# Binary Classification under Label Differential Privacy Using Randomized Response Mechanisms

Anonymous authors
Paper under double-blind review

## **Abstract**

Label differential privacy is a popular branch of  $\epsilon$ -differential privacy for protecting labels in training datasets with non-private features. In this paper, we study the generalization performance of a binary classifier trained on a dataset privatized under the label differential privacy achieved by the randomized response mechanism. Particularly, we establish minimax lower bounds for the excess risks of the deep neural network plug-in classifier, theoretically quantifying how privacy guarantee  $\epsilon$  affects its generalization performance. Our theoretical result shows: (1) the randomized response mechanism slows down the convergence of excess risk by lessening the multiplicative constant term compared with the non-private case ( $\epsilon = \infty$ ); (2) as  $\epsilon$  decreases, the optimal structure of the neural network should be smaller for better generalization performance; (3) the convergence of its excess risk is guaranteed even if  $\epsilon$  is adaptive to the size of training sample n at a rate slower than  $O(n^{-1/2})$ . Our theoretical results are validated by extensive simulated examples and two real applications.

# 1 Introduction

In the past decade, differential privacy (DP; Dwork, 2008) has emerged as a standard statistical framework to protect sensitive data before releasing it to an external party. The rationale behind differential privacy is to ensure that information obtained by an external party is robust enough to the presence or absence of a single record in a dataset. Generally, the privacy protection in differential privacy is achieved by injecting noise implicitly to raw data, which inevitably distorts the raw data and hence reduces data utility for downstream learning tasks (Alvim et al., 2012; Kairouz et al., 2016). To achieve better privacy-utility tradeoffs, various research efforts have been devoted to analyzing the effect of differential privacy on learning algorithms in the machine learning community (Ghazi et al., 2021; Bassily et al., 2022; Esfandiari et al., 2022). Depending on whether the data receiver is trusted or not, differential privacy can be categorized into two main classes in the literature, including central differential privacy (CDP; Erlingsson et al. 2019; Girgis et al. 2021) and local differential privacy (LDP; Wang et al. 2017; Arachchige et al. 2019). CDP relies on a trusted curator to protect all data simultaneously, whereas LDP perturbs data on the users' side. LDP becomes a more popular solution to privacy protection due to its successful applications, including Google Chrome browser (Erlingsson et al., 2014) and macOS (Tang et al., 2017).

An important variant of differential privacy is label differential privacy (Label DP; Chaudhuri & Hsu, 2011), which is a relaxation of differential privacy for some real-life scenarios where input features are assumed to be publicly available and labels are highly sensitive and should be protected. Label DP has been gaining increasing attention in recent years due to the emerging demands in some real applications. For example, in recommender systems, users' ratings are sensitive for revealing users' preferences that can be utilized for advertising purposes (McSherry & Mironov, 2009; Xin & Jaakkola, 2014). In online advertising, a user click behavior is usually treated as a sensitive label whereas the product description for the displayed advertisement is publicly available (McMahan et al., 2013; Chapelle et al., 2014). These real scenarios motivate various research efforts to develop mechanisms for achieving label differential privacy and understanding the fundamental tradeoffs between data utility and privacy protection. In literature, label DP can be divided into two main classes depending on whether labels are protected in a local or central manner. In central label DP, privacy protection is guaranteed by ensuring the output of a randomized learning algorithm is

robust to the presence or absence of a single label in the dataset (Chaudhuri & Hsu, 2011; Ghazi et al., 2021; Bassily et al., 2022; Ghazi et al., 2022). In local label DP, labels are altered at the users' side before they are released to learning algorithms, ensuring that it is difficult to infer the true labels based on the released labels (Busa-Fekete et al., 2021; Cunningham et al., 2022).

A critical challenge in the domain of differential privacy is to understand the essential privacy-utility tradeoff that sheds light on the fundamental utility limit for a specific problem. For example, Duchi et al. (2018) established minimax bounds of estimation risks for several estimation problems under local differential privacy constraints, such as mean estimation and median estimation. These minimax bounds in estimation risks quantify the utility of privatized data any estimation procedure can utilize at most for estimation. Built upon this idea, we intend to study the generalization performance of binary classifiers under local label DP, aiming to theoretically understand how the generalization performance of differentially private classifiers are affected by the local label DP under the margin assumption (Tsybakov, 2004). To this end, we not only provide the convergence rate of the excess risk but also establish its minimax rate to show the optimal performance that can be achieved under local label DP constraint.

There are two lines of research that are intimately related to but apparently distinguished from the problem we address. The first line of research studies label differential privacy with the aim to develop more complicated mechanisms to achieve label DP efficiently and analyze their essential privacy-utility tradeoffs in downstream learning tasks. The original way to achieve label DP is via randomized response mechanisms (Warner, 1965), which alters observed labels in a probabilistic manner. For binary labels, the randomized response mechanism flips labels onto the other side with a pre-determined probability (Nayak & Adeshiyan, 2009; Wang et al., 2016b; Busa-Fekete et al., 2021). Ghazi et al. (2021) proposed a multi-stage training algorithm called randomized response with prior, which flips training labels via a prior distribution learned by the trained model in the previous stage. Such a multi-stage training algorithm significantly improves the generalization performance of the trained model under the same privacy guarantee. Malek Esmaeili et al. (2021) proposed to apply Laplace noise addition to one-hot encodings of labels and utilized the iterative Bayesian inference to de-noise the outputs of the privacy-preserving mechanism. Bassily et al. (2022) developed a private learning algorithm under the central label DP and established a dimension-independent deviation margin bound for the generalization performance of several differentially private classifiers, showing that the margin guarantees that are independent of the input dimension. However, their developed learning algorithm relies on the partition of the hypothesis and is hence computationally inefficient.

The second related line of research focuses on the fundamental impacts of local privacy on subsequent inference tasks by establishing minimax risks of estimation. Duchi et al. (2018) derived minimax bounds for several canonical families of problems under local differential privacy constraints, including mean estimation, median estimation, generalized linear models, and non-parametric density estimation, where classical classification problem is not included. Duchi & Ruan (2018) adopt a local minimax risk to measure the optimal risk of estimation since the worst-case nature of minimax lower bounds are sometimes too pessimistic in practice. Wang & Xu (2019) studied the sparse linear regression problem under local label DP and establish the minimax risk for the estimation error under label DP, which theoretically quantifies how label DP affects the asymptotic behavior of the differentially private linear model.

In literature, few attempts are made to theoretically quantify the generalization performance of classifiers under local label DP even though binary classification has already become an indispensable part of the machine learning community. An important characteristic distinguishing binary classification problem is that the convergence of the generalization performance depends on the behavior of data in the vicinity of the decision boundary, which is known as the margin assumption (Tsybakov, 2004). Therefore, our first contribution is that we theoretically quantify how local label DP alters the margin assumption, which allows us to bridge the connection between the local label DP and the generalization performance. Additionally, we mainly consider two scenarios for the function class of classifiers in this paper. First, we consider the large margin classifier with its hypothesis space being a parametric function class and the deep neural network plug-in classifier. For these two scenarios, we establish their upper bound and the minimax lower bound for their excess risks, which theoretically quantifies how  $\epsilon$  affects the generalization performance. The implications of our theoretical results are three-fold. First, the Bayes classifier stays invariant to the randomized response mechanism with any small  $\epsilon$ , which permits the possibility of learning the optimal

classifier from the privatized dataset. Second, the local label DP achieved via the randomized response mechanism implicitly reduces the information for estimation. Specifically, we theoretically prove that the convergence rate of excess risk is slowed down with an additional multiplicative constant depending on  $\epsilon$ . Third, based on our theoretical results, we show that the excess risk fails to converge when the  $\epsilon$  is adaptive to the training sample size n at the order  $O(n^{-1/2})$ , which is independent of the margin assumption. To the best of our knowledge, no existing literature related to classification under DP or corrupted labels (Cannings et al., 2020; van Rooyen & Williamson, 2018) has investigated the effects of noise on neural network structures. Our theoretical results are supported by extensive simulations and two real applications.

The rest of this paper is organized as follows. After introducing some necessary notations in Section 1.1, Section 2 introduces the backgrounds of the binary classification problem, neural network. and the local label differential privacy. In Section 3, we introduce the framework of the differentially private learning under the label differential privacy and present theoretical results regarding the asymptotic behavior of the differentially private classifier. Section 4 quantifies how  $\epsilon$  affects the generalization performance of the deep neural network plug-in classifier by establishing a minimax lower bound. Section 5 and Section 6 conduct a series of simulations and real applications to support our theoretical results. All technical proofs are provided in the Appendix.

#### 1.1 Notation

For a vector  $\mathbf{x} \in \mathbb{R}^p$ , we denote its  $l_1$ -norm and  $l_2$ -norm as  $\|\mathbf{x}\|_1 = \sum_{i=1}^p |x_i|$  and  $\|\mathbf{x}\|_2 = \left(\sum_{i=1}^p |x_i|^2\right)^{1/2}$ , respectively. For a function  $f: \mathcal{X} \to \mathbb{R}$ , we denote its  $L_p$ -norm with respect to the probability measure  $\mu$  as  $\|f\|_{L^p(\mu)} = \left(\int_{\mathcal{X}} |f(\mathbf{x})|^p d\mu(\mathbf{x})\right)^{1/p}$ . For a real number a, we let  $\lfloor a \rfloor$  denote the largest integer not larger than a. For a set S, we define  $\mathcal{N}(\xi, S, \|\cdot\|)$  as the minimal number of  $\xi$ -balls needed to cover S under a generic metric  $\|\cdot\|$ . For two given sequences  $\{A_n\}_{n\in\mathbb{N}}$  and  $\{B_n\}_{n\in\mathbb{N}}$ , we write  $A_n \gtrsim B_n$  if there exists a constant C > 0 such that  $A_n \geq CB_n$  for any  $n \in \mathbb{N}$ . Additionally, we write  $A_n \approx B_n$  if  $A_n \gtrsim B_n$  and  $A_n \lesssim B_n$ .

# 2 Preliminaries

## 2.1 Binary Classification

The goal in binary classification is to learn a discriminant function f, which well characterizes the functional relationship between the feature vector  $X \in \mathcal{X}$  and its associated label  $Y \in \{-1, 1\}$ . To measure the quality of f, the 0-1 risk is usually employed,

$$R(f) = \mathbb{E}(I(f(X)Y < 0)) = \mathbb{P}(\operatorname{sign}(f(X)) \neq Y),$$

where  $I(\cdot)$  denotes the indicator function and the expectation is taken with respect to the joint distribution of (X,Y).

Let  $f^* = \inf_f R(f)$  denote the minimizer of R(f), which refers to as the Bayes decision rule. Generally,  $f^*$  is obtained by minimizing R(f) in a point-wise manner and given as  $f^*(X) = \operatorname{sign} (\eta(X) - 1/2)$  with  $\eta(X) = \mathbb{P}(Y = 1|X)$ . The minimal risk  $R(f^*)$  can be written as  $R(f^*) = \mathbb{E}_X (\min\{\eta(X), 1 - \eta(X)\})$ . In practice, the underlying joint distribution on (X, Y) is unavailable, but a set of i.i.d. realizations  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  is given. Therefore, it is a common practice to consider the estimation procedure based on minimizing the sample average of a surrogate loss, which is given as

$$\widehat{R}_{\phi}(f) = \frac{1}{n} \sum_{i=1}^{n} \phi(f(\boldsymbol{x}_i)y_i),$$

where  $\phi(\cdot)$  is the surrogate loss function replacing the 0-1 loss since the 0-1 loss is computationally intractable (Arora et al., 1997).

Let  $R_{\phi}(f) = \mathbb{E}(\phi(f(\boldsymbol{X})Y))$  denote the  $\phi$ -risk. Given that  $\phi$  is classification calibrated, the minimizer  $f_{\phi}^* = \arg\min_f R_{\phi}(f)$  is consistent with the Bayes decision rule (Lin, 2004), i.e.,  $\operatorname{sign}(f_{\phi}^*(\boldsymbol{x})) = \operatorname{sign}(\eta(\boldsymbol{x}) - 1/2)$  for any  $\boldsymbol{x} \in \mathcal{X}$ . In literature, there are various classification-calibrated loss functions (Zhang, 2004; Bartlett et al., 2006), such as exponential loss, hinge loss, logistic loss, and  $\psi$ -loss (Shen et al., 2003).

#### 2.2 Deep Neural Network and Function Class

Let  $f(x;\Theta)$  be an L-layer neural network with Rectified Linear Unit (ReLU) activation function, that is,

$$f(\boldsymbol{x};\Theta) = \boldsymbol{A}_{L+1} (\boldsymbol{h}_L \circ \boldsymbol{h}_{L-1} \circ \cdots \boldsymbol{h}_1(\boldsymbol{x})) + \boldsymbol{b}_{L+1},$$

where  $\circ$  denotes function composition,  $h_l(x) = \sigma(A_l x + b_l)$  denotes the l-th layer, and  $\Theta = \{(A_l, b_l)\}_{l=1,\dots,L+1}$  denotes all the parameters. Here  $A_l \in \mathbb{R}^{p_l \times p_{l-1}}$  is the weight matrix,  $b_l \in \mathbb{R}^{p_l}$  is the bias term,  $p_l$  is the number of neurons in the l-th layer, and  $\sigma(x) = \max\{0, x\}$  is the ReLU function. To characterize the network architecture of f, we denote the number of layers in  $\Theta$  as  $\Upsilon(\Theta)$ , the maximum number of nodes as  $\Delta(\Theta)$ , the number of non-zero parameters as  $\|\Theta\|_0$ , the largest absolute value in  $\Theta$  as  $\|\Theta\|_{\infty}$ . For a given n, we denote by  $\mathcal{F}_n^{NN}(L_n, N_n, P_n, B_n, V_n)$  a class of neural networks, which is defined as

$$\mathcal{F}_{n}^{NN}(L_{n}, N_{n}, P_{n}, B_{n}, V_{n}) = \{ f(\boldsymbol{x}; \Theta) : \Upsilon(\Theta) \leq L_{n}, \Delta(\Theta) \leq N_{n}, \\ \|\Theta\|_{0} \leq P_{n}, \|\Theta\|_{\infty} \leq B_{n}, \|f\|_{\infty} \leq V_{n}, \}.$$

Let  $\beta > 0$  be a degree of smoothness, then the Hölder space is defined as

$$\mathcal{H}(\beta, \mathcal{X}) = \{ f \in \mathcal{C}^{\lfloor \beta \rfloor}(\mathcal{X}) : ||f||_{\mathcal{H}(\beta, \mathcal{X})} < \infty \},$$

where  $\mathcal{C}^{\lfloor \beta \rfloor}(\mathcal{X})$  the class of  $\lfloor \beta \rfloor$  times continuously differentiable functions on the open set  $\mathcal{X}$  and the Hölder norm  $\|f\|_{\mathcal{H}(\beta,\mathcal{X})}$  is given as

$$||f||_{\mathcal{H}(\beta,\mathcal{X})} = \max_{\boldsymbol{m}: ||\boldsymbol{m}||_1 \le \lfloor \beta \rfloor} \sup_{\boldsymbol{x} \in \mathcal{X}} |\partial^{\boldsymbol{m}} f(\boldsymbol{x})| + \max_{\boldsymbol{m}: ||\boldsymbol{m}||_1 = \lfloor \beta \rfloor} \sup_{\boldsymbol{x}, \boldsymbol{y} \in \mathcal{X}, \boldsymbol{x} \neq \boldsymbol{y}} \frac{|\partial^{\boldsymbol{m}} f(\boldsymbol{x}) - \partial^{\boldsymbol{m}} f(\boldsymbol{y})|}{||\boldsymbol{x} - \boldsymbol{y}||_2^{\beta - \lfloor \beta \rfloor}},$$

where  $\partial^{\boldsymbol{m}} f(\boldsymbol{x}) = \frac{\partial^{m_1 + \dots + m_p}}{\partial x_1^{m_1} \dots \partial x_p^{m_p}} f(\boldsymbol{x})$  denotes the partial derivative of order  $\boldsymbol{m}$  with respect to  $\boldsymbol{x}$  and  $\boldsymbol{m} = (m_1, \dots, m_p) \in \mathbb{N}_0^p$  is a multi-index with  $\mathbb{N}_0 = \mathbb{N} \cup \{0\}$ . Further, we let  $\mathcal{H}(\beta, \mathcal{X}, M) = \{f \in \mathcal{H}(\beta, \mathcal{X}) : \|f\|_{\mathcal{H}(\beta, \mathcal{X})} \leq M\}$  be a closed ball of radius M in  $\mathcal{H}(\beta, \mathcal{X})$ .

#### 2.3 Local Label Differential Privacy

Label differential privacy (Label DP; Chaudhuri & Hsu, 2011) is proposed as a relaxation of differential privacy (Dwork, 2008), aiming to protect the privacy of labels in the dataset, whereas training features are non-sensitive and hence publicly available. An effective approach to Label DP is local label differential privacy (LLDP; Busa-Fekete et al., 2021). As its name suggests, it protects the privacy of labels via some local randomized response mechanisms under the framework of local differential privacy.

**Definition 1.** (LLDP; Ghazi et al., 2021) Let  $\epsilon > 0$ . A randomized mechanism  $\mathcal{A} : \{-1,1\} \to \{-1,1\}$  is label-locally differential private if  $\max_{y \in \{-1,1\}} \log \left| \frac{\mathbb{P}(\mathcal{A}(Y) = y|Y = 1)}{\mathbb{P}(\mathcal{A}(Y) = y|Y = 1)} \right| \leq \epsilon$ , where the probability is taken with respect to the randomness of  $\mathcal{A}$ .

A direct way to achieve  $\epsilon$ -LLDP in binary classification is to employ the randomized response mechanism (Warner, 1965; Nayak & Adeshiyan, 2009; Wang et al., 2016b; Karwa et al., 2017). The main idea of the binary randomized response mechanism is to flip observed labels with a fixed probability. Specifically, let  $\mathcal{A}_{\theta}$  denote the randomized response mechanism parametrized by  $\theta$ . For an input label Y, we define

$$\mathcal{A}_{\theta}(Y) = \begin{cases} Y, & \text{with probability } \theta, \\ -Y, & \text{with probability } 1 - \theta, \end{cases}$$

where  $\theta > 1/2$  denotes the probability that the value of Y stays unchanged. It is straightforward to verify that the randomized response mechanism satisfies  $\epsilon$ -LLDP with  $\epsilon = \log(\theta/(1-\theta))$  (Ghazi et al., 2021; Busa-Fekete et al., 2021).

# 3 Differentially Private Learning

## 3.1 Effect of Locally-Label Differential Privacy

Under the setting of local differential privacy, users do not trust the data curator and hence privatize their sensitive data via some randomized mechanisms locally before releasing them to central servers. We let  $\mathcal{D} = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$  denote the original dataset containing n i.i.d. realizations of the random pair  $(\boldsymbol{X}, Y)$ . As illustrated in Figure 1, users' labels are privatized by a locally differentially private protocol  $\mathcal{A}_{\theta}$ , and then the untrusted curator receives a privatized dataset, which we denote as  $\widetilde{\mathcal{D}} = \{(\boldsymbol{x}_i, \widetilde{y}_i)\}_{i=1}^n$ .

A natural question is what the quantitative relation between  $\epsilon$ -LLDP and the discrepancy between the distributions of  $\mathcal{D}$  and  $\widetilde{\mathcal{D}}$  is. This is useful to analyze how privacy parameter  $\epsilon$  deteriorates the utility of data for downstream learning tasks.

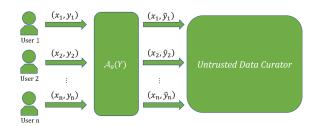


Figure 1: The framework of local label differential privacy.

Notice that  $\widetilde{y}_i$ 's are generated locally and independently, therefore the randomized response mechanism only alters the conditional distribution of Y given X, whereas the marginal distribution of X stays unchanged. Hence, we can assume  $\widetilde{\mathcal{D}}$  is a set of i.i.d. realizations of  $(X, \widetilde{Y})$ , and it is straightforward to verify that the conditional distribution of  $\widetilde{Y}$  given X = x can be written as

$$\widetilde{\eta}(\boldsymbol{x}) = \mathbb{P}(\widetilde{Y} = 1 | \boldsymbol{X} = \boldsymbol{x}) = \theta \eta(\boldsymbol{x}) + (1 - \theta)(1 - \eta(\boldsymbol{x})).$$

Clearly,  $\mathcal{A}_{\theta}$  amplifies the uncertainty of the observed labels by shrinking  $\widetilde{\eta}(\boldsymbol{x})$  towards 1/2. Particularly,  $\widetilde{\eta}(\boldsymbol{X}) = 1/2$  almost surely when  $\theta = 1/2$ . In this case, the privatized dataset  $\widetilde{\mathcal{D}}$  conveys no information for learning the decision function.

**Lemma 1.** Suppose the randomized response mechanism  $\widetilde{Y} = \mathcal{A}_{\theta}(Y)$  satisfies  $\epsilon$ -LLDP with  $\theta = \exp(\epsilon)/(1 + \exp(\epsilon))$ , then for any  $\mathbf{x} \in \mathcal{X}$  it holds that

$$\mathcal{L}(\epsilon, \boldsymbol{x}) \leq D_{KL} \big( \mathbb{P}_{Y|\boldsymbol{X} = \boldsymbol{x}} \big| \mathbb{P}_{\widetilde{Y}|\boldsymbol{X} = \boldsymbol{x}} \big) \leq \mathcal{U}(\epsilon, \boldsymbol{x}),$$

where  $D_{KL}(\mathbb{P}_{Y|\mathbf{X}=\mathbf{x}}|\mathbb{P}_{\widetilde{Y}|\mathbf{X}=\mathbf{x}})$  denotes the Kullback-Leibler divergence (KL) divergence between  $\mathbb{P}_{\widetilde{Y}|\mathbf{X}=\mathbf{x}}$  and  $\mathbb{P}_{Y|\mathbf{X}=\mathbf{x}}$ ,  $\mathcal{L}(\epsilon, \mathbf{x}) = 2(2\eta(\mathbf{x}) - 1)^2(1 + \exp(\epsilon))^{-2}$ , and  $\mathcal{U}(\epsilon, \mathbf{x}) = \min\{U_1(\epsilon, \mathbf{x}), U_2(\epsilon, \mathbf{x})\}$  with  $U_1(\epsilon, \mathbf{x}) = (2\eta(\mathbf{x}) - 1)^2\eta^{-1}(\mathbf{x})(1 - \eta(\mathbf{x}))^{-1}(1 + \exp(\epsilon))^{-2}$  and  $U_2(\epsilon, \mathbf{x}) = (2\eta(\mathbf{x}) - 1)^2\exp(-\epsilon)$ .

