

# Supplementary Material for “Deep regression learning with optimal loss function”

Xuancheng Wang\*, Ling Zhou\* and Huazhen Lin<sup>†‡</sup>

New Cornerstone Science Laboratory,  
Center of Statistical Research and School of Statistics,  
Southwestern University of Finance and Economics, Chengdu, China

This supplement contains four sections. Section [S.1](#) presents the notations and conditions needed in the proofs of Theorems 1-2. Section [S.2](#) presents Proposition [S.1](#), the proofs of Theorems 1-2 and Proposition 1. Section [S.3](#) presents related lemmas used in the proof of Theorems and Corollaries. Section [S.4](#) presents some results of simulation studies.

## Contents

<a href="#">S.1 Notations and conditions</a>	<a href="#">2</a>
<a href="#">S.1.1 Notations</a> . . . . .	<a href="#">2</a>
<a href="#">S.1.2 Conditions</a> . . . . .	<a href="#">4</a>

---

\*Co-first authors.

<sup>†</sup>Corresponding author. Email address: [linhz@swufe.edu.cn](mailto:linhz@swufe.edu.cn).

<sup>‡</sup>The research was supported by National Key R&D Program of China (No.2022YFA1003702), National Natural Science Foundation of China (Nos. 11931014 and 12271441), and New Cornerstone Science Foundation.

<b>S.2 Propositions and Proofs</b>	<b>6</b>
S.2.1 Proposition S.1 . . . . .	6
S.2.2 Proof of Theorem 1 . . . . .	7
S.2.3 Proof of Theorem 2 . . . . .	12
S.2.4 Proof of Proposition 1 . . . . .	27
<b>S.3 Lemmas</b>	<b>32</b>
<b>S.4 Results in numerical studies</b>	<b>34</b>

## S.1 Notations and conditions

### S.1.1 Notations

**Feedforward neural network.** Let  $\mathcal{G}$  be a function class consisting of ReLU neural networks, that is,  $\mathcal{G} := \mathcal{G}_{\mathcal{D}, \mathcal{U}, \mathcal{W}, \mathcal{S}, \mathcal{B}}$ , where the input data is the predictor  $X$ , forming the first layer, and the output is the last layer of the network; Such a network  $\mathcal{G}$  has  $\mathcal{D}$  hidden layers and a total of  $(\mathcal{D} + 2)$  layers. Denote the width of layer  $j$  by  $d_j$ ,  $j = 0, \dots, \mathcal{D}, \mathcal{D} + 1$  with  $d_0 = d$  representing the dimension of the input  $X$ , and  $d_{\mathcal{D}+1} = 1$  representing the dimension of the response  $Y$ . The width  $\mathcal{W}$  is defined as the maximum width among the hidden layers, i.e.,  $\mathcal{W} = \max(d_1, \dots, d_{\mathcal{D}})$ . The size  $\mathcal{S}$  is defined as the total number of parameters in the network  $\mathcal{G}$ , given by  $\mathcal{S} = \sum_{i=0}^{\mathcal{D}} d_{i+1} \times (d_i + 1)$ ; The number of neurons  $\mathcal{U}$  is defined as the total number of computational units in the hidden layers, given by  $\mathcal{U} = \sum_{i=1}^{\mathcal{D}} d_i$ . Further, we assume every function  $g \in \mathcal{G}$  satisfies  $|g|_{\infty} \leq \mathcal{B}$  with  $\mathcal{B}$  being a positive constant.

