the conditional probability for square loss consists of projections of the classifier outputs to all the simplex vertices, some of which are sure to be non-positive. The sum of the class probabilities from Equation 4.1 is always 1 but unlike that from softmax function, the summand can be negative. By maximizing the square loss, the resulting “adversarial” image can stay the same class but more confidently. To illustrate, if  $f(\mathbf{x}) = \mathbf{v}_y$ , the predicted confidence for label  $y$  will be 100% and 0 for other classes. The “adversarial” image may be such that  $f(\mathbf{x} + \delta) = 2\mathbf{v}_y$ , where the predicted label remains unchanged but with an updated confidence of  $2 - 1/K$  for label  $y$  and  $(1/K - 1)/(K - 1) < 0$  for all other classes. This is obviously not a successful attack.

To this end, we devise a special attacking scheme for classifier trained with square loss and simplex coding. The key idea is to choose attack directions tangent to the sphere inscribed by the simplex. Instead of

$$L(f(\mathbf{x}), y) = \|f(\mathbf{x}) - \mathbf{v}_y\|_2^2,$$

we choose

$$L(f(\mathbf{x}), y) = \theta(f(\mathbf{x}), \mathbf{v}_y),$$

where  $\theta(\mathbf{v}_1, \mathbf{v}_2)$  denotes the cosine similarity between  $\mathbf{v}_1$  and  $\mathbf{v}_2$ . We refer to this attack as angle attack.

Empirically, we found our angle attack to significantly outperform the naive attack by maximizing the square loss. For square loss, let the predicted probabilities from Equation 4.1 be  $\hat{\mathbf{p}}$ . Similar to cross entropy, we have also tried two cases of  $L(f(\mathbf{x}), y)$ , which corresponds to

$$L_1(f(\mathbf{x}), y) = -\log(\text{softmax}(\hat{\mathbf{p}}_y(\mathbf{x}))) \quad \text{and} \quad L_2(f(\mathbf{x}), y) = -\log(\hat{\mathbf{p}}_y(\mathbf{x})).$$

Interestingly for PGD-100,  $L_1$  performs the best, beating angle attack for the majority cases, except for attacking WRN-16-10 on CIFAR-100 with strength 2/255. The reported adversarial accuracy for square loss trained classifiers in Table 1 is by  $L_1(f(\mathbf{x}), y) = -\log(\text{softmax}(\hat{\mathbf{p}}_y(\mathbf{x})))$ .

The PGD attack results may be further improved for square loss. Nonetheless, the AutoAttack still provides convincing results, as it includes both white-box and black-box attacks. We used the standard version which includes 4 types of attacks, APGD-CE, APGD-DLR, FAB and Square Attack as in Croce & Hein (2020).

**Robustness to Gaussian Noise** To make the robustness evaluation more comprehensive, beyond the adversarial robustness, we also investigate the classifier’s robustness to Gaussian noise injections. With the image pixels’ value normalized to 0 and 1, we consider injecting Gaussian noises to test images and report the test accuracy. The noise standard deviation ranges from 0.1 to 0.4. The test accuracy results for both CIFAR-10 and CIFAR-100 are listed in Table G.3.

Table G.3: Black-box Gaussian noise robustness results. The reported accuracy is the average of 5 replications.

| Dataset   | Network   | Loss | Gaussian noise standard deviation |              |              |              |              |
|-----------|-----------|------|-----------------------------------|--------------|--------------|--------------|--------------|
|           |           |      | 0.00                              | 0.10         | 0.20         | 0.30         | 0.40         |
| CIFAR-10  | ResNet-18 | SL   | 95.04                             | <b>90.07</b> | <b>70.16</b> | <b>42.13</b> | <b>25.38</b> |
|           |           | CE   | <b>95.15</b>                      | <b>90.03</b> | 69.71        | 41.08        | 24.66        |
|           | WRN-16-10 | SL   | <b>95.02</b>                      | <b>88.49</b> | <b>60.91</b> | <b>35.78</b> | <b>24.04</b> |
|           |           | CE   | 93.94                             | 84.78        | 56.63        | 33.70        | 22.41        |
| CIFAR-100 | ResNet-50 | SL   | <b>78.91</b>                      | <b>63.06</b> | <b>36.64</b> | <b>17.78</b> | <b>9.47</b>  |
|           |           | CE   | 79.82                             | 62.72        | 34.42        | 16.69        | 9.11         |
|           | WRN-16-10 | SL   | <b>79.65</b>                      | <b>62.01</b> | <b>30.69</b> | <b>15.11</b> | <b>8.88</b>  |
|           |           | CE   | 77.89                             | 60.14        | 26.47        | 10.26        | 5.57         |

**Simplex coding vs. one-hot coding** The one-hot coding is the usual choice for applying square loss to classification. However, it is empirically observed to struggle when the number of classes are large. For a single training data point  $\mathbf{x}$  and label  $k$ , Hui & Belkin (2020) proposed to modify the training objective from the typical  $(f_k(\mathbf{x}) - 1)^2 + \sum_{i \neq k} f_i(\mathbf{x})^2$  to  $J \cdot (f_k(\mathbf{x}) - M)^2 + \sum_{i \neq k} f_i(\mathbf{x})^2$ , where  $J, M$  are hyperparameters to make  $f_k$  more prominent in the loss. Similar modification is also proposed in Demirkaya et al. (2020). The scaling trick involves two hyperparameters, which can be hard to tune. We evaluate the two coding schemes in our experiment setting and the results are summarized in Table G.4. The test accuracy for scaled one-hot coding scheme performs comparably for ResNet-18 on CIFAR-10 and ResNet-50 on CIFAR-100. For WRN-16-10, the simplex coding performs better.

Table G.4: Test accuracy for square loss with one-hot coding (scaled) (OC) vs. simplex coding (SC). Accuracy with an asteroid sign (\*) denotes cases where the training accuracy doesn't overfit after 200 training epochs.

| Dataset   | Network   | One-hot scaling | SGD parameters | OC clean acc(%) | SC clean acc(%) |
|-----------|-----------|-----------------|----------------|-----------------|-----------------|
| CIFAR-10  | ResNet-18 | k=1, M=1        | lr=0.01, bs=32 | 94.95           | 95.04           |
|           |           |                 | lr=0.1, bs=128 | 10*             | 10*             |
|           | WRN-16-10 | k=1, M=1        | lr=0.01, bs=32 | 89.75*          | 95.02           |
|           |           |                 | lr=0.1, bs=128 | 88.43*          | 95.03           |
| CIFAR-100 | ResNet-50 | k=5, M=15       | lr=0.01, bs=32 | 79.06           | 78.91           |
|           |           |                 | lr=0.1, bs=128 | 1*              | 1*              |
|           | WRN-16-10 | k=5, M=15       | lr=0.01, bs=32 | 78.42           | 78.06           |
|           |           |                 | lr=0.1, bs=128 | 78.39           | 79.65           |