Lemma 1 quantifies the effect of the randomized response mechanism on the discrepancy between the conditional distributions of Y and  $\widetilde{Y}$  given the feature  $\boldsymbol{x}$  while  $\epsilon$ -Label DP is met. The lower bound  $\mathcal{L}(\epsilon, \boldsymbol{x})$  and the upper bound  $\mathcal{U}(\epsilon, \boldsymbol{x})$  share a factor  $(1 + \exp(-\epsilon))^2$ , indicating that  $D_{KL}(\mathbb{P}_{Y|\boldsymbol{X}=\boldsymbol{x}}|\mathbb{P}_{\widetilde{Y}|\boldsymbol{X}=\boldsymbol{x}})$  decreases exponentially with respect to  $\epsilon$ .

#### 3.2 Differentially Private Classifier

On the side of the untrusted curator, inference tasks vary according to the purpose of data collector, such as estimation of population statistics (Joseph et al., 2018; Yan et al., 2019) and supervised learning (Ghazi et al., 2021; Esfandiari et al., 2022). In this paper, we suppose that the computation based on  $\widetilde{\mathcal{D}}$  implemented

by the untrusted curator is formulated as the following regularized empirical risk minimization task,

$$\min_{f \in \mathcal{F}} L_n(f) = \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \phi(f(\boldsymbol{x}_i)\widetilde{y}_i) + \lambda_n J(f), \tag{1}$$

where  $\phi$  is a surrogate loss,  $\lambda_n$  is a tuning parameter,  $J(\cdot)$  is a penalty term, and  $\mathcal{F}$  is a pre-specified hypothesis space.

We denote by  $\widetilde{R}(f) = \mathbb{E}\big[\operatorname{sign}(f(\boldsymbol{X})) \neq \widetilde{Y}\big]$  the risk with expectation taken with respect to the joint distribution of  $(\boldsymbol{X}, \widetilde{Y})$  and let  $\widetilde{f}^* = \arg\min_f \widetilde{R}(f)$  denote the Bayes decision rule under the distribution of  $(\boldsymbol{X}, \widetilde{Y})$ . The excess risk of f under the distributions of  $(\boldsymbol{X}, Y)$  and  $(\boldsymbol{X}, \widetilde{Y})$  is denoted as  $D(f, f^*) = R(f) - R(f^*)$  and  $\widetilde{D}(f, \widetilde{f}^*) = \widetilde{R}(f) - \widetilde{R}(\widetilde{f}^*)$ , respectively.

**Lemma 2.** If  $\widetilde{Y} = \mathcal{A}_{\theta}(Y)$  with  $\theta > 1/2$ , then  $f^*(\boldsymbol{x}) = \widetilde{f}^*(\boldsymbol{x})$  for any  $\boldsymbol{x} \in \mathcal{X}$  and  $\widetilde{D}(f, \widetilde{f}^*) = (2\theta - 1)D(f, f^*)$  for any f.

Lemma 2 shows that the Bayes decision rule stays invariant under the randomized response mechanism. It is clear to see that  $\widetilde{D}(f,f^*)$  diminishes as  $\theta$  gets close to 1/2. Particularly,  $\widetilde{D}(f,f^*)$  is equal to 0 when  $\theta=1/2$ . This is as expected since  $\theta=1/2$  implies that  $\widetilde{\eta}(\boldsymbol{X})=1/2$  almost surely, where all classifiers deteriorates in performance simultaneously. Additionally, Lemma 2 theoretically validates the empirical concern (Busa-Fekete et al., 2021) that privatized dataset generated by the randomized response mechanism might still maintains the utility for learning a consistent classifier in downstream binary classification tasks (Beimel et al., 2013). Therefore, the estimated classifier can be used to infer the labels in the training dataset and leads to privacy leakage of labels (Busa-Fekete et al., 2021; Wu et al., 2022).

Let  $f_{\phi}^* = \arg\min_f R_{\phi}(f)$  and  $\widetilde{f}_{\phi}^* = \arg\min_f \widetilde{R}_{\phi}(f)$  be the minimizers of  $\phi$ -risk under the joint distributions of (X,Y) and  $(X,\widetilde{Y})$ , respectively. For illustration, we only consider hinge loss  $\phi(x) = \max\{1-x,0\}$  in this paper, and similar results can be easily obtained for other loss functions by the comparison theorem (Bartlett et al., 2006; Bartlett & Wegkamp, 2008). With  $\phi(\cdot)$  being the hinge loss, we have  $\widetilde{f}_{\phi}^* = f_{\phi}^* = f^* = \widetilde{f}^*$  (Bartlett et al., 2006; Lecué, 2007). With a slight abuse of notation, we use  $f^*$  to refer to these four functions simultaneously in the sequel.

## 3.3 Consistency under the Randomized Response Mechanism

In this section, we establish the asymptotic behavior of the classifier trained on privatized datasets. Specifically, we theoretically quantify how the randomized response mechanism affects the convergence rate of excess risk under the low-noise assumption (Lecué, 2007).

We suppose that the randomized response mechanism satisfies the  $\epsilon$ -LLDP, which indicates that  $\epsilon = \log(\theta/(1-\theta))$ . Furthermore, we denote that  $\widetilde{f}_n = \arg\min_{f \in \mathcal{F}} L_n(f)$  and  $H(f, f^*) = R_{\phi}(f) - R_{\phi}(f^*)$ . In classification problems,  $H(f, f^*)$  is an important metric that admits the decomposition into the estimation error and the approximation error (Bartlett et al., 2006),

$$H(f, f^*) = H(f, f_{\mathcal{F}}^*) + H(f_{\mathcal{F}}^*, f^*),$$

where  $f_{\mathcal{F}}^* = \arg\min_{f \in \mathcal{F}} R_{\phi}(f)$ . Here  $H(f, f_{\mathcal{F}}^*)$  is the estimation error and  $H(f_{\mathcal{F}}^*, f^*)$  is the approximation error. The estimation error depends on the learning algorithm in finding  $f_{\mathcal{F}}^*$  based on a dataset with finite samples, where the searching difficulty is unavoidably affected by the complexity of  $\mathcal{F}$  as the approximation error. Generally, the richness of  $\mathcal{F}$  can be measured by VC dimension (Blumer et al., 1989; Vapnik & Chervonenkis, 2015), Rademacher complexity (Bartlett & Mendelson, 2002; Smale & Zhou, 2003), and metric entropy methods (Zhou, 2002; Shalev-Shwartz & Ben-David, 2014).

**Assumption 1.** (Low-noise assumption) There exists a constant c>0 and  $0<\gamma\leq +\infty$  such that  $\mathbb{P}\Big(|2\eta(\boldsymbol{X})-1|\leq t\Big)\leq ct^{\gamma}$ , for any  $t\in[0,1)$ .

Assumption 1 is known as the low-noise assumption in the binary classification (Lecué, 2007; Shen et al., 2003; Bartlett et al., 2006), which characterizes the behavior of  $2\eta(\mathbf{x}) - 1$  around the decision boundary

 $\eta(\boldsymbol{x}) = 1/2$ . Particularly, the case with  $\gamma = +\infty$  and c = 1 implies that the labels  $y_i$ 's are deterministic in that  $\eta(\boldsymbol{X})$  takes values in  $\{0,1\}$  almost surely, resulting in the fastest convergence rate of the estimation error.

**Lemma 3.** Denote that  $\kappa_{\epsilon} = (e^{\epsilon} - 1)/(e^{\epsilon} + 1)$ . Under Assumption 1, for any  $\theta \in (1/2, 1]$ , it holds that

(1) 
$$\mathbb{P}\Big(|2\widetilde{\eta}(\boldsymbol{X})-1| \le t\Big) \le c\kappa_{\epsilon}^{-\gamma}t^{\gamma}$$
, for any  $t \in [0,1)$ ,

(2) 
$$\widetilde{R}_{\phi}(f) - \widetilde{R}_{\phi}(f^*) \ge \kappa_{\epsilon}(4c)^{-1/\gamma} \Big( \mathbb{E}\big[ |f(\boldsymbol{X}) - f^*(\boldsymbol{X})| \big] \Big)^{\frac{\gamma+1}{\gamma}} \text{ for any } f \in \mathcal{F},$$

Lemma 3 presents some insights regarding the influence of the randomized response mechanism on the low-noise assumption and the margin relation (Lecué, 2007), showing that the low-noise structure and margin relation are both invariant to the randomized response mechanism in the sense that only the multiplicative terms are enlarged by the effect of the privacy guarantee  $\epsilon$ .

**Assumption 2.** We assume that  $\mathcal{F}$  is properly chosen satisfying that  $||f||_{\infty} \leq 1$  for any  $f \in \mathcal{F}$  and

$$\log \mathcal{N}(\xi, \mathcal{F}, \|\cdot\|_{L^2(\mu)}) \le \mathcal{V}_1(\Theta) \log(1 + \xi^{-1} \mathcal{V}_2(\Theta)),$$

where  $\mu$  denotes the marginal distribution of X,  $\Theta$  denotes the parameters of f, and  $\mathcal{V}_1(\Theta)$  and  $\mathcal{V}_2(\Theta)$  are some functions depending on  $\Theta$ .

Assumption 2 characterizes the complexity of function class  $\mathcal{F}$  through the metric entropy (Zhou, 2002; Bousquet et al., 2003; Lei et al., 2016), where  $\mathcal{V}_1(\Theta)$  and  $\mathcal{V}_2(\Theta)$  are quantities increasing with the size of  $\Theta$ . Assumption 2 generally holds for function classes of parametric models (Wang et al., 2016a; Xu et al., 2021), most notably for deep neural networks (Bartlett et al., 2019; Schmidt-Hieber, 2020). Additionally, Assumption 2 also holds for those VC classes with VC dimensions increasing with the size of  $\Theta$  (Wellner et al., 2013; Bartlett et al., 2019; Lee et al., 1994).

**Theorem 1.** Under Assumptions 1 and 2, for any minimizer  $\tilde{f}_n$  of (1), there exist some positive constants  $A_1$  and  $A_2$  such that

$$A_2\Big\{\Big(\frac{\mathcal{V}_1(\Theta)}{n\kappa_{\epsilon}^2}\Big)^{\frac{\gamma+1}{\gamma+2}} + \tau_n\Big\} \leq \sup_{\pi \in \mathcal{P}_{\gamma}} \mathbb{E}_{\widetilde{\mathcal{D}}}\big[R(\widetilde{f}_n) - R(f^*)\big] \leq A_1\Big\{\Big(\frac{\mathcal{V}_1(\Theta)\log(n)}{n\kappa_{\epsilon}^2}\Big)^{\frac{\gamma+1}{\gamma+2}} + s_n\Big\},$$

where  $\mathcal{P}_{\gamma}$  be a class of distributions of  $(\mathbf{X}, Y)$  satisfying Assumption 1,  $s_n = \sup_{\pi \in \mathcal{P}_{\gamma}} \inf_{f \in \mathcal{F}} H(f, f^*)$ , and  $\tau_n = \sup_{\pi \in \mathcal{P}_{\gamma}} \inf_{f \in \mathcal{F}} D(f, f^*)$ .

Theorem 1 quantifies the asymptotic behavior of  $\widetilde{f}_n$  by establishing its upper and lower bounds, which explicitly demonstrates the quantitative relation between the effect of privacy guarantee  $\epsilon$  and the excess risk. Particularly, the upper bound matches the lower bound except for a logarithmic factor when the approximation error term  $s_n \lesssim \left(\frac{\mathcal{V}_1(\Theta)}{n\kappa_\epsilon^2}\right)^{\frac{\gamma+1}{\gamma+2}}$ . This shows that the randomized response mechanism slows down the convergence rate by enlarging the multiplicative constant. Moreover, based on Theorem 1, we can obtain the optimal convergence rate of the excess risk in classification problem under the low-noise assumption (Lecué, 2007) by setting  $\epsilon = \infty$  and  $|\mathcal{F}| < \infty$ .

# 4 Deep Learning with Label Differential Privacy

A typical class of models that are popularly considered in the domain of differential privacy is deep neural network (Ghazi et al., 2021; Yuan et al., 2021) due to its success in various applications in the past decade. Unlike Section 3.3 considering the estimation and approximation errors separately, we establish theoretical results regarding the convergence rate of the excess risk of the deep neural network plug-in classifier trained from  $\widetilde{\mathcal{D}}$ , which is obtained by making a tradeoff between the estimation and approximation errors of deep neural networks (Schmidt-Hieber, 2020). Our theoretical results not only quantify how the optimal structure

of the deep neural network changes with the privacy parameter  $\epsilon$ , but also derive the optimal privacy guarantee we can achieve for the deep neural network classifier.

Remind that  $\widetilde{D} = \{(\boldsymbol{x}_i, \widetilde{y}_i)\}_{i=1}^n$  is the privatized dataset with  $\widetilde{y}_i = \mathcal{A}_{\theta}(y_i)$  and  $\theta = \exp(\epsilon)/(1 + \exp(\epsilon))$ . The deep neural network is solved as

$$\widetilde{f}_{nn} = \operatorname*{arg\,min}_{f \in \mathcal{F}_n^{NN}(L_n, N_n, P_n, B_n, V_n)} \frac{1}{n} \sum_{i=1}^n \left( f(\boldsymbol{x}_i) - \widetilde{z}_i \right)^2, \tag{2}$$

where  $\tilde{z}_i = (\tilde{y}_i + 1)/2$  and  $\mathcal{F}_n^{NN}(L_n, N_n, P_n, B_n, V_n)$  is a class of multilayer perceptrons defined in Section 2.2. The plug-in classifier based on  $\tilde{f}_{nn}$  can be obtained as  $\tilde{s}_{nn} = \text{sign}(\tilde{f}_{nn} - 1/2)$ . To quantify the asymptotic behavior of  $\tilde{s}_{nn}$ , we further assume that the support of  $\boldsymbol{x}$  is  $[0,1]^p$ , which is a common assumption for deep neural network s(Yarotsky, 2017; Nakada & Imaizumi, 2020)

**Theorem 2.** Let  $\mathcal{P}_{\gamma,\beta}$  be a class of probability measures on  $\mathcal{X} \times \{-1,1\}$  satisfying Assumption 1 and  $\eta(\mathbf{X}) \in \mathcal{H}(\beta, [0,1]^p, M)$ . For any minimizer  $\widetilde{f}_{nn}$  in (2) with  $L_n \asymp \log(\kappa_{\epsilon} n/\log(n))$ ,  $N_n \asymp (\kappa_{\epsilon} n/\log(n))^{\frac{2p}{2\beta+p}}$ ,  $B_n = 1$ , and  $P_n \asymp N_n \log(\kappa_{\epsilon} n/\log(n))$ . Then we have

$$\left(\frac{1}{n\kappa_{\epsilon}^{2}}\right)^{\frac{\beta(\gamma+1)}{\beta(\gamma+2)+p}} \lesssim \sup_{\pi \in \mathcal{P}_{\gamma,\beta}} \mathbb{E}_{\widetilde{\mathcal{D}}}\left[R(\widetilde{s}_{nn}) - R(f^{*})\right] \lesssim \left(\frac{\log n}{n\kappa_{\epsilon}^{2}}\right)^{\frac{2\beta(\gamma+1)}{2\beta(\gamma+2)+p(\gamma+2)}}.$$
(3)

Particularly,  $\sup_{\pi \in \mathcal{P}_{\gamma,\beta}} \mathbb{E}_{\widetilde{\mathcal{D}}}[R(\widetilde{s}_{nn}) - R(f^*)] = o(1)$  given that  $\epsilon \gtrsim \log^{\zeta}(n)/\sqrt{n}$  for any  $\zeta > 0$ .

In Theorem 2, we quantify the asymptotic behavior of the excess risk of  $\widetilde{s}_{nn}$  by providing upper and lower bounds for  $\sup_{\pi \in \mathcal{P}_{\gamma,\beta}} \mathbb{E}_{\widetilde{\mathcal{D}}}[R(\widetilde{s}_{nn}) - R^*]$ . It should be noted that if  $\epsilon = \infty$  which refers to the non-private case, the upper and lower bounds in (3) match with existing theoretical results established in Audibert & Tsybakov (2007). Additionally, Theorem 2 explicitly characterizes the effect of  $\epsilon$  on the convergence of  $\widetilde{s}_{nn}$  to the Bayes decision rule and specifies how the structure of the optimal neural network shrinks as  $\epsilon$  decreases for achieving the fastest convergence rate. Furthermore, from Theorem 2, we can derive the fastest adaptive rate of  $\epsilon$  under the consistency of  $\widetilde{s}_{nn}$ , i.e.,  $\epsilon \gtrsim \log^{\zeta}(n)n^{-1/2}$  for any  $\zeta > 0$ .

# 5 Simulated Experiments

This section aims to validate our theoretical results through extensive simulated examples. Specifically, we show that the excess risk of a differentially private classifier converges to 0 for any fixed  $\epsilon$ , whereas the convergence is not achievable as long as  $\epsilon$  is adaptive to the size of the training dataset with some properly chosen orders as shown in Theorems 1 and 2.

# 5.1 Support Vector Machine

This simulation experimentally analyzes the effect of label DP on the SVM classifier. The generation of simulated datasets is as follows. First, we set the regression function in classification as  $\eta(x) = 1/(1 + \exp(-\beta_0^T x))$ , where  $\beta_0 \in \mathbb{R}^p$  and x are both p-dimensional vectors generated via  $\beta_{0i}, x_i \sim \text{Unif}(-1, 1)$  for  $i = 1, \ldots, p$ . For each feature x, its label y is chosen from  $\{1, -1\}$  with probabilities  $\eta(x)$  and  $1 - \eta(x)$ , respectively. Repeating the above process n times, we obtain a non-private training dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ . Subsequently, we apply the randomized response mechanism to generate  $\tilde{y}_i$  as  $\tilde{y}_i = \mathcal{A}_{\theta}(y_i)$ , where  $\theta = \exp(\epsilon)/(1 + \exp(\epsilon))$  and  $\epsilon$  is the privacy guarantee. Therefore, the obtained privatized training dataset  $\tilde{\mathcal{D}} = \{(x_i, \tilde{y}_i)\}_{i=1}^n$  satisfies  $\epsilon$ -LLDP. Based on  $\tilde{\mathcal{D}}$ , we obtain the SVM classifier (Cortes & Vapnik, 1995),

$$\widetilde{f}_n = \underset{\boldsymbol{\beta}}{\operatorname{arg\,min}} \frac{1}{n} \sum_{i=1}^n (1 - \widetilde{y}_i \boldsymbol{\beta}^T \boldsymbol{x}_i)_+ + \lambda \|\boldsymbol{\beta}\|_2^2,$$

where  $(x)_{+} = \max\{x, 0\}$ . Next, we evaluate the performance of  $\widetilde{f}_n$  in terms of the empirical excess risk and the classification error,

$$\widehat{E}(\widetilde{f}_n) = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} I\Big(\operatorname{sign}(\widetilde{f}_n(\boldsymbol{x}_i')) \neq \operatorname{sign}(\eta(\boldsymbol{x}_i') - 1/2)\Big) |2\eta(\boldsymbol{x}_i') - 1|,$$

$$\operatorname{CE}(\widetilde{f}_n) = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} I\Big(\operatorname{sign}(\widetilde{f}_n(\boldsymbol{x}_i')) \neq \operatorname{sign}(\eta(\boldsymbol{x}_i') - 1/2)\Big).$$

where  $x_i'$ 's are testing samples generated in the same way as  $x_i$ 's.

**Scenario I.** In the first scenario, we aim to verify that  $\widehat{E}(\widetilde{f}_n)$  will converge to 0 as sample size n increases when the privacy parameter is a fixed constant. To this end, we consider cases  $(n, \epsilon) = \{100 \times 2^i, i = 0, 1, \dots, 8\} \times \{1, 2, 3, 4, \infty\}$ .

**Scenario II**. In the second scenario, we explore the asymptotic behavior of  $\widehat{E}(\widetilde{f}_n)$  with  $\epsilon$  adaptive to the sample size n. Specifically, we set  $\epsilon = 5n^{-\zeta}$  and consider cases  $(n,\zeta) = \{100 \times 2^i, i = 0, 1, \dots, 8\} \times \{1/5, 1/4, 1/3, 1/2, 2/3, 1\}$ . We also include the worst case  $\epsilon = 0$  as a baseline.

**Scenario III.** In the third scenario, we intend to verify that  $\epsilon \approx n^{-1/2}$  is the dividing line between whether or not the excess risk converges. To this end, we consider three kinds of adaptive  $\epsilon$ , including  $\epsilon \approx n^{-1/2}$ ,  $\epsilon \approx \log(n)n^{-1/2}$ , and  $\epsilon \approx \log^{-1}(n)n^{-1/2}$ . The size of  $\mathcal{D}$  is set as  $\{100 \times 3^i, i = 0, 1, \dots, 7\}$ . For all cases, we report the averaged empirical excess risk in 1,000 replications as well as their 95% confidence intervals.

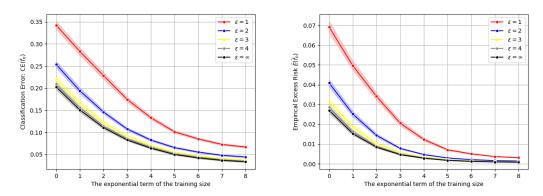
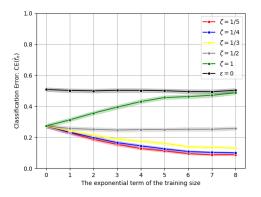


Figure 2: The averaged classification errors (Left) and averaged empirical excess risks (Right) of all settings with  $n_{test} = 50,000$  in Scenario I.

For Scenario I and Scenario II, we report the averaged empirical excess risk and the classification error in 1,000 replications for each setting in Figure 2 and Figure 3, respectively. From the left panel of Figure 2, we can see that the empirical excess risks and the classification errors with a fixed  $\epsilon$  converge to 0 regardless of the value of  $\epsilon$ , showing that the randomized response mechanism with a fixed  $\epsilon$  fails to prevent the third party from learning the optimal classifier based on  $\widetilde{\mathcal{D}}$ . Moreover, as seen from Figure 3, when  $\zeta < 1/2$  the estimated excess risks present a decreasing pattern as the sample size increases, whereas that of the case  $\zeta = 1$  deteriorates steadily and the curve finally overlaps with that of the worst case  $\epsilon = 0$ . It is also interesting to observe that the curve of  $\zeta = 1/2$  remains unaffected by the sample size. All these phenomenons are in accordance with the results of Theorem 1.

As can be seen in Figure 4, the curve of the case  $\epsilon \approx n^{-1/2}$  remains unchanged as sample size increases as in Scenario II. This is due to the offset of information gain yielded by increasing the sample size and the information loss in the label-flipping mechanism. As expected, the additional logarithmic term significantly alters the original curve pattern. Specifically, for the case with  $\epsilon \approx \log^{-1}(n)n^{-1/2}$ , the performance of  $\tilde{f}_n$  deteriorates significantly and approaches the worst case with  $\epsilon = 0$ . On the contrary, the performance of  $\tilde{f}_n$  in the case  $\epsilon \approx \log(n)n^{-1/2}$  improves significantly with  $\hat{E}(\tilde{f}_n)$  converging to 0. Therefore,  $\epsilon \approx n^{-1/2}$  appears



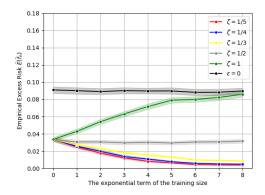
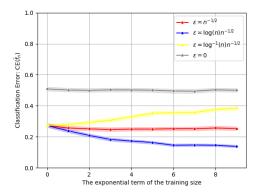


Figure 3: The averaged classification errors (Left) and averaged empirical excess risks (Right) of all settings with  $n_{test} = 50,000$  in Scenario II.



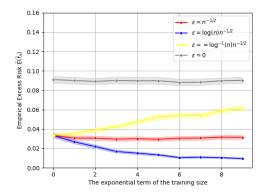


Figure 4: The averaged classification errors (Left) and averaged empirical excess risks (Right) of all settings with  $n_{test} = 50,000$  in Scenario III.

to be the dividing line that determines whether the excess risk converges, which completely matches our theoretical results.

#### 5.2 Deep Neural Network Classifier

The simulation in this section aims to verify our theoretical results in Theorem 2, which mainly lies in two aspects. First, we intend to verify the effect of  $\epsilon$  of Label DP on the optimal structure of deep neural networks for classification problems. Specifically, as stated in Theorem 2, if label noise yielded by the randomized response mechanism increases, a smaller deep neural network should be employed to strike a better balance between the approximation error and the estimation error to achieve a better generalization error. Second, the consistency in estimating the decision boundary is prohibited provided that  $\epsilon$  is adaptive to the training size n at some specific orders. To these ends, we consider generating training datasets  $\widetilde{\mathcal{D}} = \{(x_i, \widetilde{y}_i)\}_{i=1}^n$  as follows. First, we set the regression function as  $\eta_{nn}(x_i) = \sum_{j=1}^4 \sin(2\pi x_{ij})/8 + 1/2$  with  $x_{ij} \sim \text{Unif}(0,1)$  for any i, j. Then we generate  $y_i$  from  $\{1, -1\}$  with probabilities  $\eta_{nn}(x_i)$  and  $1 - \eta_{nn}(x_i)$ , respectively. As in the last simulation, we then apply the randomized response mechanism  $\mathcal{A}_{\theta}$  to each  $y_i$  to generate  $\widetilde{y}_i = \mathcal{A}_{\theta}(y_i)$  with  $\theta = \exp(\epsilon)/(1 + \exp(\epsilon))$ . Then we set the hypothesis space to be the class of L-layer fully connected neural network with equal width and the ReLU activation function, where h denotes the widths in all hidden layers.

The overall training process of the neural network is implemented in Tensorflow (Abadi et al., 2016) with the Adam optimizer and learning rate being 0.001. Additionally, we employ the early-stopping technique to

monitor the training error with patience 10 and maintain the parameter with the smallest training error. Let  $\tilde{f}_{nn}$  denote the resultant neural network obtained from minimizing (2). We construct the associated plug-in classifier as  $\tilde{s}_{nn} = \text{sign}(\tilde{f}_{nn} - 1/2)$  and evaluate its performance by the empirical excess risk and the classification error as in Section 5.1.

Scenario I. In the first scenario, we consider privacy guarantees  $\epsilon \in \{1, 2, \infty\}$  and neural network structures with L=2 and  $h \in \{8, 12, 16, 20, 24\}$  with . We report the averaged empirical excess risks and the classification errors of all cases in 100 replications as well as their 95% confidence intervals in Figure 5. Clearly, if  $\epsilon=1$ , the optimal neural network structure is h=8, whereas those of the cases  $\epsilon=2$  and  $\epsilon=\infty$  are h=20 and h=24, respectively. Such results show that a smaller neural network is preferred when stronger privacy protection of labels (smaller  $\epsilon$ ) is considered, which coincides with our theoretical results in Theorem 2 that the optimal neural network structure to achieve the fastest convergence rate of excess risk should diminish as  $\epsilon$  decreases. Moreover, as the training sample size n increases from 2,000 to 4,000, the optimal structure of the neural network enlarges for the cases  $\epsilon=1$  and  $\epsilon=2$  due to their new tradeoffs between the estimation and approximation errors.

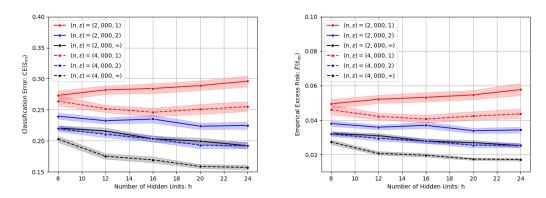


Figure 5: The averaged classification errors (Left) and averaged empirical excess risks (Right) of all cases with  $n_{test} = 50,000$ .

Scenario II. In the second scenario, we fix the neural network structure as (L,h)=(2,16) and consider training sample sizes  $n \in \{2^i \times 10^3; i=0,1,2,3,4\}$ . To verify our theoretical results, we compare two privacy schemes  $\epsilon \simeq n^{-1/2}$  and  $\epsilon \simeq \log(n)n^{-1/2}$ . We also include the case with invariant privacy scheme  $\epsilon=4$  as a baseline. For comparing their difference in trend patterns, the multiplicative constants of two adaptive privacy schemes are chosen such that they have the same starting point as  $\epsilon=4$  when n=1,000. We report the averaged empirical excess risks, the classification errors of all cases in 100 replications, and their 95% confidence intervals in Figure 6. Clearly, the performances of the cases  $\epsilon=4$  and  $\epsilon \simeq \log(n)n^{-1/2}$  improve as n increases. However, the performance of case  $\epsilon \simeq n^{-1/2}$  presents a different pattern in generalization performance as n increases. Most notably, as n increases from 8,000 to 16,000, it performance deteriorates while the other two cases still observe significant improvements, showing that the additional logarithmic term plays a deterministic role in the convergence of excess risk, which perfectly aligns with our theoretical findings in Theorem 2.

# 6 Real Applications - Mnist Dataset

This experiment considers similar settings of the privacy parameter  $\epsilon$  as Section 5.2 in order to verify our theoretical findings of DNN on the MNIST dataset (LeCun, 1998). The MNIST dataset consists of 60,000 training images and 10,000 testing images, and each sample is a  $28 \times 28$  grey-scale pixel image of one of the 10 digits. In this experiment, we consider a binary classification problem by only including samples of digits 2 and 3 for training and testing due to their similarity in appearance. The resultant training dataset contains 12,089 training samples and 2,042 testing samples.

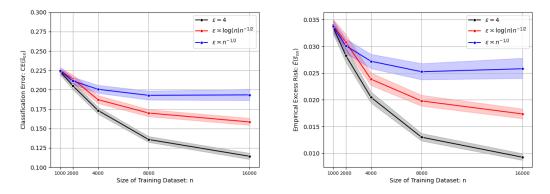


Figure 6: The averaged classification errors (Left) and averaged empirical excess risks (Right) of all cases with  $n_{test} = 50,000$ .

For the hyperparameters setting, we consider the neural network with three convolution layers with ReLU activation and two fully-connected layers. For the three convolution layers, we set their kernel sizes and numbers of channels as  $4 \times 4$  and 4, respectively. Additionally, each convolution layer is followed by a max pooling layer with size  $2 \times 2$ . The first fully-connected layer has 10 hidden units with ReLU activation, and the last layer outputs the probability of an image being digit 2. As in Section 5.2, the neural network is trained with the Adam optimizer and learning rate being 0.001, and the early-stopping technique to monitor the training error with patience 10 and maintain the parameter with the smallest training error. We evaluate of the trained model by the testing error on 2,042 testing samples. We consider training sample size  $n \in \{2 \times i \times 10^3; i = 1, 2, 3, 4, 5\}$ . We mainly consider two scenarios, including the fixed privacy guarantee with  $\epsilon \in \{1, 2, \infty\}$  and the adaptive privacy schemes with  $\epsilon \approx \log(n)\sqrt{n^{-1}}$ ,  $\epsilon \approx n^{-1/2}$ , and  $\epsilon \approx \log^{-1}(n)n^{-1/2}$ . The averaged testing error of each case in 50 replications is reported in Figure 7.

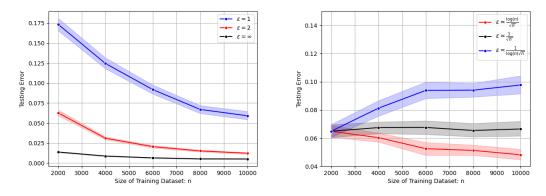


Figure 7: The averaged testing errors with fixed  $\epsilon$  (Left) and adaptive  $\epsilon$  (Right) under different training sample sizes in MNIST dataset

Figure 7 presents similar results as in Section 5.2. First, when  $\epsilon$  is fixed, differentially private classifiers improve in generalization performance as the training sample size increases and attain competitive performance as the non-private classifier ( $\epsilon = \infty$ ) when the training sample size is large enough. This result accords with our theoretical findings in Theorem 1 that fixed privacy guarantee in the label DP slows down the convergence to the optimal classifier with a multiplicative constant. In stark contrast, as shown in the right plot of Figure 7, the convergence to the optimal classifier is prevented if  $\epsilon \approx n^{-1/2}$ , as boosting the training size from 2,000 to 10,000 fails to significantly improve the testing error. This again aligns with our Theorem 1

# References

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. {TensorFlow}: a system for {Large-Scale} machine learning. In 12th USENIX symposium on operating systems design and implementation (OSDI 16), pp. 265–283, 2016.
- Mário S Alvim, Miguel E Andrés, Konstantinos Chatzikokolakis, Pierpaolo Degano, and Catuscia Palamidessi. Differential privacy: on the trade-off between utility and information leakage. In *International Workshop on Formal Aspects in Security and Trust*, pp. 39–54. Springer, 2012.
- Pathum Chamikara Mahawaga Arachchige, Peter Bertok, Ibrahim Khalil, Dongxi Liu, Seyit Camtepe, and Mohammed Atiquzzaman. Local differential privacy for deep learning. *IEEE Internet of Things Journal*, 7(7):5827–5842, 2019.
- Sanjeev Arora, László Babai, Jacques Stern, and Z Sweedyk. The hardness of approximate optima in lattices, codes, and systems of linear equations. *Journal of Computer and System Sciences*, 54(2):317–331, 1997.
- 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 and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- Peter L Bartlett and Marten H Wegkamp. Classification with a reject option using a hinge loss. *Journal of Machine Learning Research*, 9(8), 2008.
- 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.
- Peter L Bartlett, Nick Harvey, Christopher Liaw, and Abbas Mehrabian. Nearly-tight vc-dimension and pseudodimension bounds for piecewise linear neural networks. *The Journal of Machine Learning Research*, 20(1):2285–2301, 2019.
- Raef Bassily, Mehryar Mohri, and Ananda Theertha Suresh. Differentially private learning with margin guarantees. arXiv preprint arXiv:2204.10376, 2022.
- Amos Beimel, Kobbi Nissim, and Uri Stemmer. Private learning and sanitization: Pure vs. approximate differential privacy. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pp. 363–378. Springer, 2013.
- Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K Warmuth. Learnability and the vapnik-chervonenkis dimension. *Journal of the ACM (JACM)*, 36(4):929–965, 1989.
- Olivier Bousquet, Stéphane Boucheron, and Gábor Lugosi. Introduction to statistical learning theory. In Summer school on machine learning, pp. 169–207. Springer, 2003.
- Robert Istvan Busa-Fekete, Umar Syed, Sergei Vassilvitskii, et al. Population level privacy leakage in binary classification with label noise. In NeurIPS 2021 Workshop Privacy in Machine Learning, 2021.
- Timothy I Cannings, Yingying Fan, and Richard J Samworth. Classification with imperfect training labels. Biometrika, 107(2):311-330, 04 2020. doi: 10.1093/biomet/asaa011. URL https://doi.org/10.1093/biomet/asaa011.
- Olivier Chapelle, Eren Manavoglu, and Romer Rosales. Simple and scalable response prediction for display advertising. ACM Transactions on Intelligent Systems and Technology (TIST), 5(4):1–34, 2014.
- Kamalika Chaudhuri and Daniel Hsu. Sample complexity bounds for differentially private learning. In *Proceedings of the 24th Annual Conference on Learning Theory*, pp. 155–186. JMLR Workshop and Conference Proceedings, 2011.

- Corinna Cortes and Vladimir Vapnik. Support-vector networks. Machine learning, 20(3):273–297, 1995.
- Teddy Cunningham, Konstantin Klemmer, Hongkai Wen, and Hakan Ferhatosmanoglu. Geopointgan: Synthetic spatial data with local label differential privacy. arXiv preprint arXiv:2205.08886, 2022.
- John C Duchi and Feng Ruan. The right complexity measure in locally private estimation: It is not the fisher information. arXiv preprint arXiv:1806.05756, 2018.
- John C Duchi, Michael I Jordan, and Martin J Wainwright. Minimax optimal procedures for locally private estimation. *Journal of the American Statistical Association*, 113(521):182–201, 2018.
- Cynthia Dwork. Differential privacy: A survey of results. In *International Conference on Theory and Applications of Models of Computation*, pp. 1–19. Springer, 2008.
- Ülfar Erlingsson, Vasyl Pihur, and Aleksandra Korolova. Rappor: Randomized aggregatable privacy-preserving ordinal response. In Proceedings of the 2014 ACM SIGSAC conference on computer and communications security, pp. 1054–1067, 2014.
- Úlfar Erlingsson, Vitaly Feldman, Ilya Mironov, Ananth Raghunathan, Kunal Talwar, and Abhradeep Thakurta. Amplification by shuffling: From local to central differential privacy via anonymity. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 2468–2479. SIAM, 2019.
- Hossein Esfandiari, Vahab Mirrokni, Umar Syed, and Sergei Vassilvitskii. Label differential privacy via clustering. In *International Conference on Artificial Intelligence and Statistics*, pp. 7055–7075. PMLR, 2022.
- Badih Ghazi, Noah Golowich, Ravi Kumar, Pasin Manurangsi, and Chiyuan Zhang. Deep learning with label differential privacy. Advances in Neural Information Processing Systems, 34:27131–27145, 2021.
- Badih Ghazi, Pritish Kamath, Ravi Kumar, Ethan Leeman, Pasin Manurangsi, Avinash Varadarajan, and Chiyuan Zhang. Regression with label differential privacy. arXiv preprint arXiv:2212.06074, 2022.
- Antonious Girgis, Deepesh Data, Suhas Diggavi, Peter Kairouz, and Ananda Theertha Suresh. Shuffled model of differential privacy in federated learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 2521–2529. PMLR, 2021.
- Matthew Joseph, Aaron Roth, Jonathan Ullman, and Bo Waggoner. Local differential privacy for evolving data. Advances in Neural Information Processing Systems, 31, 2018.
- Peter Kairouz, Sewoong Oh, and Pramod Viswanath. Extremal mechanisms for local differential privacy. The Journal of Machine Learning Research, 17(1):492–542, 2016.
- Vishesh Karwa, Pavel N Krivitsky, and Aleksandra B Slavković. Sharing social network data: differentially private estimation of exponential family random-graph models. *Journal of the Royal Statistical Society:* Series C (Applied Statistics), 66(3):481–500, 2017.
- Vladimir Koltchinskii. Oracle inequalities in empirical risk minimization and sparse recovery problems: École D'Été de Probabilités de Saint-Flour XXXVIII-2008, volume 2033. Springer Science & Business Media, 2011.
- Guillaume Lecué. Optimal rates of aggregation in classification under low noise assumption. *Bernoulli*, 13 (4):1000–1022, 2007.
- Yann LeCun. The mnist database of handwritten digits. http://yann. lecun. com/exdb/mnist/, 1998.
- Wee Sun Lee, Peter L Bartlett, and Robert C Williamson. Lower bounds on the vc-dimension of smoothly parametrized function classes. In *Proceedings of the seventh annual conference on Computational learning theory*, pp. 362–367, 1994.

- Yunwen Lei, Lixin Ding, and Yingzhou Bi. Local rademacher complexity bounds based on covering numbers. Neurocomputing, 218:320–330, 2016.
- Yi Lin. A note on margin-based loss functions in classification. Statistics & Probability Letters, 68(1):73–82, 2004.
- Mani Malek Esmaeili, Ilya Mironov, Karthik Prasad, Igor Shilov, and Florian Tramer. Antipodes of label differential privacy: Pate and alibi. *Advances in Neural Information Processing Systems*, 34:6934–6945, 2021.
- H Brendan McMahan, Gary Holt, David Sculley, Michael Young, Dietmar Ebner, Julian Grady, Lan Nie, Todd Phillips, Eugene Davydov, Daniel Golovin, et al. Ad click prediction: a view from the trenches. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 1222–1230, 2013.
- Frank McSherry and Ilya Mironov. Differentially private recommender systems: Building privacy into the netflix prize contenders. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 627–636, 2009.
- Ryumei Nakada and Masaaki Imaizumi. Adaptive approximation and generalization of deep neural network with intrinsic dimensionality. J. Mach. Learn. Res., 21(174):1–38, 2020.
- Tapan K Nayak and Samson A Adeshiyan. A unified framework for analysis and comparison of randomized response surveys of binary characteristics. *Journal of Statistical Planning and Inference*, 139(8):2757–2766, 2009.
- Igal Sason and Sergio Verdú. f-divergence inequalities. *IEEE Transactions on Information Theory*, 62(11): 5973–6006, 2016.
- Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with relu activation function. *The Annals of Statistics*, 48(4):1875–1897, 2020.
- Shai Shalev-Shwartz and Shai Ben-David. Understanding machine learning: From theory to algorithms. Cambridge university press, 2014.
- Xiaotong Shen and Wing Hung Wong. Convergence rate of sieve estimates. *The Annals of Statistics*, pp. 580–615, 1994.
- Xiaotong Shen, George C Tseng, Xuegong Zhang, and Wing Hung Wong. On  $\psi$ -learning. Journal of the American Statistical Association, 98(463):724–734, 2003.
- Steve Smale and Ding-Xuan Zhou. Estimating the approximation error in learning theory. Analysis and Applications, 1(01):17–41, 2003.
- Jun Tang, Aleksandra Korolova, Xiaolong Bai, Xueqiang Wang, and Xiaofeng Wang. Privacy loss in apple's implementation of differential privacy on macos 10.12. arXiv preprint arXiv:1709.02753, 2017.
- Alexander B Tsybakov. Optimal aggregation of classifiers in statistical learning. *The Annals of Statistics*, 32(1):135–166, 2004.
- Brendan van Rooyen and Robert C. Williamson. A theory of learning with corrupted labels. *Journal of Machine Learning Research*, 18(228):1-50, 2018. URL http://jmlr.org/papers/v18/16-315.html.
- Vladimir N Vapnik and A Ya Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. In *Measures of complexity*, pp. 11–30. Springer, 2015.
- Di Wang and Jinhui Xu. On sparse linear regression in the local differential privacy model. In *International Conference on Machine Learning*, pp. 6628–6637. PMLR, 2019.

Junhui Wang, Xiaotong Shen, Yiwen Sun, and Annie Qu. Classification with unstructured predictors and an application to sentiment analysis. *Journal of the American Statistical Association*, 111(515):1242–1253, 2016a.

Tianhao Wang, Jeremiah Blocki, Ninghui Li, and Somesh Jha. Locally differentially private protocols for frequency estimation. In 26th USENIX Security Symposium (USENIX Security 17), pp. 729–745, 2017.

Yue Wang, Xintao Wu, and Donghui Hu. Using randomized response for differential privacy preserving data collection. In *EDBT/ICDT Workshops*, volume 1558, pp. 0090–6778, 2016b.

Stanley L Warner. Randomized response: A survey technique for eliminating evasive answer bias. *Journal* of the American Statistical Association, 60(309):63–69, 1965.

Jon Wellner et al. Weak convergence and empirical processes: with applications to statistics. Springer Science & Business Media, 2013.

Ruihan Wu, Jin Peng Zhou, Kilian Q Weinberger, and Chuan Guo. Does label differential privacy prevent label inference attacks? arXiv preprint arXiv:2202.12968, 2022.

Yu Xin and Tommi Jaakkola. Controlling privacy in recommender systems. Advances in neural information processing systems, 27, 2014.

Shirong Xu, Ben Dai, and Junhui Wang. Sentiment analysis with covariate-assisted word embeddings. *Electronic Journal of Statistics*, 15(1):3015–3039, 2021.

Ziqi Yan, Qiong Wu, Meng Ren, Jiqiang Liu, Shaowu Liu, and Shuo Qiu. Locally private jaccard similarity estimation. *Concurrency and Computation: Practice and Experience*, 31(24):e4889, 2019.

Dmitry Yarotsky. Error bounds for approximations with deep relu networks. *Neural Networks*, 94:103–114, 2017.

Sen Yuan, Milan Shen, Ilya Mironov, and Anderson CA Nascimento. Practical, label private deep learning training based on secure multiparty computation and differential privacy. Cryptology ePrint Archive, 2021.

Tong Zhang. Statistical behavior and consistency of classification methods based on convex risk minimization. *The Annals of Statistics*, 32(1):56–85, 2004.

Ding-Xuan Zhou. The covering number in learning theory. Journal of Complexity, 18(3):739–767, 2002.

## A Appendix

## A.1 Proofs

**Proof of Lemma 1**: We first prove the lower bound of  $D_{KL}\left(\mathbb{P}_{Y|X=x}|\mathbb{P}_{\widetilde{Y}|X=x}\right)$ , which is mainly based on the Pinsker's inequality (Sason & Verdú, 2016). Without loss of generality, we assume that  $\eta(x) > 1/2$ . Define

$$S(p,q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q} - 2(p-q)^2,$$

where  $p, q \in [0, 1]$ . Let S(p, q) take the partial derivative with respect to q, we get

$$\frac{\partial S(p,q)}{\partial q} = -\frac{p}{q} + \frac{1-p}{1-q} + 4(p-q) = -(p-q)(\frac{1}{q(1-q)} - 4) \le 0, \text{ for } q \le p,$$

where the last inequality follows from that fact that  $q(1-q) \leq 1/4$  for  $q \in [0,1]$  and the equality holds when p=q. Suppose that  $\mathcal{A}_{\theta}$  satisfies  $\epsilon$ -LLDP, which is equivalent to set  $\theta = \exp(\epsilon)/(1 + \exp(\epsilon))$ . Therefore, it holds that

$$D_{KL}\Big(\mathbb{P}_{Y|\boldsymbol{X}=\boldsymbol{x}}\big|\mathbb{P}_{\widetilde{Y}|\boldsymbol{X}=\boldsymbol{x}}\Big) \geq 2(\eta(\boldsymbol{x}) - \widetilde{\eta}(\boldsymbol{x}))^2 = 2(2\eta(\boldsymbol{x}) - 1)^2(1 + \exp(\epsilon))^{-2} \triangleq \mathcal{L}(\epsilon, \boldsymbol{x}).$$

Next, we proceed to prove the upper bound. For any pair of distribution  $\mathbb{P}$  and  $\mathbb{Q}$ , we have

$$D_{KL}(\mathbb{P}|\mathbb{Q}) = D_{KL}(\mathbb{P}|\mathbb{Q}) + 1 - 1 \le \exp(D_{KL}(\mathbb{P}|\mathbb{Q})) - 1.$$

Recall that Y and  $\widetilde{Y}$  take values in  $\{-1,1\}$ , we have

$$D_{KL}\left(\mathbb{P}_{Y|\mathbf{X}=\mathbf{x}}\middle|\mathbb{P}_{\widetilde{Y}|\mathbf{X}=\mathbf{x}}\right) \leq \exp\left(D_{KL}\left(\mathbb{P}_{Y|\mathbf{X}=\mathbf{x}}\middle|\mathbb{P}_{\widetilde{Y}|\mathbf{X}=\mathbf{x}}\right)\right) - 1$$

$$\leq \eta(\mathbf{x})\frac{\eta(\mathbf{x})}{\widetilde{\eta}(\mathbf{x})} + (1 - \eta(\mathbf{x}))\frac{1 - \eta(\mathbf{x})}{1 - \widetilde{\eta}(\mathbf{x})} - 1 = \eta(\mathbf{x})\left(\frac{\eta(\mathbf{x})}{\widetilde{\eta}(\mathbf{x})} - 1\right) + (1 - \eta(\mathbf{x}))\left(\frac{1 - \eta(\mathbf{x})}{1 - \widetilde{\eta}(\mathbf{x})} - 1\right)$$

$$= \eta(\mathbf{x})\left(\frac{\eta(\mathbf{x}) - \widetilde{\eta}(\mathbf{x})}{\widetilde{\eta}(\mathbf{x})}\right) + (1 - \eta(\mathbf{x}))\left(\frac{\widetilde{\eta}(\mathbf{x}) - \eta(\mathbf{x})}{1 - \widetilde{\eta}(\mathbf{x})}\right) = \frac{\left(\eta(\mathbf{x}) - \widetilde{\eta}(\mathbf{x})\right)^{2}}{\widetilde{\eta}(\mathbf{x})(1 - \widetilde{\eta}(\mathbf{x}))}.$$

Note that

$$\widetilde{\eta}(\boldsymbol{x})(1-\widetilde{\eta}(\boldsymbol{x})) = \left(\theta\eta(\boldsymbol{x}) + (1-\theta)(1-\eta(\boldsymbol{x}))\right) \left((1-\theta)\eta(\boldsymbol{x}) + \theta(1-\eta(\boldsymbol{x}))\right)$$

$$= \theta(1-\theta)\eta^2(\boldsymbol{x}) + \eta(\boldsymbol{x})(1-\eta(\boldsymbol{x}))\left(\theta^2 + (1-\theta)^2\right) + \theta(1-\theta)(1-\eta(\boldsymbol{x}))^2$$

$$\geq \max\{\eta(\boldsymbol{x})(1-\eta(\boldsymbol{x})), \theta(1-\theta)\}.$$

It then follows that

$$D_{KL}\left(\mathbb{P}_{Y|X=x}\middle|\mathbb{P}_{\widetilde{Y}|X=x}\right) \le \frac{\left(\eta(x) - \widetilde{\eta}(x)\right)^2}{\max\{\eta(x)(1 - \eta(x)), \theta(1 - \theta)\}}.$$
(4)

Plugging  $\theta = \exp(\epsilon)/(1 + \exp(\epsilon))$  into (4) yields that

$$D_{KL}\Big(\mathbb{P}_{Y|\boldsymbol{X}=\boldsymbol{x}}\big|\mathbb{P}_{\widetilde{Y}|\boldsymbol{X}=\boldsymbol{x}}\Big) \leq \frac{(2\eta(\boldsymbol{x})-1)^2}{\max\{\eta(\boldsymbol{x})(1-\eta(\boldsymbol{x}))(1+\exp(\epsilon))^2,\exp(\epsilon)\}} \triangleq \mathcal{U}(\theta,\boldsymbol{x}).$$

This completes the proof.

**Proof of Lemma** 2: By the definition of  $\widetilde{\eta}(x)$ , we can obtain

$$2\widetilde{\eta}(\boldsymbol{x}) - 1 = (2\theta - 1)(2\eta(\boldsymbol{x}) - 1). \tag{5}$$

Clearly,  $2\eta(\boldsymbol{x}) - 1$  indicates that  $2\widetilde{\eta}(\boldsymbol{x}) - 1$  provided that  $\theta > 1/2$ . By the fact that  $f^*(\boldsymbol{x}) = \text{sign}(\eta(\boldsymbol{x}) - 1/2)$ , it holds that  $f^*(\boldsymbol{x}) = \widetilde{f}^*(\boldsymbol{x})$  for any  $\boldsymbol{x} \in \mathcal{X}$ . Recall that the excess risk in terms of 0-1 loss can be written as

$$R(f) - R(f^*) = \mathbb{E}[|\operatorname{sign}(f(\boldsymbol{X})) \neq \operatorname{sign}(f^*(\boldsymbol{X}))||2\eta(\boldsymbol{X}) - 1|].$$

Combined with (5), it holds that

$$\begin{split} \widetilde{R}(f) - \widetilde{R}(\widetilde{f}^*) = & \mathbb{E}\big[|\operatorname{sign}(f(\boldsymbol{X})) \neq \operatorname{sign}(\widetilde{f}^*(\boldsymbol{X}))||2\widetilde{\eta}(\boldsymbol{X}) - 1|\big] \\ = & (2\theta - 1)\mathbb{E}\big[|\operatorname{sign}(f(\boldsymbol{X})) \neq \operatorname{sign}(\widetilde{f}^*(\boldsymbol{X}))||2\eta(\boldsymbol{X}) - 1|\big] \\ = & (2\theta - 1)\big(R(f) - R(f^*)\big). \end{split}$$

This completes the proof.

**Proof of Lemma 3**: By the assumption that  $||f||_{\infty} \le 1$  and the fact that 0-1 loss is upper bounded by hinge loss, we obtain  $R_{\phi}(f) - R_{\phi}(f^*) \ge R(f) - R(f^*)$ . Since  $\phi$  is set as hinge loss, it holds that

$$R_{\phi}(f) - R_{\phi}(f^*) = \mathbb{E}_{\boldsymbol{X}} \left[ |f(\boldsymbol{X}) - f^*(\boldsymbol{X})| |1 - 2\eta(\boldsymbol{X})| \right]$$
$$\geq t \mathbb{E}_{\boldsymbol{X}} \left[ |f(\boldsymbol{X}) - f^*(\boldsymbol{X})| I(|1 - 2\eta(\boldsymbol{X})| > t) \right].$$

By the fact that  $I(|1-2\eta(X)|>t)=1-I(|1-2\eta(X)|\leq t)$ , it then follows that

$$R_{\phi}(f) - R_{\phi}(f^{*}) \ge t \Big( \mathbb{E}_{\boldsymbol{X}} \big[ |f(\boldsymbol{X}) - f^{*}(\boldsymbol{X})| \big] - 2\mathbb{P} \big( |1 - 2\eta(\boldsymbol{X})| \le t \big) \Big)$$

$$\ge t \Big( \mathbb{E}_{\boldsymbol{X}} \big[ |f(\boldsymbol{X}) - f^{*}(\boldsymbol{X})| \big] - 2ct^{\gamma} \Big)$$

$$= t \mathbb{E}_{\boldsymbol{X}} \big[ |f(\boldsymbol{X}) - f^{*}(\boldsymbol{X})| \big] - 2ct^{\gamma+1}$$

where the last inequality follows from Assumption 1. Choosing t such that  $t\mathbb{E}[|f(\mathbf{X}) - f^*(\mathbf{X})|] = 4ct^{\gamma+1}$ , we get

$$R_{\phi}(f) - R_{\phi}(f^*) \ge (4c)^{-1/\gamma} \Big( \mathbb{E}\big[ |f(\boldsymbol{X}) - f^*(\boldsymbol{X})| \big] \Big)^{\frac{\gamma+1}{\gamma}}.$$

Next, we proceed to establish the relation between  $\widetilde{R}_{\phi}(f) - \widetilde{R}_{\phi}(f^*)$  and  $\mathbb{E}[|f(\boldsymbol{X}) - f^*(\boldsymbol{X})|]$ . For each  $\boldsymbol{x} \in \mathcal{X}$ , the conditional risk is given as

$$\begin{split} &\mathbb{E}_{\widetilde{Y}}\Big(\phi(f(\boldsymbol{X})\widetilde{Y})|\boldsymbol{X}=\boldsymbol{x}\Big) = \widetilde{\eta}(\boldsymbol{x})\phi(f(\boldsymbol{x})) + (1-\widetilde{\eta}(\boldsymbol{x}))\phi(-f(\boldsymbol{x})) \\ &= [\theta\eta(\boldsymbol{x}) + (1-\theta)(1-\eta(\boldsymbol{x}))]\phi(f(\boldsymbol{x})) + [\theta(1-\eta(\boldsymbol{x})) + (1-\theta)\eta(\boldsymbol{x})]\phi(-f(\boldsymbol{x})) \\ &= \Big(\theta\eta(\boldsymbol{x})\phi(f(\boldsymbol{x})) + \theta(1-\eta(\boldsymbol{x}))\phi(-f(\boldsymbol{x}))\Big) + \Big((1-\theta)(1-\eta(\boldsymbol{x}))\phi(f(\boldsymbol{x})) + (1-\theta)\eta(\boldsymbol{x})\phi(-f(\boldsymbol{x}))\Big). \end{split}$$

Taking the expectation with respect to X yields that

$$\widetilde{R}_{\phi}(f) = \mathbb{E}_{(\boldsymbol{X},\widetilde{Y})} \left( \phi(f(\boldsymbol{X})\widetilde{Y}) \right) = \theta R_{\phi}(f) + (1 - \theta)R_{\phi}(-f)$$
(6)

Notice that  $\phi$  is hinge loss, hence

$$R_{\phi}(f) - R_{\phi}(f^*) = R_{\phi}(-f^*) - R_{\phi}(-f),$$

for any f with  $||f||_{\infty} \leq 1$ . It follows that

$$\widetilde{R}_{\phi}(f) - \widetilde{R}_{\phi}(f^{*}) = \theta(R_{\phi}(f) - R_{\phi}(f^{*})) + (1 - \theta)(R_{\phi}(-f) - R_{\phi}(-f^{*}))$$

$$= (2\theta - 1)(R_{\phi}(f) - R_{\phi}(f^{*}))$$

$$\geq (2\theta - 1)(4c)^{-1/\gamma} \Big( \mathbb{E}[|f(\mathbf{X}) - f^{*}(\mathbf{X})|] \Big)^{\frac{\gamma + 1}{\gamma}}.$$

This completes the proof.

**Proof of Theorem 1**: We first denote that  $\widetilde{H}(f, f^*) = \widetilde{R}_{\phi}(f) - \widetilde{R}_{\phi}(f^*)$ . For any  $\delta_n > 0$ , we let  $\mathcal{F}_{\delta_n} = \{f \in \mathcal{F} : \widetilde{H}(f, f^*) \geq \delta_n\}$ , then we have

$$\mathbb{P}\Big(\widetilde{H}(\widetilde{f}_n, f^*) \ge \delta_n\Big)$$

$$\leq \mathbb{P}\Big(\sup_{f \in \mathcal{F}_{\delta_n}} \frac{1}{n} \sum_{i=1}^n \phi(f_{\mathcal{F}}^*(\boldsymbol{x}_i)\widetilde{y}_i) + \lambda_n J(f_{\mathcal{F}}^*) - \frac{1}{n} \sum_{i=1}^n \phi(f(\boldsymbol{x}_i)\widetilde{y}_i) - \lambda_n J(f) \ge 0\Big) \equiv I,$$

where the first inequality follows from the optimality of  $\widetilde{f}$  in minimizing (1).

Define that  $\mathcal{H}_{ij} = \{ f \in \mathcal{F} : 2^{i-1}\delta_n \leq \widetilde{H}(f, f^*) \leq 2^i\delta_n, 2^{j-1}J_0 \leq J(f) \leq 2^jJ_0 \}$  for  $i \geq 1$  and  $j \geq 1$  and  $\mathcal{H}_{i0} = \{ f \in \mathcal{F} : 2^{i-1}\delta_n \leq \widetilde{H}(f, f^*) \leq 2^i\delta_n, J(f) \leq J_0 \}$ . It can be easily verified that  $\mathcal{F}_{\delta_n}$  admits the decomposition as  $\mathcal{F}_{\delta_n} = \bigcup_{i=1}^{\infty} \bigcup_{j=0}^{\infty} \mathcal{H}_{ij}$ . With this, I can be upper bounded as

$$I = \mathbb{P}\left(\sup_{f \in \mathcal{F}_{\delta_n}} \frac{1}{n} \sum_{i=1}^n \phi(f_{\mathcal{F}}^*(\boldsymbol{x}_i) \widetilde{y}_i) + \lambda_n J(f_{\mathcal{F}}^*) - \frac{1}{n} \sum_{i=1}^n \phi(f(\boldsymbol{x}_i) \widetilde{y}_i) - \lambda_n J(f) \ge 0\right)$$

$$\leq \sum_{i=1}^{\infty} \sum_{j=0}^{\infty} \mathbb{P}\left(\sup_{f \in \mathcal{H}_{ij}} \frac{1}{n} \sum_{i=1}^n \phi(f_{\mathcal{F}}^*(\boldsymbol{x}_i) \widetilde{y}_i) + \lambda_n J(f_{\mathcal{F}}^*) - \frac{1}{n} \sum_{i=1}^n \phi(f(\boldsymbol{x}_i) \widetilde{y}_i) - \lambda_n J(f) \ge 0\right)$$

$$\leq \sum_{i=1}^{\infty} \sum_{j=0}^{\infty} \mathbb{P}\left(\sup_{f \in \mathcal{H}_{ij}} \frac{1}{n} \sum_{i=1}^n l_{\phi}(f_{\mathcal{F}}^*, z_i) - \frac{1}{n} \sum_{i=1}^n l_{\phi}(f, z_i) \ge \lambda_n \left(\inf_{f \in \mathcal{H}_{ij}} J(f) - J_0\right) + \inf_{f \in \mathcal{H}_{ij}} \mathbb{E}(\phi(f(\boldsymbol{X}_i) \widetilde{Y}_i)) - \mathbb{E}(\phi(f_{\mathcal{F}}^*(\boldsymbol{X}_i) \widetilde{Y}_i))\right),$$

where  $z_i = (\boldsymbol{x}_i, \widetilde{y}_i)$  and  $l_{\phi}(f, z_i) = \phi(f(\boldsymbol{x}_i)\widetilde{y}_i) - \mathbb{E}(\phi(f(\boldsymbol{X}_i)\widetilde{Y}_i))$ . Here it is important to note that

$$\inf_{f \in \mathcal{H}_{ij}} \mathbb{E}(\phi(f(\boldsymbol{X}_i)\widetilde{Y}_i)) - \mathbb{E}(\phi(f_{\mathcal{F}}^*(\boldsymbol{X}_i)\widetilde{Y}_i)) = \inf_{f \in \mathcal{H}_{ij}} \widetilde{R}_{\phi}(f) - \widetilde{R}_{\phi}(f^*) - \widetilde{R}_{\phi}(f_{\mathcal{F}}^*) + \widetilde{R}_{\phi}(f^*)$$

$$= \inf_{f \in \mathcal{H}_{ij}} \widetilde{H}(f, f^*) + \widetilde{H}(f_{\mathcal{F}}^*, f^*) = \inf_{f \in \mathcal{H}_{ij}} \widetilde{H}(f, f^*) + (2\theta - 1)H(f_{\mathcal{F}}^*, f^*)$$

Let  $V(i,j) = \lambda_n \left( \inf_{f \in \mathcal{H}_{ij}} J(f) - J_0 \right) + \inf_{f \in \mathcal{H}_{ij}} \mathbb{E}(l_{\phi}(f,Z)) - \mathbb{E}(l_{\phi}(f_{\mathcal{F}}^*,Z))$ . By the definition of  $\mathcal{H}_{ij}$ , we get

$$V(i,j) \ge M(i,j) = \lambda_n (2^{j-1} - 1)J_0 + (2^{i-1} - 1/4)\delta_n, \text{ for } i, j \ge 1,$$
(7)

where the inequality follows by assuming that  $(2\theta - 1)s_n \le 1/4\delta_n$ . Next, we suppose that  $\lambda_n J_0 \le 1/4\delta_n$ , we further have

$$M(i,0) \ge (2^{i-1} - 1/2)\delta_n \ge 2^{i-2}\delta_n$$
, for  $i \ge 1$ . (8)

Plugging (7) and (8) into I, it follows that

$$I \leq \sum_{i=1}^{\infty} \sum_{j=0}^{\infty} \mathbb{P} \Big( \sup_{f \in \mathcal{H}_{ij}} \frac{1}{n} \sum_{i=1}^{n} l_{\phi}(f_{\mathcal{F}}^{*}, z_{i}) - \frac{1}{n} \sum_{i=1}^{n} l_{\phi}(f, z_{i}) \geq M(i, j) \Big) = \sum_{i=1}^{\infty} \sum_{j=0}^{\infty} P_{ij}.$$

Therefore, bounding I reduces to bounding  $P_{ij}$ . Let  $Q_n = \frac{1}{n} \sum_{i=1}^n \left( l_{\phi}(f_{\mathcal{F}}^*, z_i) - l_{\phi}(f, z_i) \right)$ , then  $P_{ij}$  can be written as

$$P_{ij} = \mathbb{P}\Big(\sup_{f \in \mathcal{H}_{ij}} Q_n - \mathbb{E}\big[\sup_{f \in \mathcal{H}_{ij}} Q_n\big] \ge M(i,j) - \mathbb{E}\big[\sup_{f \in \mathcal{H}_{ij}} Q_n\big]\Big).$$

Next, we proceed to bound  $P_{ij}$  by the Talagrand's inequality (see Theorem 2.6 in Koltchinskii, 2011). To this end, we first establish the relation between  $\mathbb{E}\big[\sup_{f\in\mathcal{H}_{ij}}Q_n\big]$  and M(i,j). Let  $q(f,z_i)=\phi(f_{\mathcal{F}}^*(\boldsymbol{x}_i)y_i)-\phi(f(\boldsymbol{x}_i)y_i)$  and  $\boldsymbol{z}'=(z_1',\ldots,z_n')$  be a ghost sample.

$$\mathbb{E}\left[\sup_{f\in\mathcal{H}_{ij}}Q_{n}\right] = \mathbb{E}\left[\sup_{f\in\mathcal{H}_{ij}}\frac{1}{n}\sum_{i=1}^{n}q(f,z_{i}) - \mathbb{E}(q(f,Z))\right]$$

$$=\mathbb{E}_{\boldsymbol{z}}\left[\sup_{f\in\mathcal{H}_{ij}}\mathbb{E}_{\boldsymbol{z}'}\left(\frac{1}{n}\sum_{i=1}^{n}q(f,z_{i}) - \frac{1}{n}\sum_{i=1}^{n}q(f,z'_{i})|\boldsymbol{z}\right)\right]$$

$$\leq \mathbb{E}_{\boldsymbol{z},\boldsymbol{z}'}\left[\sup_{f\in\mathcal{H}_{ij}}\left(\frac{1}{n}\sum_{i=1}^{n}q(f,z_{i}) - \frac{1}{n}\sum_{i=1}^{n}q(f,z'_{i})\right)\right]$$

$$=\mathbb{E}_{\boldsymbol{z},\boldsymbol{z}',\boldsymbol{\sigma}}\left[\sup_{f\in\mathcal{H}_{ij}}\left(\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}(q(f,z_{i}) - q(f,z'_{i}))\right)\right]$$

$$\leq 2\mathbb{E}_{\boldsymbol{z}}\mathbb{E}_{\boldsymbol{\sigma}}\left[\sup_{f\in\mathcal{H}_{ij}}\left(\frac{1}{n}\sum_{i=1}^{n}\sigma_{i}q(f,z_{i})\right)\right] = 2\mathcal{R}_{n}(\mathcal{H}_{ij}).$$

By Theorem 3.11 in Koltchinskii (2011), it follows that there exists some constant  $C_2$  such that

$$\mathcal{R}_{n}(\mathcal{H}_{ij}) \leq \frac{C_{2}}{\sqrt{n}} \mathbb{E} \int_{0}^{\sigma_{ij}} \sqrt{\log \mathcal{N}(u, \mathcal{H}_{ij}, L_{2}(P_{n}))} du$$
$$\leq \frac{C_{2}}{\sqrt{n}} \mathbb{E} \int_{0}^{\sigma_{ij}} \sqrt{\mathcal{V}_{1}(\Theta) \log(u^{-1}\mathcal{V}_{2}(\Theta))} du$$

where  $\sigma_{ij} = \sqrt{\sup_{f \in \mathcal{H}_{ij}} \frac{1}{n} \sum_{i=1}^{n} q^2(f, z_i)}$ ,  $||f||_{L_2(P_n)} = \sqrt{\frac{1}{n} \sum_{i=1}^{n} q^2(f, z_i)}$ , and  $\mathcal{N}(\mathcal{H}_{ij}, L_2(P_n), u)$  is the minimal number of  $L_2(P_n)$ -balls of radius u to cover  $\mathcal{H}_{ij}$ .

Notice that  $\int_0^{\sigma_{ij}} \sqrt{\mathcal{V}_1(\Theta) \log(u^{-1}\mathcal{V}_2(\Theta))} du$  is concave function with respect to  $\sigma_{ij}$ , therefore

$$\mathbb{E} \int_{0}^{\sigma_{ij}} \sqrt{\mathcal{V}_{1}(\Theta) \log(u^{-1}\mathcal{V}_{2}(\Theta))} du \leq \int_{0}^{\mathbb{E}\sigma_{ij}} \sqrt{\mathcal{V}_{1}(\Theta) \log(u^{-1}\mathcal{V}_{2}(\Theta))} du$$
$$\leq \int_{0}^{\sqrt{\mathbb{E}\sigma_{ij}^{2}}} \sqrt{\mathcal{V}_{1}(\Theta) \log(u^{-1}\mathcal{V}_{2}(\Theta))} du.$$

By the definition of  $\sigma_{ij}$ , we have  $\mathbb{E}[\sigma_{ij}^2] = \mathbb{E}[\sup_{f \in \mathcal{H}_{ij}} \frac{1}{n} \sum_{i=1}^n q^2(f, z_i)]$ . Using symmetrization and contraction inequalities, we get

$$\mathbb{E}[\sigma_{ij}^2] \le \sup_{f \in \mathcal{H}_{ij}} \mathbb{E}[q^2(f, Z)] + 8\mathcal{R}_n(f, \mathcal{H}_{ij}).$$

Next, we proceed to bound  $\mathbb{E}[q^2(f,Z)]$ .

$$\mathbb{E}[q^{2}(f,Z)] = \mathbb{E}\left[\phi(f_{\mathcal{F}}^{*}(\boldsymbol{X})\widetilde{Y}) - \phi(f(\boldsymbol{X})\widetilde{Y})\right]^{2}$$

$$\leq 2\mathbb{E}\left[\phi(f_{\mathcal{F}}^{*}(\boldsymbol{X})\widetilde{Y}) - \phi(f_{\phi}^{*}(\boldsymbol{X})\widetilde{Y})\right]^{2} + 2\mathbb{E}\left[\phi(f_{\mathcal{F}}^{*}(\boldsymbol{X})\widetilde{Y}) - \phi(f(\boldsymbol{X})\widetilde{Y})\right]^{2}$$

$$\leq 2\mathbb{E}\left[\phi(f_{\mathcal{F}}^{*}(\boldsymbol{X})\widetilde{Y}) - \phi(f^{*}(\boldsymbol{X})\widetilde{Y})\right] + 2\mathbb{E}\left[\left|\phi(f(\boldsymbol{X})\widetilde{Y}) - \phi(f_{\mathcal{F}}^{*}(\boldsymbol{X})\widetilde{Y})\right|\right]$$

$$\leq 2s_{n} + 2\mathbb{E}\left[\left|f(\boldsymbol{X}) - f_{\mathcal{F}}^{*}(\boldsymbol{X})\right|\right] \leq 2^{-1}\delta_{n} + 2\mathbb{E}\left[\left|f(\boldsymbol{X}) - f_{\mathcal{F}}^{*}(\boldsymbol{X})\right|\right], \tag{9}$$

where the second inequality follows from the assumption that  $||f||_{\infty} \leq 1$ . Combining this with Lemma 3 yields that

$$\sup_{f \in \mathcal{H}_{ij}} \mathbb{E}[q^2(f,Z)] \le 2^{-1} \delta_n + 2C_1 \left( (2\theta - 1)^{-1} (\widetilde{R}_{\phi}(f) - \widetilde{R}_{\phi}(f^*)) \right)^{\frac{\gamma}{\gamma+1}} = 4C_1 (2\theta - 1)^{-\frac{\gamma}{\gamma+1}} (2^i \delta_n)^{\frac{\gamma}{\gamma+1}},$$

where  $C_1 = (4c)^{\frac{1}{\gamma+1}}$ . Consequently, we get

$$\mathcal{R}_n(\mathcal{H}_{ij}) \le \frac{C_2}{\sqrt{n}} \int_0^{U_{ij}(f)} \sqrt{\mathcal{V}_1(\Theta) \log(u^{-1}\mathcal{V}_2(\Theta))} du, \tag{10}$$

where  $U_{ij}(f) = \min \left\{ \sqrt{4C_1(2\theta - 1)^{-\frac{\gamma}{\gamma+1}}(2^i\delta_n)^{\frac{\gamma}{\gamma+1}} + 8\mathcal{R}_n(\mathcal{H}_{ij})}, 1 \right\}$  due to the fact that  $|q(f, Z)| \leq 1$ . Then, the right-hand side of (10) can be upper bounded as

$$\frac{C_2}{\sqrt{n}} \int_0^{U_{ij}(f)} \sqrt{\mathcal{V}_1(\Theta)} \log(u^{-1}\mathcal{V}_2(\Theta)) du$$

$$= \frac{C_2\mathcal{V}_2(\Theta)\sqrt{\mathcal{V}_1(\Theta)}}{\sqrt{n}} \int_0^{U_{ij}(f)/\mathcal{V}_2(\Theta)} \sqrt{\log\left(\frac{1}{u}\right)} du$$

$$= \frac{C_2\mathcal{V}_2(\Theta)\sqrt{\mathcal{V}_1(\Theta)}}{\sqrt{n}} \int_{\frac{\mathcal{V}_2(\Theta)}{U_{ij}(f)}}^{+\infty} \frac{1}{u^2} \sqrt{\log\left(u\right)} du$$

$$\leq \frac{2C_2\sqrt{\mathcal{V}_1(\Theta)}U_{ij}(f)}{\sqrt{n}} \sqrt{\log\left(\frac{\mathcal{V}_2(\Theta)}{U_{ij}(f)}\right)}.$$
(11)

Combining (11) with (10) yields that

$$\left(\mathcal{R}_n(\mathcal{H}_{ij})\right)^2 \le \frac{4C_2^2 \mathcal{V}_1(\Theta) \left(4C_1(2\theta - 1)^{-\frac{\gamma}{\gamma+1}} \left(2^i \delta_n\right)^{\frac{\gamma}{\gamma+1}} + 8\mathcal{R}_n(\mathcal{H}_{ij})\right)}{n} \log\left(\frac{\mathcal{V}_2(\Theta)}{U_{ij}(f)}\right). \tag{12}$$

Solving (12) gives  $\mathcal{R}_n(f,\mathcal{H}_{ij}) \lesssim \max\{\mathcal{V}_1(\Theta)n^{-1}, n^{-1/2}\mathcal{V}_1(\Theta)^{1/2}(2\theta-1)^{-\frac{\gamma}{2(\gamma+1)}}\delta_n^{\frac{\gamma}{2(\gamma+1)}}\}$ . Therefore, we get  $\mathcal{R}_n(f,\mathcal{H}_{ij}) \leq 1/4\delta_n$  provided that  $(\mathcal{V}_1(\Theta)/n)^{\frac{\gamma+1}{\gamma+2}}(2\theta-1)^{-\frac{\gamma}{\gamma+2}}\log(n/\mathcal{V}_1(\Theta)) \leq C\delta_n$  for some large constants C. Hence, as n goes to infinity, it follows that

$$\mathcal{R}_n(f, \mathcal{H}_{ij}) \le 1/4\delta_n \le 1/4M(i, j).$$

With this, we get

$$\mathbb{E}\big[\sup_{f\in\mathcal{H}_{ij}}Q_n\big]\leq \mathcal{R}_n(f,\mathcal{H}_{ij})\leq 1/4M(i,j).$$

Therefore,  $P_{ij}$  can be further bounded as

$$P_{ij} \le \mathbb{P}\Big(\sup_{f \in \mathcal{H}_{ij}} Q_n - \mathbb{E}\big[\sup_{f \in \mathcal{H}_{ij}} Q_n\big] \ge 1/2M(i,j)\Big). \tag{13}$$

Applying Talagrand's inequality to the right-hand side, it follows that there exists some positive constants  $C_4$  such that

$$P_{ij} \le C_4 \exp\left(-\frac{nM(i,j)}{2C_4}\log\left(1 + \frac{M(i,j)}{2\mathbb{E}[\sigma_{ij}^2]}\right)\right).$$

As proved above, we can verify that there exists some constant  $C_5$ 

$$\mathbb{E}[\sigma_{ij}^2] \le \sup_{f \in \mathcal{H}_{ij}} \mathbb{E}[q^2(f,Z)] + 8\mathcal{R}_n(f,\mathcal{H}_{ij}) \le C_5\Big((2\theta - 1)^{-\frac{\gamma}{\gamma+1}}(2^i\delta_n)^{\frac{\gamma}{\gamma+1}} + M(i,j)\Big).$$

Notice that  $\frac{M(i,j)}{2\mathbb{E}[\sigma_{ij}^2]}$  converges to 0 as n increases, therefore there exists some constant  $0 < C_7 < 1$  such that  $\log(1+x) \ge C_7 x$  for  $x \in [0,1/(2C_5)]$ . It then follows that

$$P_{ij} \le C_4 \exp\Big(-\frac{C_7 n M^2(i,j)}{4C_4 C_5 \Big((2\theta-1)^{-\frac{\gamma}{\gamma+1}} (M(i,j))^{\frac{\gamma}{\gamma+1}} + M(i,j)\Big)}\Big).$$

Since  $M(i,j) \ll (2\theta-1)^{-\frac{\gamma}{\gamma+1}} (M(i,j))^{\frac{\gamma}{\gamma+1}}$  when  $\theta-1/2=o(1)$  and M(i,j)=o(1), we have

$$P_{ij} \le C_4 \exp\left(-C_8 n(2\theta - 1)^{\frac{\gamma}{\gamma+1}} M^{\frac{\gamma+2}{\gamma+1}}(i,j)\right),$$

where  $C_8 = C_7/(4C_4C_5)$ . Therefore, we have

$$\sum_{i=1}^{n} \sum_{j=0}^{n} P_{ij} \leq \sum_{i=1}^{n} \sum_{j=1}^{n} C_4 \exp\left(-\frac{C_8 n M^{\frac{\gamma+2}{\gamma+1}}(i,j)}{(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\right) + \sum_{i=1}^{n} C_4 \exp\left(-\frac{C_8 n M^{\frac{\gamma+2}{\gamma+1}}(i,0)}{(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\right) \equiv I_1 + I_2.$$

Note that for any  $i, j \geq 1$ ,

$$M^{\frac{\gamma+2}{\gamma+1}}(i,j) \ge \left( (2^{i-1} - 1/4)\delta_n + \lambda_n (2^{j-1} - 1)J_0 \right)^{\frac{\gamma+2}{\gamma+1}}$$

$$\ge (2^{i-1} - 1/4)\delta_n^{\frac{\gamma+2}{\gamma+1}} + (2^{j-1} - 1)(\lambda_n J_0)^{\frac{\gamma+2}{\gamma+1}}$$

$$\ge i/2\delta_n^{\frac{\gamma+2}{\gamma+1}} + (j-1)(\lambda_n J_0)^{\frac{\gamma+2}{\gamma+1}}.$$

Hence  $I_1$  is upper bounded as

$$I_{1} \leq \sum_{i=1}^{n} \sum_{j=1}^{n} C_{4} \exp\left(-C_{8} n \frac{i/2 \delta_{n}^{\frac{\gamma+2}{\gamma+1}} + (j-1)(\lambda_{n} J_{0})^{\frac{\gamma+2}{\gamma+1}}}{(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\right)$$

$$= C_{4} \frac{\exp\left(\frac{-C_{8} n \delta_{n}^{\frac{\gamma+2}{\gamma+1}}}{2(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\right)}{1 - \exp\left(\frac{-C_{8} n \delta_{n}^{\frac{\gamma+2}{\gamma+1}}}{2(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\right)} \cdot \frac{1}{1 - \exp\left(\frac{-C_{8} n(\lambda_{n} J_{0})^{\frac{\gamma+2}{\gamma+1}}}{2(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\right)}$$

$$\leq 4C_{4} \exp\left(\frac{-C_{8} n \delta_{n}^{\frac{\gamma+2}{\gamma+1}}}{2(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\right),$$

where the last inequality holds when  $\max\{\exp\left(\frac{-C_8n\delta_n^{\frac{\gamma+2}{\gamma+1}}}{2(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\right), \exp\left(\frac{-C_8n\delta_n^{2-1/\gamma}}{2(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\right)\} \le 1/2$ , which holds true when n goes to infinity. Similarly, for  $I_2$ , we get

$$I_2 \le \sum_{i=1}^n C_4 \exp\Big(-\frac{i/4C_8 n \delta_n^{\frac{\gamma+2}{\gamma+1}}}{(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\Big) \le 2C_4 \exp\Big(-\frac{C_8 n \delta_n^{\frac{\gamma+2}{\gamma+1}}}{4(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\Big).$$

Combining  $I_1$  and  $I_2$ , it follows that

$$I \le 6C_4 \exp\left(-\frac{C_8 n \delta_n^{\frac{\gamma+2}{\gamma+1}}}{4(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\right) = T_1 \exp\left(-T_2 \frac{n \delta_n^{\frac{\gamma+2}{\gamma+1}}}{(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\right),$$

where  $T_1 = 6C_4$  and  $T_2 = C_8/4$ . Therefore, we conclude that

$$\mathbb{P}\Big(\widetilde{H}(\widetilde{f}_n, f^*) \ge \delta_n\Big) \le T_1 \exp\Big(-T_2 \frac{n\delta_n^{\frac{\gamma+2}{\gamma+1}}}{(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\Big).$$

With a choice of  $\delta_n$  such that  $\delta_n \geq C_9((2\theta-1)^{-\frac{\gamma}{\gamma+1}}|\Theta|n^{-1}\log(n/|\Theta|))^{\frac{\gamma+1}{\gamma+2}}$  for some positive constants  $C_9 > 0$ , we have  $\widetilde{H}(\widetilde{f}_n, f^*) = o_p(1)$ , which implies that  $\mathbb{E}(\widetilde{H}(\widetilde{f}_n, f^*)) = O_p(\delta_n)$ . Notice that  $T_1$  and  $T_2$  are both independent of  $\pi$ , therefore we further have

$$\sup_{\pi \in \mathcal{P}_{\gamma}} \mathbb{P}\Big(\widetilde{H}(\widetilde{f}_n, f^*) \ge \delta_n\Big) \le T_1 \exp\Big(-T_2 \frac{n\delta_n^{\frac{\gamma+2}{\gamma+1}}}{(2\theta-1)^{-\frac{\gamma}{\gamma+1}}}\Big).$$

By the relation between excess risk and excess  $\phi$ -risk  $\mathbb{E}_{\widetilde{D}}(\widetilde{D}(\widetilde{f}_n, f^*)) \leq \mathbb{E}_{\widetilde{D}}(\widetilde{H}(\widetilde{f}_n, f^*))$  (Bartlett et al., 2006), we further have

$$\sup_{\pi \in \mathcal{P}_{\gamma}} \mathbb{E}_{\widetilde{\mathcal{D}}} \Big( \widetilde{D}(\widetilde{f}_n, f^*) \Big) = O(\delta_n).$$

Combining this with the Lemme 2 that  $D(\widetilde{f}_n, f^*) = (2\theta - 1)^{-1} \widetilde{D}(\widetilde{f}_n, f^*)$ , we get

$$\sup_{\pi \in \mathcal{P}_{\gamma}} \mathbb{E}_{\widetilde{\mathcal{D}}} \Big( D(\widetilde{f}_n, f^*) \Big) = O(\delta_n (2\theta - 1)^{-1}).$$

**Proof of lower bound**: We first define the minimax excess risk under the noisy distribution and true distribution as

$$W_n = \inf_{\widetilde{f}_n} \sup_{\pi \in \mathcal{P}_{\gamma}} \mathbb{E}_{\widetilde{\mathcal{D}}} [R(\widetilde{f}_n) - R(f^*)],$$

which admits the decomposition as

$$W_n = \inf_{\widetilde{f}_n} \sup_{\pi \in \mathcal{P}_n} \Big\{ \mathbb{E}_{\widetilde{\mathcal{D}}} \big[ R(\widetilde{f}_n) - \inf_{f \in \mathcal{F}} R(f) \big] + \inf_{f \in \mathcal{F}} R(f) - R(f^*) \Big\},$$

where the second term on the right-hand side denotes the approximation error under the 0-1 risk. In what follows, we proceed to consider a sub-family of  $\mathcal{P}_{\gamma}$  such that  $\inf_{f \in \mathcal{F}} R(f) = R(f^*)$ . For example, let  $\mathcal{P}' \subset \mathcal{P}_{\gamma}$  be such a sub-family, then we have

$$\begin{split} W_n &\geq \inf_{\widetilde{f}_n} \sup_{\pi \in \mathcal{P}'} \left\{ \mathbb{E}_{\widetilde{\mathcal{D}}} \big[ R(\widetilde{f}_n) - \inf_{f \in \mathcal{F}} R(f) \big] \right\} + \sup_{\pi \in \mathcal{P}_{\gamma}} \left( \inf_{f \in \mathcal{F}} R(f) - R(f^*) \right) \\ &= \inf_{\widetilde{f}_n} \sup_{\pi \in \mathcal{P}'} \left\{ \mathbb{E}_{\widetilde{\mathcal{D}}} \big[ R(\widetilde{f}_n) - R(f^*) \big] \right\} + \sup_{\pi \in \mathcal{P}_{\gamma}} \left( \inf_{f \in \mathcal{F}} R(f) - R(f^*) \right) \\ &= (2\theta - 1)^{-1} \inf_{\widetilde{f}_n} \sup_{\pi \in \widetilde{\mathcal{P}}'(\theta)} \left\{ \mathbb{E}_{\widetilde{\mathcal{D}}} \big[ \widetilde{R}(\widetilde{f}_n) - \widetilde{R}(f^*) \big] \right\} + \sup_{\pi \in \mathcal{P}_{\gamma}} \left( \inf_{f \in \mathcal{F}} R(f) - R(f^*) \right), \end{split}$$

where the last inequality follows from Lemma 2 and  $\widetilde{\mathcal{P}}'(\theta)$  is a set of probability measures  $\widetilde{\pi}$  on  $(X, \widetilde{Y})$  such that any probability distribution  $\widetilde{\pi} \in \widetilde{\mathcal{P}}'(\theta)$  is associated with an  $\pi \in \mathcal{P}'$  with  $\widetilde{\pi}$  and  $\pi$  having the same marginal distribution on X and  $\widetilde{Y} = \mathcal{A}_{\theta}(Y)$ .

Denote that  $\widetilde{W_n} = \inf_{\widetilde{f_n}} \sup_{\widetilde{\pi} \in \widetilde{\mathcal{P}}'(\theta)} \mathbb{E}_{\widetilde{\mathcal{D}}} \big[ \widetilde{R}(\widetilde{f_n}) - \widetilde{R}(f^*) \big]$  and  $\widetilde{\mathcal{P}}_{\gamma}(\theta) = \{ \widetilde{\pi}(\boldsymbol{X}, \widetilde{Y}) : \widetilde{Y} = \mathcal{A}_{\theta}(Y), \pi \in \mathcal{P}_{\gamma} \}$ . Then we proceed to construct  $\widetilde{\mathcal{P}}'(\theta) \subset \widetilde{\mathcal{P}}_{\gamma}(\theta)$ . First, following from Lemma 3, we have for any  $\widetilde{\pi} \in \widetilde{\mathcal{P}}_{\gamma}(\theta)$ 

$$\mathbb{P}\Big(|2\widetilde{\eta}(\boldsymbol{X}) - 1| \le t\Big) \le c(2\theta - 1)^{-\gamma}t^{\gamma},$$

where the probability is measured under  $\tilde{\pi}$ .

By Assumption 2, we know that the VC dimension of  $\mathcal{G}_{\mathcal{F}}$  is lower bounded by  $C_V$  such that  $VC(\mathcal{G}_{\mathcal{F}}) \geq C_V \Lambda$ . Therefore, for any  $N \leq C_V \Lambda$ , there exist N distinct points  $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N$  such that  $\{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N\}$  is shattered by  $\mathcal{G}_{\mathcal{F}}$ . Therefore, we consider distribution supported on  $\{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_N\}$ . Let  $w \in (0,1)$  be a number satisfying  $(N-1)w \leq 1$ . Let Q be the probability measure on  $\mathcal{X}$  such that

$$Q(x_i) = \begin{cases} w, & i = 1, \dots, N - 1, \\ 1 - (N - 1)w, & i = N. \end{cases}$$

In what follows, we consider the hypercube  $C = \{-1, 1\}^{N-1}$ . For any  $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_{N-1}) \in C$ , we define the regression function as

$$\widetilde{\eta}_{\sigma}(\boldsymbol{x}_i) = \begin{cases} \frac{1+\sigma_j h}{2}, & i = 1, \dots, N-1, \\ 1, & i = N. \end{cases}$$

Next, we let  $\tilde{\pi}_{\sigma}$  denote the associated probability measure on  $\mathcal{X} \times \{-1, 1\}$  with Q being marginal distribution on  $\mathcal{X}$  and  $\tilde{\eta}_{\sigma}(\boldsymbol{x})$  being the regression function. With this, we obtain

$$\mathbb{P}(|2\widetilde{\eta}_{\sigma}(\boldsymbol{X}) - 1| \le t) = (N - 1)wI(h \le t).$$

By assuming  $(N-1)w \leq c(2\theta-1)^{-\gamma}h^{\gamma}$ , we have  $\mathbb{P}(|2\eta_{\sigma}(X)-1|\leq t)=c(2\theta-1)^{-\gamma}t^{\gamma}$  for any  $t\in[0,1)$ , which indicates that  $\widetilde{\pi}_{\sigma}\in\widetilde{\mathcal{P}}_{\gamma}(\theta)$ . By setting  $\widetilde{\mathcal{P}}'(\theta)=\{\widetilde{\pi}_{\sigma}:\sigma\in\mathcal{C}\}$ , we have

$$\widetilde{W_n} = \inf_{\widetilde{f}_n} \sup_{\widetilde{\pi}_{\sigma} \in \widetilde{\mathcal{P}}'(\theta)} \mathbb{E} \big[ \widetilde{R}(\widetilde{f}_n) - \widetilde{R}(f_{\sigma}^*) \big],$$

where  $f_{\boldsymbol{\sigma}}^*(\boldsymbol{x}_i) = \sigma_i$  for  $i = 1, \dots, N-1$ .

Under the probability measure  $\widetilde{\pi}_{\sigma}$ 

$$\begin{split} \widetilde{R}(\widetilde{f}_n) - \widetilde{R}(f_{\boldsymbol{\sigma}}^*) &= \mathbb{E}_{\widetilde{\pi}_{\boldsymbol{\sigma}}} \big[ |\operatorname{sign}(\widetilde{f}_n(\boldsymbol{X})) \neq f_{\boldsymbol{\sigma}}^*(\boldsymbol{X}) | |1 - 2\widetilde{\eta}(\boldsymbol{X})| \big] \\ &= \frac{1}{2} \mathbb{E}_{\widetilde{\pi}_{\boldsymbol{\sigma}}} \big[ |\operatorname{sign}(\widetilde{f}_n(\boldsymbol{X})) - f_{\boldsymbol{\sigma}}^*(\boldsymbol{X}) | |1 - 2\widetilde{\eta}(\boldsymbol{X})| \big] \\ &\geq \frac{t}{2} \mathbb{E}_{\widetilde{\pi}_{\boldsymbol{\sigma}}} \Big[ |\operatorname{sign}(\widetilde{f}_n(\boldsymbol{X})) - f_{\boldsymbol{\sigma}}^*(\boldsymbol{X}) | \Big( 1 - \mathbb{P} \Big( |1 - 2\widetilde{\eta}(\boldsymbol{X})| \leq t \Big) \Big) \Big] \\ &\geq (2\theta - 1) (4^{1-\gamma}c)^{-1/\gamma} \Big( \mathbb{E}_{\widetilde{\pi}_{\boldsymbol{\sigma}}} \big[ |\operatorname{sign}(\widetilde{f}_n(\boldsymbol{X})) - f_{\boldsymbol{\sigma}}^*(\boldsymbol{X}) | \big] \Big)^{\frac{\gamma+1}{\gamma}} \\ &\geq (2\theta - 1) (4^{1-\gamma}c)^{-1/\gamma} \Big( w \sum_{i=1}^{N-1} |\operatorname{sign}(\widetilde{f}_n(\boldsymbol{x}_i)) - f_{\boldsymbol{\sigma}}^*(\boldsymbol{x}_i) | \big] \Big)^{\frac{\gamma+1}{\gamma}}. \end{split}$$

Under the distribution  $\widetilde{\pi}_{\sigma}$ , we have  $f^*(\boldsymbol{x}_i) = \sigma_i$  for  $i = 1, \dots, N-1$ , which then implies that

$$\widetilde{R}(\widetilde{f}_n) - \widetilde{R}(f_{\boldsymbol{\sigma}}^*) \ge (2\theta - 1)(4^{1-\gamma}c)^{-1/\gamma} w^{\frac{\gamma+1}{\gamma}} \Big( \sum_{i=1}^{N-1} |\operatorname{sign}(\widetilde{f}_n(\boldsymbol{x}_i)) - \sigma_i| \Big)^{\frac{\gamma+1}{\gamma}}.$$

Taking the expectation of both sides with respect to  $\widetilde{\mathcal{D}}$  yields that

$$\mathbb{E}_{\widetilde{\mathcal{D}}}\left(\widetilde{R}(\widetilde{f}_n) - \widetilde{R}(f_{\boldsymbol{\sigma}}^*)\right) \ge (2\theta - 1)(4^{1-\gamma}c)^{-1/\gamma}w^{\frac{\gamma+1}{\gamma}}\mathbb{E}_{\widetilde{\mathcal{D}}}\left[\left(\sum_{i=1}^{N-1}|\operatorname{sign}(\widetilde{f}_n(\boldsymbol{x}_i)) - \sigma_i|\right)^{\frac{\gamma+1}{\gamma}}\right]$$

$$\ge (2\theta - 1)(4^{1-\gamma}c)^{-1/\gamma}w^{\frac{\gamma+1}{\gamma}}\mathbb{E}_{\widetilde{\mathcal{D}}}^{\frac{\gamma+1}{\gamma}}\left[\left(\sum_{i=1}^{N-1}|\operatorname{sign}(\widetilde{f}_n(\boldsymbol{x}_i)) - \sigma_i|\right)\right],$$

where the second inequality follows from Jensen's inequality.

For ease of notation, we let  $X(w, \gamma, \theta) = (2\theta - 1)(4^{1-\gamma}c)^{-1/\gamma}w^{\frac{\gamma+1}{\gamma}}$ 

$$\sup_{\widetilde{\pi}_{\sigma} \in \widetilde{\mathcal{P}}'(\theta)} \mathbb{E}_{\widetilde{\mathcal{D}}} \left( \widetilde{R}(\widetilde{f}_{n}) - \widetilde{R}(f_{\sigma}^{*}) \right) \geq \sup_{\widetilde{\pi}_{\sigma} \in \widetilde{\mathcal{P}}'(\theta)} X(w, \gamma, \theta) \mathbb{E}_{\widetilde{\mathcal{D}}}^{\frac{\gamma+1}{\gamma}} \left[ \left( \sum_{i=1}^{N-1} |\operatorname{sign}(\widetilde{f}_{n}(\boldsymbol{x}_{i})) - \sigma_{i}| \right) \right],$$

$$\geq \frac{1}{2^{N}} \sum_{\sigma \in \mathcal{C}} X(w, \gamma, \theta) \mathbb{E}_{\widetilde{\mathcal{D}}}^{\frac{\gamma+1}{\gamma}} \left[ \left( \sum_{i=1}^{N-1} |\operatorname{sign}(\widetilde{f}_{n}(\boldsymbol{x}_{i})) - \sigma_{i}| \right) \right]$$

$$= \frac{1}{2^{N-1}} \sum_{\sigma \in \mathcal{C}} X(w, \gamma, \theta) \mathbb{E}_{\widetilde{\mathcal{D}}}^{\frac{\gamma+1}{\gamma}} \left[ \sum_{i=1}^{N-1} I\left(\operatorname{sign}(\widetilde{f}_{n}(\boldsymbol{x}_{i})) \neq \sigma_{i}\right) \right]$$

$$= X(w, \gamma, \theta) \left( \frac{1}{2^{N-1}} \sum_{\sigma \in \mathcal{C}} \sum_{i=1}^{N-1} \mathbb{E}_{\widetilde{\mathcal{D}}} \left[ I\left(\operatorname{sign}(\widetilde{f}_{n}(\boldsymbol{x}_{i})) \neq \sigma_{i}\right) \right] \right)^{\frac{\gamma+1}{\gamma}}.$$

$$= X(w, \gamma, \theta) \left( \sum_{i=1}^{N-1} \frac{1}{2^{N-1}} \sum_{\sigma \in \mathcal{C}} \mathbb{P}_{\widetilde{\pi}_{\sigma}} \left(\operatorname{sign}(\widetilde{f}_{n}(\boldsymbol{x}_{i})) \neq \sigma_{i}\right) \right)^{\frac{\gamma+1}{\gamma}}.$$

For each i, we observe that

$$\frac{1}{2^{N-1}} \sum_{\boldsymbol{\sigma} \in \mathcal{C}} \mathbb{P}_{\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}}} \left( \operatorname{sign}(\widetilde{f}_{n}(\boldsymbol{x}_{i})) \neq \sigma_{i} \right) \\
= \frac{1}{2^{N-1}} \sum_{\boldsymbol{\sigma}: \sigma_{i} = 1} \mathbb{P}_{\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}}} \left( \operatorname{sign}(\widetilde{f}_{n}(\boldsymbol{x}_{i})) \neq \sigma_{i} \right) + \frac{1}{2^{N-1}} \sum_{\boldsymbol{\sigma}: \sigma_{i} = -1} \mathbb{P}_{\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}}} \left( \operatorname{sign}(\widetilde{f}_{n}(\boldsymbol{x}_{i})) \neq \sigma_{i} \right) \\
= 2\mathbb{P}_{+i} \left( \operatorname{sign}(\widetilde{f}_{n}(\boldsymbol{x}_{i})) \neq 1 \right) + 2\mathbb{P}_{-i} \left( \operatorname{sign}(\widetilde{f}_{n}(\boldsymbol{x}_{i})) \neq -1 \right) \geq 2 - 2\operatorname{TV}(\mathbb{P}_{+i}^{\otimes n}, \mathbb{P}_{-i}^{\otimes n}).$$

where  $\mathbb{P}_{+i} = \frac{1}{2^{N-1}} \sum_{\boldsymbol{\sigma}: \sigma_i = 1} \mathbb{P}_{\widetilde{\pi}_{\boldsymbol{\sigma}}}, \ \mathbb{P}_{-i} = \frac{1}{2^{N-1}} \sum_{\boldsymbol{\sigma}: \sigma_i = -1} \mathbb{P}_{\widetilde{\pi}_{\boldsymbol{\sigma}}}, \ \text{and} \ \mathrm{TV}(\mathbb{P}^n_{+i}, \mathbb{P}^n_{-i}) \ \text{denotes the total variation}$  between  $\mathbb{P}^{\otimes n}_{+i}$  and  $\mathbb{P}^{\otimes n}_{-i}$ . Notice that for any probability measures P and Q, we have

$$\begin{aligned} \text{TV}(P,Q) &= \frac{1}{2} \int |p(\boldsymbol{x}) - q(\boldsymbol{x})| d\boldsymbol{x} = \frac{1}{2} \int |\sqrt{p(\boldsymbol{x})} - \sqrt{q(\boldsymbol{x})}| |\sqrt{p(\boldsymbol{x})} + \sqrt{q(\boldsymbol{x})}| d\boldsymbol{x} \\ &\leq \frac{1}{2} \Big( \int \left(\sqrt{p(\boldsymbol{x})} - \sqrt{q(\boldsymbol{x})}\right)^2 d\boldsymbol{x} \Big)^{1/2} \Big( \int \left(\sqrt{p(\boldsymbol{x})} + \sqrt{q(\boldsymbol{x})}\right)^2 d\boldsymbol{x} \Big)^{1/2} \\ &\leq \sqrt{H^2(P,Q)}, \end{aligned}$$

where  $H^2(P,Q)$  denotes the Hellienger distance between P and Q. Let  $\sigma$  and  $\sigma'$  be two indexes such that  $\sigma_l = \sigma'_l$  for  $l \neq i$  and  $\sigma_i = -\sigma'_i = 1$ , then we have

$$\mathrm{TV}(\mathbb{P}_{+i}^{\otimes n}, \mathbb{P}_{-i}^{\otimes n}) \leq \sum_{\sigma, \sigma'} \frac{1}{2^{N-2}} \mathrm{TV}\big(\mathbb{P}_{\widetilde{\pi}_{\sigma}}^{\otimes n}, \mathbb{P}_{\widetilde{\pi}_{\sigma'}}^{\otimes n}\big) \leq \max_{\sigma, \sigma'} \mathrm{TV}\big(\mathbb{P}_{\widetilde{\pi}_{\sigma}}^{\otimes n}, \mathbb{P}_{\widetilde{\pi}_{\sigma'}}^{\otimes n}\big) \leq \max_{\sigma, \sigma'} H(\mathbb{P}_{\widetilde{\pi}_{\sigma}}^{\otimes n}, \mathbb{P}_{\widetilde{\pi}_{\sigma'}}^{\otimes n})$$

By the definition of  $\widetilde{\pi}_{\sigma}$  and  $\widetilde{\pi}_{\sigma'}$ , we get

$$H^{2}(\mathbb{P}_{\widetilde{\pi}_{\sigma}}, \mathbb{P}_{\widetilde{\pi}_{\sigma'}}) = w \sum_{i=1}^{N-1} \left( \sqrt{\eta_{\sigma}(\boldsymbol{x}_{i})} - \sqrt{\eta_{\sigma'}(\boldsymbol{x}_{i})} \right)^{2} + \left( \sqrt{1 - \eta_{\sigma}(\boldsymbol{x}_{i})} - \sqrt{1 - \eta_{\sigma'}(\boldsymbol{x}_{i})} \right)^{2}$$
$$= 2w(1 - \sqrt{1 - h^{2}}) \leq 2wh^{2},$$

for any  $h \in [0,1]$ . By the fact that  $H^2(P^{\otimes n}, Q^{\otimes n}) = 2 - 2(1 - 2^{-1}H^2(P,Q))^n$ , we have

$$\begin{split} H^2(\mathbb{P}^{\otimes n}_{\widetilde{\pi}_{\pmb{\sigma}}}, \mathbb{P}^{\otimes n}_{\widetilde{\pi}_{\pmb{\sigma}'}}) = & 2 - 2 \big[ 1 - w(1 - \sqrt{1 - h^2}) \big]^n \leq 2 - 2 \big[ 1 - nw(1 - \sqrt{1 - h^2}) \big] \\ = & 2nw(1 - \sqrt{1 - h^2}) \leq 1/16, \end{split}$$

where the last inequality holds by choosing w and h such that  $wh^2 \leq n^{-1}/32$ , from which it follows that  $TV(\mathbb{P}_{+i}^n, \mathbb{P}_{-i}^n) \leq 1/4$ . Notice that w and h are chosen to satisfy  $(N-1)w \leq c(2\theta-1)^{-\gamma}h^{\gamma}$  and  $wh^2 \leq n^{-1}/32$ . For simplicity, we obtain the following solution by considering the case equalities hold.

$$h = \frac{(N-1)^{\frac{1}{\gamma+2}} (2\theta-1)^{\frac{\gamma}{\gamma+2}}}{(32cn)^{\frac{1}{\gamma+2}}} \quad \text{and} \quad w = \frac{c^{\frac{2}{\gamma+2}} (N-1)^{-\frac{2}{\gamma+2}} (2\theta-1)^{-\frac{2\gamma}{\gamma+2}}}{(32n)^{\frac{\gamma}{\gamma+2}}}.$$

Therefore, we can conclude that

$$\begin{split} \sup_{\widetilde{\pi}_{\sigma} \in \widetilde{\mathcal{P}}'(\theta)} \mathbb{E}_{\widetilde{\mathcal{D}}} \Big( \widetilde{R}(\widetilde{f}_n) - \widetilde{R}(f_{\sigma}^*) \Big) \geq & (2\theta - 1)(4^{1-\gamma})^{-1/\gamma} c^{-\frac{1}{\gamma+2}} \frac{(N-1)^{-\frac{2\gamma+2}{\gamma(\gamma+2)}} (2\theta - 1)^{-\frac{2\gamma+2}{\gamma+2}}}{(32n)^{\frac{\gamma+1}{\gamma+2}}} \Big( \frac{3(N-1)}{2} \Big)^{\frac{\gamma+1}{\gamma}} \\ = & (2\theta - 1)(4^{1-\gamma})^{-1/\gamma} c^{-\frac{1}{\gamma+2}} \frac{(N-1)^{\frac{\gamma+1}{\gamma+2}} (2\theta - 1)^{-\frac{2\gamma+2}{\gamma+2}}}{(32n)^{\frac{\gamma+1}{\gamma+2}}} \Big( \frac{3}{2} \Big)^{\frac{\gamma+1}{\gamma}} \\ = & (2\theta - 1)(4^{1-\gamma})^{-1/\gamma} c^{-\frac{1}{\gamma+2}} \Big( \frac{N-1}{32n(2\theta - 1)^2} \Big)^{\frac{\gamma+1}{\gamma+2}} \Big( \frac{3}{2} \Big)^{\frac{\gamma+1}{\gamma}} \,. \end{split}$$

Choosing  $N = C_V \Lambda$  yields that

$$\inf_{\widetilde{f}_n} \sup_{\widetilde{\pi}_{\sigma} \in \widetilde{\mathcal{P}}'(\theta)} \mathbb{E}_{\widetilde{\mathcal{D}}} \left( \widetilde{R}(\widetilde{f}_n) - \widetilde{R}(f^*) \right) \ge (2\theta - 1) A_2 \left( \frac{\Lambda}{n(2\theta - 1)^2} \right)^{\frac{\gamma + 1}{\gamma + 2}},$$

where  $A_2 = C_V^{\frac{\gamma+1}{\gamma+2}} (4^{1-\gamma})^{-1/\gamma} c^{-\frac{1}{\gamma+2}} (2^{-1})^{\frac{6\gamma+6}{\gamma+2}} (3/2)^{\frac{\gamma+1}{\gamma}}$ .

Following from Lemma 2, it holds that

$$\inf_{\widetilde{f}_n} \sup_{\pi \in \mathcal{P}_{\gamma}} \mathbb{E}_{\widetilde{\mathcal{D}}} \left( R(\widetilde{f}_n) - R(f^*) \right) \ge A_2 \left( \frac{\Lambda}{n(2\theta - 1)^2} \right)^{\frac{\gamma + 1}{\gamma + 2}} + \sup_{\pi \in \mathcal{P}_{\gamma}} \left( \inf_{f \in \mathcal{F}} R(f) - R(f^*) \right).$$

This completes the proof.

Corollary 1. Suppose that  $\mathcal{F}$  is chosen such that  $s_n = 0$  and  $\epsilon$  is adaptive to n such that  $\epsilon = o(1)$ . There exists some constants  $A_3 > 0$  such that  $\sup_{\pi \in \mathcal{P}_{\alpha}} \mathbb{E}_{\widetilde{\Omega}}[R(\widetilde{f}_n) - R(f^*)] \geq A_3$  provided that  $\epsilon \lesssim n^{-1/2}$ .

**Proof of Corollary 1**: Without loss of generality, we assume that  $\epsilon < 1$  since  $\epsilon = o(1)$ . Then, by the definition of  $\kappa_{\epsilon}$ , there exists some constants C > 0 such that

$$\frac{\mathcal{V}_1(\Theta)}{n\kappa_{\epsilon}^2} = \frac{\mathcal{V}_1(\Theta)(\exp(\epsilon) + 1)^2}{n(\exp(\epsilon) - 1)^2} \ge \frac{\mathcal{V}_1(\Theta)}{ne^2\epsilon^2},\tag{14}$$

where the last inequality follows from the fact that  $e^x-1>ex$  for any  $x\in(0,1]$ . Further, since  $\epsilon=o\left(1/\sqrt{n}\right)$ , we have  $n\epsilon^2=o(1)$ . Therefore, it follows that  $\mathcal{V}_1(\Theta)n^{-1}\kappa_\epsilon^{-2}\gtrsim\mathcal{V}_1(\Theta)e^{-2}$ , and this completes the proof.  $\square$ 

**Proof of Theorem 2**: Let  $\widetilde{Z} = (\widetilde{Y} + 1)/2$ . Then,  $\mathbb{P}(\widetilde{Z} = 0|X) = \mathbb{P}(\widetilde{Y}|X) = 1 - \eta(X)$ . Next, we intend to establish the connection between the excess risk of  $\widetilde{f}_{nn}$  and  $\|\widetilde{f}_{nn} - \eta\|_{L^2(\mathbb{P}_X)}^2$ . Here the proof is mainly based

on Lemma 5.2 in Audibert & Tsybakov (2007).

$$\widetilde{R}(\widetilde{s}_{nn}) - \widetilde{R}(f^{*}) = \mathbb{E}\left[\left|2\widetilde{\eta}(\boldsymbol{X}) - 1\right| \cdot I\left(\widetilde{s}_{nn}(\boldsymbol{X}) \neq f^{*}(\boldsymbol{X})\right)\right]$$

$$\leq 2\mathbb{E}\left[\left|\widetilde{\eta}(\boldsymbol{X}) - 1/2\right| \cdot I\left(\widetilde{s}_{nn}(\boldsymbol{X}) \neq f^{*}(\boldsymbol{X})\right) \cdot I\left(\left|\widetilde{\eta}(\boldsymbol{X}) - 1/2\right| \leq t\right)\right]$$

$$+ 2\mathbb{E}\left[\left|\widetilde{\eta}(\boldsymbol{X}) - 1/2\right| \cdot I\left(\widetilde{s}_{nn}(\boldsymbol{X}) \neq f^{*}(\boldsymbol{X})\right) \cdot I\left(\left|\widetilde{\eta}(\boldsymbol{X}) - 1/2\right| > t\right)\right]$$

$$\leq 2\mathbb{E}\left[\left|\widetilde{\eta}(\boldsymbol{X}) - \widetilde{f}_{nn}(\boldsymbol{X})\right| \cdot I\left(\left|\widetilde{\eta}(\boldsymbol{X}) - 1/2\right| \leq t\right)\right]$$

$$+ 2\mathbb{E}\left[\left|\widetilde{\eta}(\boldsymbol{X}) - \widetilde{f}_{nn}(\boldsymbol{X})\right| \cdot I\left(\left|\widetilde{\eta}(\boldsymbol{X}) - \widetilde{f}_{nn}(\boldsymbol{X})\right| > t\right)\right], \tag{15}$$

where the last inequality follows from the fact that  $|\widetilde{\eta}(\mathbf{X}) - 1/2| \leq |\widetilde{\eta}(\mathbf{X}) - \widetilde{f}_{nn}(\mathbf{X})|$ , when  $\widetilde{s}_{nn}(\mathbf{X}) \neq f^*(\mathbf{X})$ . Next, by Cauchy–Schwarz inequality, (15) can be further bounded as

$$2\mathbb{E}\left[\left|\widetilde{\eta}(\boldsymbol{X}) - \widetilde{f}_{nn}(\boldsymbol{X})\right| \cdot I\left(\left|\widetilde{\eta}(\boldsymbol{X}) - 1/2\right| \leq t\right)\right] + 2\mathbb{E}\left[\left|\widetilde{\eta}(\boldsymbol{X}) - \widetilde{f}_{nn}(\boldsymbol{X})\right| \cdot I\left(\left|\widetilde{\eta}(\boldsymbol{X}) - 1/2\right| > t\right)\right]$$

$$\leq 2\|\widetilde{\eta} - \widetilde{f}_{nn}\|_{L^{2}(\mathbb{P}_{\boldsymbol{X}})} \sqrt{\mathbb{P}(\left|\widetilde{\eta}(\boldsymbol{X}) - 1/2\right| \leq t)} + 2\|\widetilde{\eta} - \widetilde{f}_{nn}\|_{L^{2}(\mathbb{P}_{\boldsymbol{X}})}^{2}/t$$

$$\leq 2^{1-\gamma/2}\|\widetilde{\eta} - \widetilde{f}_{nn}\|_{L^{2}(\mathbb{P}_{\boldsymbol{X}})} c\kappa_{\epsilon}^{-\gamma/2} t^{\gamma/2} + 2\|\widetilde{\eta} - \widetilde{f}_{nn}\|_{L^{2}(\mathbb{P}_{\boldsymbol{X}})}^{2}/t.$$

Choosing  $t = \|\widetilde{f}_{nn} - \widetilde{\eta}\|_{L^2(\mathbb{P}_X)}^{\frac{2}{\gamma+2}} \kappa_{\epsilon}^{\frac{\gamma}{\gamma+2}}$ , we get

$$\widetilde{R}(\widetilde{s}_{nn}) - \widetilde{R}(f^*) \lesssim \kappa_{\epsilon}^{-\frac{\gamma}{\gamma+2}} \|\widetilde{\eta} - \widetilde{f}_{nn}\|_{L^2(\mathbb{P}_{\mathbf{X}})}^{\frac{2\gamma+2}{\gamma+2}}$$

Subsequently, by Lemma 2, it follows that

$$R(\widetilde{s}_{nn}) - R(f^*) = \kappa_{\epsilon}^{-1} \left( \widetilde{R}(\widetilde{s}_{nn}) - \widetilde{R}(f^*) \right) \lesssim \left( \kappa_{\epsilon}^{-2} \| \widetilde{\eta} - \widetilde{f}_{nn} \|_{L^2(\mathbb{P}_{\boldsymbol{X}})}^2 \right)^{\frac{\gamma+1}{\gamma+2}}.$$
 (16)

Next, we proceed to establish the convergence rate of  $\|\widetilde{\eta} - \widetilde{f}_{nn}\|_{L^2(\mathbb{P}_X)}^2$ . Let  $\mathcal{F}_{n,\delta_n}^{NN} = \{f \in \mathcal{F}_n^{NN} : \|\widetilde{\eta} - f\|_{L^2(\mathbb{P}_X)}^2 \ge \delta_n\}$ . For any  $\delta_n > 0$ ,

$$\mathbb{P}\Big(\|\widetilde{\eta} - \widetilde{f}_{nn}\|_{L^2(\mathbb{P}_{\mathbf{X}})}^2 \ge \delta_n\Big) \le \mathbb{P}\Big(\inf_{f \in \mathcal{F}_{n,\delta_n}^{NN}} n^{-1} \sum_{i=1}^n (f_{nn}^*(\boldsymbol{x}_i) - \widetilde{z}_i)^2 - n^{-1} \sum_{i=1}^n (f(\boldsymbol{x}_i) - \widetilde{z}_i)^2 \ge \delta_n\Big),$$

where  $f_{nn}^* = \arg\min_{f \in \mathcal{F}_n^{NN}} \|f - \widetilde{\eta}\|_{L^2(\mathbb{P}_{\mathbf{X}})}^2$ .

For ease of notation, we denote that  $U_n(f) = n^{-1} \sum_{i=1}^n (f(\boldsymbol{x}_i) - \widetilde{z}_i)^2$  and  $U(f) = \mathbb{E}(f(\boldsymbol{X}) - \widetilde{Z})^2$ . Then, we have

$$\mathbb{P}\Big(\|\widetilde{\eta} - \widetilde{f}_{nn}\|_{L^2(\mathbb{P}_{\mathbf{X}})}^2 \ge \delta_n\Big) \le \mathbb{P}\Big(\inf_{f \in \mathcal{F}_{nn,\delta_n}} U_n(f_{nn}^*) - U_n(f) \ge 0\Big).$$

Notice that  $\mathcal{F}_{n,\delta_n}^{NN}$  admits the decomposition as  $\mathcal{F}_{n,\delta_n}^{NN} = \bigcup_{i=1}^n \mathcal{H}_i$  with  $\mathcal{H}_i = \{ f \in \mathcal{F}_n^{NN} : 2^{i-1}\delta_n \leq \|\widetilde{\eta} - f\|_{L^2(\mathbb{P}_X)}^2 \leq 2^i\delta_n \}$ . Therefore, we further have

$$\mathbb{P}\Big(\|\widetilde{\eta} - \widetilde{f}_{nn}\|_{L^2(\mathbb{P}_{\mathbf{X}})}^2 \ge \delta_n\Big) \le \sum_{i=1}^n \mathbb{P}\Big(\inf_{f \in \mathcal{H}_i} U_n(f_{nn}^*) - U_n(f) \ge 0\Big).$$

Clearly, it suffices to bound  $\mathbb{P}\Big(\inf_{f\in\mathcal{H}_i}U_n(f_{nn}^*)-U_n(f)\geq 0\Big)$  for upper bounding  $\mathbb{P}\Big(\|\widetilde{\eta}-\widetilde{f}_{nn}\|_{L^2(\mathbb{P}_X)}^2\geq \delta_n\Big)$ . For  $i\geq 1$ ,

$$\mathbb{P}\left(\inf_{f \in \mathcal{H}_{i}} U_{n}(f_{nn}^{*}) - U_{n}(f) \geq 0\right) \\
\leq \mathbb{P}\left(\inf_{f \in \mathcal{H}_{i}} \left[U_{n}(f_{nn}^{*}) - U(f_{nn}^{*})\right] - \left[U_{n}(f) - U(f)\right] \geq \inf_{f \in \mathcal{H}_{i}} U(f) - U(f_{nn}^{*})\right) \\
= \mathbb{P}\left(\inf_{f \in \mathcal{H}_{i}} \left[U_{n}(f_{nn}^{*}) - U(f_{nn}^{*})\right] - \left[U_{n}(f) - U(f)\right] \geq \inf_{f \in \mathcal{H}_{i}} U(f) - U(\widetilde{\eta}) + U(\widetilde{\eta}) - U(f_{nn}^{*})\right) \\
\leq \mathbb{P}\left(\inf_{f \in \mathcal{H}_{i}} \left[U_{n}(f_{nn}^{*}) - U(f_{nn}^{*})\right] - \left[U_{n}(f) - U(f)\right] \geq 2^{i-1}\delta_{n} - \|f_{nn}^{*} - \widetilde{\eta}\|_{L^{2}(\mathbb{P}_{\mathbf{X}})}^{2}\right).$$

Assuming  $||f_{nn}^* - \widetilde{\eta}||_{L^2(\mathbb{P}_X)}^2 \le \delta_n/2$  yields that

$$\mathbb{P}\Big(\inf_{f\in\mathcal{H}_i}U_n(f_{nn}^*)-U_n(f)\geq 0\Big)\leq \mathbb{P}\Big(\inf_{f\in\mathcal{H}_i}\left[U_n(f_{nn}^*)-U(f_{nn}^*)\right]-\left[U_n(f)-U(f)\right]\geq 2^{i-2}\delta_n\Big).$$

Denote that  $M_i = 2^{i-2}\delta_n$ . We turn to establish the relation between the variance of  $(f_{nn}^*(\boldsymbol{X}) - \widetilde{Z})^2 - (f(\boldsymbol{X}) - \widetilde{Z})^2$  and  $M_i$ .

$$\sup_{f \in \mathcal{H}_{i}} \operatorname{Var} \left[ \left( f_{nn}^{*}(\boldsymbol{X}) - \widetilde{Z} \right)^{2} - \left( f(\boldsymbol{X}) - \widetilde{Z} \right)^{2} \right] \\
= \sup_{f \in \mathcal{H}_{i}} \operatorname{Var} \left[ \left( f_{nn}^{*}(\boldsymbol{X}) - f(\boldsymbol{X}) \right) \left( f_{nn}^{*}(\boldsymbol{X}) + f(\boldsymbol{X}) - 2\widetilde{Z} \right) \right] \\
\leq \sup_{f \in \mathcal{H}_{i}} \mathbb{E} \left[ \left( f_{nn}^{*}(\boldsymbol{X}) - f(\boldsymbol{X}) \right)^{2} \left( f_{nn}^{*}(\boldsymbol{X}) + f(\boldsymbol{X}) - 2\widetilde{Z} \right)^{2} \right] \leq \sup_{f \in \mathcal{H}_{i}} 4V_{n}^{2} \| f_{nn}^{*} - f \|_{L^{2}(\mathbb{P}_{\boldsymbol{X}})}^{2} \\
\leq \sup_{f \in \mathcal{H}_{i}} 8V_{n}^{2} \left( \| f_{nn}^{*} - \widetilde{\eta} \|_{L^{2}(\mathbb{P}_{\boldsymbol{X}})}^{2} + \| \widetilde{\eta} - f \|_{L^{2}(\mathbb{P}_{\boldsymbol{X}})}^{2} \right) \leq 64V_{n}^{2} M_{i} \equiv V_{i}.$$

In the following, we proceed to verify conditions (4.5)-(4.7) in Shen & Wong (1994). First, the relation between  $M_i$  and  $V_i$  directly implies (4.6) with  $T = 32V_n^2$  and  $\epsilon = 1/2$ . Second, by Lemma 5 of Schmidt-Hieber (2020),

$$\log \mathcal{N}\left(\epsilon, \mathcal{F}_n^{NN}(L_n, N_n, P_n, B_n, V_n), \|\cdot\|_{L^{\infty}(\mathbb{P}(\mathbf{X}))}\right) \le 2L_n(P_n + 1)\log\left(\epsilon^{-1}(L_n + 1)(N_n + 1)\max\{B_n, 1\}\right).$$

It then follows that

$$\int_{\frac{\epsilon}{32}M_i}^{V_i^{1/2}} \sqrt{\log \mathcal{N}\left(\epsilon, \mathcal{F}_n^{NN}(L_n, N_n, P_n, B_n, V_n), \|\cdot\|_{L^{\infty}(\mathbb{P}(\mathbf{X})}\right)} d\epsilon/M_i$$

$$\leq \int_{\frac{\epsilon}{32}M_i}^{V_i^{1/2}} \sqrt{2L_n(P_n+1)\log\left(\epsilon^{-1}(L_n+1)(N_n+1)\max\{B_n, 1\}\right)} d\epsilon/M_i. \tag{17}$$

Notice that the right-hand side of (17) is non-increasing in i and  $M_i$ , it then follows that

$$\int_{\frac{\epsilon}{32}M_{i}}^{V_{i}^{1/2}} \sqrt{2L_{n}(P_{n}+1)\log\left(\epsilon^{-1}(L_{n}+1)(N_{n}+1)\max\{B_{n},1\}\right)} d\epsilon/M_{i}$$

$$\leq \int_{\frac{\epsilon}{22}M_{i}}^{V_{i}^{1/2}} \sqrt{2L_{n}(P_{n}+1)\log\left(\epsilon^{-1}(L_{n}+1)(N_{n}+1)\max\{B_{n},1\}\right)} d\epsilon/M_{1}.$$

With this, condition (4.7) can be satisfied by imposing

$$\int_{\frac{\epsilon}{32}M_1}^{V_1^{1/2}} \sqrt{2L_n(P_n+1)\log\left(\epsilon^{-1}(L_n+1)(N_n+1)\max\{B_n,1\}\right)} d\epsilon/M_1 \lesssim n^{1/2},\tag{18}$$

which directly implies the condition (4.5) by appropriate choices of  $L_n$  and  $S_n$ . By Theorem 3 in Shen & Wong (1994), we get

$$\mathbb{P}\Big(\|\widetilde{\eta} - \widetilde{f}_{nn}\|_{L^2(\mathbb{P}_{\mathbf{X}})}^2 \ge \delta_n\Big) \lesssim \sum_{i=1}^{\infty} \exp\Big(-\frac{nM_i^2}{512V_n^2 M_i + 2M_i/3}\Big) \lesssim \sum_{i=1}^{\infty} \exp\Big(-n2^{i-2}\delta_n\Big).$$

By the fact  $2^{i-2} \ge (i-1/2)$  for  $i \ge 1$ , there exists some constants C such that

$$\mathbb{P}\Big(\|\widetilde{\eta} - \widetilde{f}_{nn}\|_{L^2(\mathbb{P}_{\mathbf{X}})}^2 \ge \delta_n\Big) \lesssim \exp(-Cn\delta_n),$$

provided that  $n\delta = o(1)$ .

Next, we proceed to consider the approximation error of neural network to meet the assumption that  $||f_{nn}^* - \eta||_{L^{\infty}(\mathbb{P}_{\mathbf{X}})}^2 \leq \delta_n$ . Notice that  $\eta(\mathbf{X}) \in \mathcal{H}(\beta, [0, 1]^p, M)$ . By Theorem 5 of Schmidt-Hieber (2020), there exists a class of neural networks  $\mathcal{F}_n^{NN}(L_n, N_n, P_n, B_n, V_n)$  such that for any  $0 < \psi_n < 1$ 

$$\inf_{f \in \mathcal{F}_{\epsilon}^{NN}(L_n, N_n, P_n, B_n, V_n)} \|f - \eta\|_{L^{\infty}(\mathbb{P}_{\mathbf{X}})} \le \kappa_{\epsilon}^{-1} \psi_n,$$

where  $P_n \simeq (\kappa_{\epsilon}^{-1}\psi_n)^{-p/\beta} \log(\kappa_{\epsilon}/\psi_n)$ ,  $N_n \simeq (\kappa_{\epsilon}^{-1}\psi_n)^{-p/\beta}$ ,  $B_n = 1$ ,  $V_n \geq M+1$  and  $L_n \simeq \log(\kappa_{\epsilon}/\psi_n)$ . Let  $\eta_{nn}^*$  denote the optimal function in  $\mathcal{F}_n^{NN}(L_n, N_n, P_n, B_n, \infty)$  to approximate  $\eta$ . Suppose that  $\eta_{nn}^*$  is a L-layer neural network and formulated as

$$\eta_{nn}^*(\boldsymbol{x}) = \boldsymbol{A}_{L+1} \boldsymbol{g}_L(\boldsymbol{x}) + \boldsymbol{b}_{L+1},$$

where  $g_L(\mathbf{x}) = \mathbf{h}_L \circ \mathbf{h}_{L-1} \circ \cdots \circ h_1(\mathbf{x})$ . We construct a new neural network that  $\widetilde{\eta}_{nn}$  as

$$\widetilde{\eta}_{nn}(\boldsymbol{x}) = \theta \eta_{nn}^*(\boldsymbol{x}) + (1 - \theta)(1 - \eta_{nn}^*(\boldsymbol{x})) 
= \theta \boldsymbol{A}_{L+1} \boldsymbol{g}_L(\boldsymbol{x}) + \theta \boldsymbol{b}_{L+1} + (1 - \theta)(-\boldsymbol{A}_{L+1} \boldsymbol{g}_L(\boldsymbol{x}) + (1 - \boldsymbol{b}_{L+1})) 
= (2\theta - 1)\boldsymbol{A}_{L+1} \boldsymbol{g}_L(\boldsymbol{x}) + (2\theta - 1)\boldsymbol{b}_{L+1} + (1 - \theta)\boldsymbol{1}_{L+1}.$$

It can be easily verified that  $\widetilde{\eta}_{nn} \in \mathcal{F}_n^{NN}(L_n, N_n, P_n, B_n, V_n)$ . This along with the fact that  $\widetilde{\eta}(\boldsymbol{x}) = \theta \eta(\boldsymbol{x}) + (1 - \theta)(1 - \eta(\boldsymbol{x}))$  with  $B_n \geq 1$  results in

$$\|\widetilde{\eta}_{nn} - \widetilde{\eta}\|_{L^{\infty}(\mathbb{P}_{\mathbf{X}})} = \|\theta\eta_{nn}^{*} + (1 - \theta)(1 - \eta_{nn}^{*}) - \theta\eta - (1 - \theta)(1 - \eta)\|_{L^{\infty}(\mathbb{P}_{\mathbf{X}})}$$
$$= (2\theta - 1)\|\eta_{nn}^{*} - \eta\|_{L^{\infty}(\mathbb{P}_{\mathbf{X}})} = \kappa_{\epsilon}\|\eta_{nn}^{*} - \eta\|_{L^{\infty}(\mathbb{P}_{\mathbf{X}})} \le \psi_{n}.$$

Therefore,

$$||f_{nn}^* - \widetilde{\eta}||_{L^{\infty}(\mathbb{P}_{\mathbf{X}})}^2 = \inf_{f \in \mathcal{F}_n^{NN}(L_n, N_n, P_n, B_n, V_n)} ||f - \widetilde{\eta}||_{L^{\infty}(\mathbb{P}_{\mathbf{X}})}^2 \le ||\widetilde{\eta}_{nn} - \widetilde{\eta}||_{L^{\infty}(\mathbb{P}_{\mathbf{X}})}^2 \le \psi_n^2$$

Plugging  $P_n \simeq (\kappa_{\epsilon}^{-1}\psi_n)^{-p/\beta} \log(\kappa_{\epsilon}/\psi_n)$ ,  $N_n \simeq (\kappa_{\epsilon}^{-1}\psi_n)^{-p/\beta}$ , and  $L_n \simeq \log(\kappa_{\epsilon}/\psi_n)$  into (18) yields that  $\kappa_{\epsilon}^{p/(2\beta)}\psi_n^{-p/(2\beta)}\log(\psi_n^{-1}) \lesssim (n\delta_n)^{1/2}$ . Combining this with the assumption that  $\psi_n^2 \lesssim \delta_n$ , it follows that  $\delta_n \simeq \psi_n^2 \simeq \kappa_{\epsilon}^{2p/(2\beta+p)}(\log n/n)^{2\beta/(2\beta+p)}$ . Plugging this into (16) yields that

$$\mathbb{E}\Big[R(\widetilde{s}_{nn}) - R(f^*)\Big] \lesssim \left(\frac{\log n}{n\kappa^2}\right)^{\frac{2\beta(\gamma+1)}{2\beta(\gamma+2) + p(\gamma+2)}}.$$
(19)

Notice that the proof of (19) is independent of the distribution of X, therefore (19) holds for any distribution in  $\mathcal{P}_{\gamma,\beta}$ , which implies that

$$\sup_{\pi \in \mathcal{P}_{\gamma,\beta}} \mathbb{E} \Big[ R(\widetilde{s}_{nn}) - R(f^*) \Big] \lesssim \Big( \frac{\log n}{n \kappa_{\epsilon}^2} \Big)^{\frac{2\beta(\gamma+1)}{2\beta(\gamma+2) + p(\gamma+2)}}.$$

Next, we proceed to prove the lower bound. The proof is based on the well-known Assouad's lemma. The overall proof for the lower bound is mainly based on the proofs Theorem 3.5 and 4.1 in Audibert & Tsybakov (2007), which employs the Assouad's lemma.

We first introduce the partition  $\{\mathcal{X}_i\}_{i=0}^m$  of the cube  $[0,1]^p$  using the grid  $G_q\subseteq [0,1]^p$  defined by

$$G_q = \left\{ \left( \frac{2k_1 + 1}{2q}, \dots, \frac{2k_p + 1}{2q} \right) : k_i \in \{0, \dots, q - 1\}, i \in \{1, \dots, p\} \right\},$$

where  $q \geq 1$  is an integer. For any  $\boldsymbol{x} \in \mathbb{R}^p$ , let  $n_q(\boldsymbol{x}) \in G_q$  be the unique point which is the closest point to  $\boldsymbol{x} \in \mathbb{R}^p$  among all points in  $G_q$ . Without loss of generality, we assume the uniqueness of  $n_q(\boldsymbol{x})$  by choosing the one closest to 0. We define a partition  $\{\mathcal{X}_i'\}_{i=1}^{q^p}$  of  $\mathbb{R}^p$  as follows. For  $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^p$ ,  $\boldsymbol{x}$  and  $\boldsymbol{y}$  are in the same cell  $\mathcal{X}_i$  if and only if  $n_q(\boldsymbol{x}) = n_q(\boldsymbol{y})$ . Then, for  $m \leq q^p$ , we define  $\mathcal{X}_i$  as  $\mathcal{X}_i = \mathcal{X}_i'$  for  $1 \leq i \leq m$  and  $\mathcal{X}_0 = \mathbb{R}^p / \bigcup_{i=1}^m \mathcal{X}_i$ .

Let  $u: \mathbb{R}_+ \to \mathbb{R}_+$  be a non-increasing infinitely differentiable function such that u=1 on [0,1/4] and u=0 on  $[1/2,\infty)$ . Let  $\phi:=C_\phi u(\|\boldsymbol{x}\|)$ , where  $C_\phi\leq 1$  is taken small enough such that  $\phi\in\mathcal{H}(\beta,[0,1]^p,M)$ . For a given  $\boldsymbol{\sigma}=(\sigma_1,\cdots,\sigma_m)\in\{\pm 1\}^m$ , we construct a distribution  $\pi_{\boldsymbol{\sigma}}$  on  $\mathbb{R}^p\times\{-1,1\}$  as follows. Let  $\mu$  be the Lebesgue measure. For  $0<\omega<\frac{1}{m}$  and  $A_0\subseteq\mathcal{X}_0$  with  $\mu(A_0)>0$ , we construct the marginal distribution  $\mathbb{P}_{\boldsymbol{X}}$  on  $\mathbb{R}^p$  that has the density function

$$\mathbb{P}_{\boldsymbol{X}}(\boldsymbol{x}) = \left\{ \begin{array}{ll} \frac{\omega}{\mu(\mathcal{B}(0,q^{-1}/4))}, & \boldsymbol{x} \in \mathcal{B}(z,q^{-1}/4) \text{ for some } z \in G_q, \\ (1-m\omega)/\mu(A_0), & \boldsymbol{x} \in A_0, \\ 0, & \text{otherwise.} \end{array} \right.$$

The conditional distribution on  $\{-1,1\}$  is defined by

$$\eta_{\sigma}(\boldsymbol{x}) = \mathbb{P}(Y = 1 | \boldsymbol{X} = \boldsymbol{x}) = \begin{cases} \frac{1 + \sigma_j \varphi(\boldsymbol{x})}{2}, & \text{for } \boldsymbol{x} \in \mathcal{X}_j, j = 1, \cdots, m, \\ 1/2, & \boldsymbol{x} \in \mathcal{X}_0, \end{cases}$$

where  $\varphi(\mathbf{x}) = q^{-\beta}\phi(q(\mathbf{x} - n_q(\mathbf{x})))$ . Correspondingly, for any  $\sigma$ ,  $\widetilde{\eta}_{\sigma}(\mathbf{x})$  is given as

$$\widetilde{\eta}_{\sigma}(\boldsymbol{x}) = \theta \widetilde{\eta}_{\sigma}(\boldsymbol{x}) + (1 - \theta)(1 - \widetilde{\eta}_{\sigma}(\boldsymbol{x})) = \begin{cases} \frac{1 + (2\theta - 1)\sigma_{j}\varphi(\boldsymbol{x})}{2}, & \text{for } \boldsymbol{x} \in \mathcal{X}_{j}, j = 1, \cdots, m, \\ 1/2, & \boldsymbol{x} \in \mathcal{X}_{0}, \end{cases}$$

Notice that  $D^s \varphi = q^{|s|-\beta} D^s(q(\boldsymbol{x} - n_q(\boldsymbol{x})))$  for any  $s \in \mathbb{N}$  with  $|s| \leq \beta$ . Thus,  $\eta_{\boldsymbol{\sigma}}(\boldsymbol{x})$  belongs to  $\mathcal{H}(\beta, [0, 1]^p, M)$ .

$$\begin{split} & \mathbb{P}\Big(|\eta_{\boldsymbol{\sigma}}(\boldsymbol{x}) - 1/2| \leq t\Big) = m\mathbb{P}\Big(\phi(q(\boldsymbol{x} - n_q(\boldsymbol{x}))) \leq 2tq^{\beta}\Big) \\ = & m \int_{\mathcal{B}(\boldsymbol{x}_0, (4q)^{-1})} I\Big(\phi\big(q(\boldsymbol{x} - \boldsymbol{x}_0)\big) \leq 2tq^{\beta}\Big) \frac{w}{\mu(\mathcal{B}(\boldsymbol{0}, (4q)^{-1}))} d\boldsymbol{x} \\ = & m \int_{\mathcal{B}(\boldsymbol{0}, 1/4)} I\Big(\phi(\boldsymbol{x}) \leq 2tq^{\beta}\Big) \frac{w}{\mu(\mathcal{B}(\boldsymbol{0}, 1/4))} d\boldsymbol{x} = mwI\Big(t \geq C_{\phi}/(2q^{\beta})\Big), \end{split}$$

where  $\mathbf{x}_0 = \left(\frac{1}{2K}, \dots, \frac{1}{2K}\right)$ . Clearly, the low-noise assumption of  $\eta_{\boldsymbol{\sigma}}$  can be satisfied by setting  $mw \leq C_{\phi}^{\gamma}/(2q^{\beta})^{\gamma}$ . Let  $\mathcal{P}_{\gamma,\beta}$  denote the set of joint distributions of  $(\boldsymbol{X},Y)$  satisfying the low-noise assumption and  $\eta(\boldsymbol{x}) \in \mathcal{H}(\beta, [0,1]^p, M)$ . For any  $\boldsymbol{\sigma}$ , we have  $\pi_{\boldsymbol{\sigma}} \in \mathcal{P}_{\gamma,\beta}$ , implying  $\mathcal{P}' = \{\pi_{\boldsymbol{\sigma}} : \boldsymbol{\sigma} \in \{-1,1\}^m\} \subset \mathcal{P}_{\gamma,\beta}$ . Therefore,

$$\sup_{\pi \in \mathcal{P}_{\gamma,\beta}} \mathbb{E}_{\widetilde{\mathcal{D}}} \Big[ R(\widetilde{s}_{nn}) - R(f^*) \Big] \ge \sup_{\pi_{\sigma} \in \mathcal{P}'} \mathbb{E}_{\widetilde{\mathcal{D}}} \Big[ R(\widetilde{s}_{nn}) - R(f^*) \Big].$$

Let  $\widetilde{P}'$  be a set of probability measures on  $(\boldsymbol{X}, \widetilde{Y})$  satisfying that for each  $\pi_{\boldsymbol{\sigma}} \in \mathcal{P}'$  there exists an  $\widetilde{\pi}_{\boldsymbol{\sigma}} \in \widetilde{\mathcal{P}}'$  such that  $\pi_{\boldsymbol{\sigma}}$  and  $\widetilde{\pi}_{\boldsymbol{\sigma}}$  have the same marginal distribution of  $\boldsymbol{X}$  and  $\widetilde{\eta}_{\boldsymbol{\sigma}}(\boldsymbol{x}) = \theta \eta_{\boldsymbol{\sigma}}(\boldsymbol{x}) + (1 - \theta)(1 - \eta_{\boldsymbol{\sigma}}(\boldsymbol{x}))$ . It follows that

$$\sup_{\pi_{\sigma} \in \mathcal{P}'} \mathbb{E}_{\widetilde{\mathcal{D}}} \Big[ R(\widetilde{s}_{nn}) - R(f^*) \Big] \ge (2\theta - 1)^{-1} \sup_{\widetilde{\pi} \in \widetilde{\mathcal{P}}'} \mathbb{E}_{\widetilde{\mathcal{D}}} \Big[ \widetilde{R}(\widetilde{s}_{nn}) - \widetilde{R}(f^*) \Big].$$

Next, we proceed to bound  $\sup_{\widetilde{\pi}_{\sigma} \in \widetilde{\mathcal{P}}'} \mathbb{E}_{\widetilde{\mathcal{D}}} \left[ \widetilde{R}(\widetilde{s}_{nn}) - \widetilde{R}(f^*) \right]$ . Notice that  $f^*$  varies with the value of  $\sigma$ , therefore we use  $f_{\sigma}^*$  to characterize its dependence on  $\sigma$ .

$$\begin{split} \sup_{\widetilde{\pi}_{\sigma} \in \widetilde{\mathcal{P}}'} \mathbb{E}_{\widetilde{\mathcal{D}}} \Big[ \widetilde{R}(\widetilde{s}_{nn}) - \widetilde{R}(f^*) \Big] &= \sup_{\widetilde{\pi}_{\sigma} \in \widetilde{\mathcal{P}}'} \Big\{ \mathbb{E}_{\widetilde{\mathcal{D}}} \left[ \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ |2\widetilde{\eta}_{\boldsymbol{\sigma}}(\boldsymbol{X}) - 1| \mathbf{1}_{\{\widetilde{s}_{nn}(\boldsymbol{X}) \neq \sigma_{j}; \boldsymbol{X} \in \mathcal{X}_{j}\}} \right] \right] \Big\} \\ &= (2\theta - 1) \sup_{\widetilde{\pi}_{\sigma} \in \widetilde{\mathcal{P}}'} \Big\{ \mathbb{E}_{\widetilde{\mathcal{D}}} \left[ \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(\boldsymbol{X}) \mathbf{1}_{\{\widetilde{s}_{nn}(\boldsymbol{X}) \neq \sigma_{j}; \boldsymbol{X} \in \mathcal{X}_{j}\}} \right] \right] \Big\} \,. \end{split}$$

Let  $\Pi$  be the distribution of a Rademacher random variable  $\sigma$ , that is,  $\Pi(\sigma = 1) = \Pi(\sigma = -1) = 1/2$ . Then

$$\sup_{\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}} \in \widetilde{\mathcal{P}}'} \left\{ \mathbb{E}_{\widetilde{\mathcal{D}}} \left[ \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(\boldsymbol{X}) \mathbf{1}_{\{\widetilde{s}_{nn}(\boldsymbol{X}) \neq f_{\boldsymbol{\sigma}}^*(\boldsymbol{X})\}} \right] \right] \right\} \geq \mathbb{E}_{\Pi^m} \left\{ \mathbb{E}_{\widetilde{\mathcal{D}}} \left[ \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(\boldsymbol{X}) \mathbf{1}_{\{\widetilde{s}_{nn}(\boldsymbol{X}) \neq f_{\boldsymbol{\sigma}}^*(\boldsymbol{X})\}} \right] \right] \right\}.$$

Note that for  $x \in \mathcal{X}_0$ ,  $||x - n_q(x)|| \ge (2q)^{-1}$  and  $\varphi(x) = 0$ . Thus, we have

$$\mathbb{E}_{\Pi^m} \left\{ \mathbb{E}_{\widetilde{\mathcal{D}}} \left[ \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(\boldsymbol{X}) 1_{\{\widetilde{s}_{nn}(\boldsymbol{X}) \neq f_{\boldsymbol{\sigma}}^*(\boldsymbol{X})\}} \right] \right] \right\} = \sum_{j=1}^m \mathbb{E}_{\Pi^m} \left\{ \mathbb{E}_{\widetilde{\mathcal{D}}} \left[ \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(\boldsymbol{X}) 1_{\{\widetilde{s}_{nn}(\boldsymbol{X}) \neq \sigma_j; \boldsymbol{X} \in \mathcal{X}_j\}} \right] \right] \right\}.$$

Let  $\sigma_{j,r} = (\sigma_1, \dots, \sigma_{j-1}, r, \sigma_{j+1}, \dots, \sigma_m)$  for  $r \in \{0, \pm 1\}$ . Here  $\sigma_{j,r}$  denotes a vector deduced from  $\sigma$  by fixing its j-th element to r. We have

$$\mathbb{E}_{\Pi^{m}} \left\{ \mathbb{E}_{\widetilde{\mathcal{D}}} \left[ \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(\boldsymbol{X}) 1_{\{\widetilde{s}_{nn}(\boldsymbol{X}) \neq \sigma_{j}; \boldsymbol{X} \in \mathcal{X}_{j}\}} \right] \right] \right\}$$

$$= \mathbb{E}_{\Pi^{m}} \left\{ \mathbb{E}_{\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,0}}^{n}} \left[ \frac{d\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}}^{n}}{d\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,0}}^{n}} \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(\boldsymbol{X}) 1_{\{\widetilde{s}_{nn}(\boldsymbol{X}) \neq \sigma_{j}; \boldsymbol{X} \in \mathcal{X}_{j}\}} \right] \right] \right\}$$

$$= \mathbb{E}_{\Pi^{m-1}} \left\{ \mathbb{E}_{\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,0}}^{n}} \mathbb{E}_{\boldsymbol{\sigma}_{j} \sim \Pi} \left[ \frac{d\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}}^{n}}{d\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,0}}^{n}} \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(\boldsymbol{X}) 1_{\{\widetilde{s}_{nn}(\boldsymbol{X}) \neq \sigma_{j}; \boldsymbol{X} \in \mathcal{X}_{j}\}} \right] \right] \right\}.$$

where  $\widetilde{\pi}_{\sigma_{j,0}}$  has the same marginal distribution as  $\mathbb{P}_{X}$  and  $\widetilde{\pi}_{\sigma_{j,0}}$  and  $\widetilde{\pi}_{\sigma}$  differ in the conditional distribution over the points in  $\mathcal{X}_{j}$ . Specifically,

$$\widetilde{\eta}_{\sigma_{j,0}}(\boldsymbol{x}) = \begin{cases} 1/2, & \text{if } \boldsymbol{x} \in \mathcal{X}_j, \\ \widetilde{\eta}_{\boldsymbol{\sigma}}(\boldsymbol{x}), & \text{otherwise.} \end{cases}$$

It then follows that

$$\mathbb{E}_{\Pi^{m-1}} \left\{ \mathbb{E}_{\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,0}}^{n}} \mathbb{E}_{\boldsymbol{\sigma}_{j} \sim \Pi} \left[ \frac{d\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}}^{n}}{d\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,0}}^{n}} \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(\boldsymbol{X}) \mathbf{1}_{\{\widetilde{\boldsymbol{s}}_{nn}(\boldsymbol{X}) \neq \sigma_{j}; \boldsymbol{X} \in \mathcal{X}_{j}\}} \right] \right] \right\}$$

$$= \mathbb{E}_{\Pi^{m-1}} \left\{ \mathbb{E}_{\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,0}}^{n}} \mathbb{E}_{\boldsymbol{\sigma}_{j} \sim \Pi} \left[ \frac{d\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}}^{n}}{d\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,0}}^{n}} \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(\boldsymbol{X}) \mathbf{1}_{\{\boldsymbol{X} \in \mathcal{X}_{j}\}} \mathbf{1}_{\{\widetilde{\boldsymbol{s}}_{nn}(\boldsymbol{X}) \neq \sigma_{j}\}} \right] \right] \right\}$$

$$= \mathbb{E}_{\Pi^{m-1}} \left\{ \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(\boldsymbol{X}) \mathbf{1}_{\{\boldsymbol{X} \in \mathcal{X}_{j}\}} \mathbb{E}_{\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,0}}^{n}} \mathbb{E}_{\boldsymbol{\sigma}_{j} \sim \Pi} \left[ \frac{d\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}}^{n}}{d\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,0}}^{n}} \mathbf{1}_{\{\widetilde{\boldsymbol{s}}_{nn}(\boldsymbol{X}) \neq \sigma_{j}\}} \right] \right] \right\}$$

$$\geq \mathbb{E}_{\Pi^{m-1}} \left\{ \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(\boldsymbol{X}) \mathbf{1}_{\{\boldsymbol{X} \in \mathcal{X}_{j}\}} \left( \frac{1}{2} \mathbb{P}_{\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,1}}^{n}} (\widetilde{\boldsymbol{s}}_{nn}(\boldsymbol{X}) \neq 1) + \frac{1}{2} \mathbb{P}_{\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,-1}}^{n}} (\widetilde{\boldsymbol{s}}_{nn}(\boldsymbol{X}) \neq -1) \right) \right] \right\}$$

$$\geq \frac{1}{2} \mathbb{E}_{\Pi^{m-1}} \left\{ \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(\boldsymbol{X}) \mathbf{1}_{\{\boldsymbol{X} \in \mathcal{X}_{j}\}} \right] \left( 1 - \mathrm{TV} (\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,1}}^{n}, \widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,-1}}^{n}) \right) \right\},$$

where TV is the total variation distance between two distributions. Since  $\mathbb{P}_{\mathbf{X}}(\mathcal{X}_i) = \omega$ , we have

$$\mathbb{E}_{\Pi^m} \left\{ \mathbb{E}_{\widetilde{\mathcal{D}}} \left[ \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(\boldsymbol{X}) 1_{\{\widetilde{\boldsymbol{s}}_{nn}(\boldsymbol{X}) \neq \sigma_j; \boldsymbol{X} \in \mathcal{X}_j\}} \right] \right] \right\} \geq \frac{\omega}{2} \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left( \varphi(\boldsymbol{X}) \big| \boldsymbol{X} \in \mathcal{X}_j \right) \left( 1 - \text{TV}(\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,1}}^n, \widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{j,-1}}^n) \right),$$

Notice that the above inequality holds for any index  $j \in [m]$ . In conclusion, we obtain

$$\sum_{j=1}^{m} \mathbb{E}_{\Pi^{m}} \left\{ \mathbb{E}_{\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}}^{n}} \left[ \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left[ \varphi(X) 1_{[\widetilde{\boldsymbol{s}}_{nn}(\boldsymbol{X}) \neq \sigma_{j}; \boldsymbol{X} \in \mathcal{X}_{j}]} \right] \right] \right\} \geq \frac{m\omega}{2} \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left( \varphi(\boldsymbol{X}) | \boldsymbol{X} \in \mathcal{X}_{j} \right) \left( 1 - \operatorname{TV}(\widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{1,1}}^{n}, \widetilde{\boldsymbol{\pi}}_{\boldsymbol{\sigma}_{1,-1}}^{n}) \right).$$

Now we bound  $\mathrm{TV}(\widetilde{\pi}^n_{\sigma_{1,1}}, \widetilde{\pi}^n_{\sigma_{1,-1}})$ . First, it holds

$$\operatorname{TV}(\widetilde{\pi}_{\boldsymbol{\sigma}_{1,1}}^{n}, \widetilde{\pi}_{\boldsymbol{\sigma}_{1,-1}}^{n}) = \sum_{l=1}^{n} \binom{n}{l} \omega^{l} (1-\omega)^{n-l} \mathcal{V}_{l},$$

where  $\mathcal{V}_l = \text{TV}(\widetilde{\pi}_{-1}^l, \widetilde{\pi}_1^l)$  with  $\widetilde{\pi}_r = \widetilde{\pi}_{\sigma_{1,r}}(\cdot | \boldsymbol{X} \in \mathcal{X}_1)$ . Note that

$$\mathcal{V}_{l} \leq H(\widetilde{\pi}_{-1}^{l}, \widetilde{\pi}_{1}^{l}) = \sqrt{2\left(1 - \left[1 - \frac{H^{2}(\widetilde{\pi}_{-1}, \widetilde{\pi}_{1})}{2}\right]^{l}\right)}$$

where H is the Hellinger distance and

$$1 - \frac{H^2(\widetilde{\pi}_{-1}, \widetilde{\pi}_1)}{2} = \mathbb{E}_{\boldsymbol{X} \sim \mathbb{P}_{\boldsymbol{X}}} \left( \sqrt{1 - (2\theta - 1)^2 \varphi(\boldsymbol{X})} \middle| \boldsymbol{X} \in \mathcal{X}_1 \right) =: \sqrt{1 - b^2}$$

Since  $1-(1-x^2)^{l/2} \leq \frac{lx^2}{2}$  for  $l \geq 2$  and x>0, we have  $\mathcal{V}_l \leq b\sqrt{l}$  and

$$\operatorname{TV}(\widetilde{\pi}_{\boldsymbol{\sigma}_{1,1}}^{n}, \widetilde{\pi}_{\boldsymbol{\sigma}_{1,-1}}^{n}) \leq b \sum_{l=1}^{n} \mathbb{P}\Big(\sum_{i=1}^{n} \epsilon_{i,\omega} = l\Big) \sqrt{l} \leq b \sqrt{n\omega},$$

where  $\epsilon_i$  are i.i.d. random variables such that  $\mathbb{P}(\epsilon_i = 1) = \omega = 1 - \mathbb{P}(\epsilon_i = -1)$ . In conclusion, we have

$$\sup_{\widetilde{\pi}_{\sigma} \in \widetilde{\mathcal{P}}'} \left\{ \mathbb{E}_{\widetilde{\mathcal{D}}} \left[ \widetilde{R}(\widetilde{s}_{nn}) \right] - \widetilde{R}^* \right\} \ge m\omega b' (1 - b\sqrt{n\omega}),$$

where

$$b = \left[ 1 - \left( \mathbb{E} \left[ \sqrt{1 - (2\theta - 1)^2 \varphi^2(X)} \middle| \mathbf{X} \in \mathcal{X}_j \right] \right)^2 \right]^{1/2} \times (2\theta - 1) q^{-\beta},$$
  
$$b' = \mathbb{E} \left( (2\theta - 1) \varphi(\mathbf{X}) \middle| \mathbf{X} \in \mathcal{X}_j \right) \times (2\theta - 1) q^{-\beta}.$$

As a result, we have

$$\sup_{\pi \in \mathcal{P}_{\alpha,\beta}} \mathbb{E}_{\widetilde{\mathcal{D}}} \Big[ R(\widetilde{s}_{nn}) - R^* \Big] \gtrsim m\omega q^{-\beta} (1 - (2\theta - 1)q^{-\beta}\sqrt{n\omega}).$$

Take  $\omega = \frac{q^{2\beta}}{4n(2\theta-1)^2}$  and  $m=q^p,$  we obtain

$$\sup_{\pi \in \mathcal{P}_{\gamma,\beta}} \mathbb{E}_{\widetilde{\mathcal{D}}} \Big[ R(\widetilde{s}_{nn}) - R^* \Big] \gtrsim \frac{q^{\beta+p}}{n\kappa_{\epsilon}^2} \asymp \left( \frac{1}{n\kappa_{\epsilon}^2} \right)^{\frac{(\gamma+1)\beta}{(\gamma+2)\beta+p}}$$

by taking  $q \simeq (n(2\theta - 1)^2)^{\frac{1}{(\gamma + 2)\beta + p}}$ . This completes the proof.