**Covering number.** Given a  $\delta$ -uniform covering of  $\mathcal{G}$ , we denote the centers of the balls

by  $g_q, q = 1, \dots, \mathcal{N}_{2n}$ , where  $\mathcal{N}_{2n} = \mathcal{N}_{2n}(\delta, \|\cdot\|_\infty, \mathcal{G}|\mathbf{x})$  is the uniform covering number with radius  $\delta$  under the norm  $\|\cdot\|_\infty$ . By the definition of covering, there exists a  $q^*$  such that  $\|\hat{g} - g_{q^*}\|_\infty \leq \delta$  on  $\mathbf{x} \in (\mathbf{X}_1, \dots, \mathbf{X}_n, \mathbf{X}'_1, \dots, \mathbf{X}'_n)$ . Let  $A \preceq B$  represent  $A \leq cB$  for a positive constant  $c$ .

**The definitions of  $\hat{g}$  and  $\hat{g}_{oracle}$ .** For any independent and identically distributed (i.i.d.) samples  $D_n = \{\mathbf{X}_i, Y_i\}_{i=1}^n$  with the sample size  $n$ . With  $g \in \mathcal{G}$ , given  $f(\cdot)$  known, we define the oracle estimator as

$$\hat{g}_{oracle} = \arg \min_{g \in \mathcal{G}} \mathcal{R}_n(g) := \arg \min_{g \in \mathcal{G}} \left\{ \frac{1}{n} \sum_{i=1}^n (-\log f(Y_i - g(\mathbf{X}_i))) \right\}. \quad (\text{S.1})$$

The proposed estimator is defined by

$$\begin{aligned} \tilde{g} &= \hat{g} - \frac{1}{n} \sum_{i=1}^n \hat{g}(\mathbf{X}_i) + \frac{1}{n} \sum_{i=1}^n Y_i, \\ \hat{g} &\in \arg \min_{g \in \mathcal{G}} \hat{\mathcal{R}}_n(g) := \arg \min_{g \in \mathcal{G}} n^{-1} \sum_{i=1}^n \left( -\log \hat{f}_g(Y_i - g(\mathbf{X}_i)) \right) \\ &:= \arg \min_{g \in \mathcal{G}} n^{-1} \sum_{i=1}^n \left( -\log \frac{1}{n} \sum_{j=1}^n \mathcal{K}_h(Y_j - g(\mathbf{X}_j), Y_i - g(\mathbf{X}_i)) \right), \end{aligned} \quad (\text{S.2})$$

where  $\mathcal{K}_h(y_1, y_2) = K(\frac{y_1 - y_2}{h})/h$ ,  $h$  is a bandwidth and  $K(\cdot)$  is a kernel function.

**The following notations are needed in the proofs of Theorems 1-2.** Define  $Y_i = g(\mathbf{X}_i) + \epsilon_i$  with  $\mathbb{E}(\epsilon_i) = 0$ , and

$$S(g, \mathbf{Z}_i) = -\log f(Y_i - g(\mathbf{X}_i)) + \log f(Y_i - g^*(\mathbf{X}_i)),$$

for any  $g$  and the sample  $D_n$  where  $g^*$  is defined as

$$g^* := \arg \min_g \mathcal{R}(g) = \arg \min_g \mathbb{E}(-\log f(Y_i - g(\mathbf{X}_i))).$$

Let  $D'_n$  be another sample independent of  $D_n$ , then the excess risk of  $\hat{g}$  takes the following form:

$$\mathcal{R}(\hat{g}) - \mathcal{R}(g^*) = \mathbb{E}_{D'_n} \left( n^{-1} \sum_{i=1}^n S(\hat{g}, \mathbf{Z}'_i) \right),$$

and its expected excess risk is

$$\mathbb{E}(\mathcal{R}(\hat{g}) - \mathcal{R}(g^*)) = \mathbb{E}_{D_n} \left[ \mathbb{E}_{D'_n} \left( n^{-1} \sum_{i=1}^n S(\hat{g}, \mathbf{Z}'_i) \right) \right].$$

In the following, we write

$$\begin{aligned} L_r(g - g_0) &= \int (g(\mathbf{X}) - g_0(\mathbf{X}))^r f_{\mathbf{X}}(\mathbf{X}) d\mathbf{X}, \\ L(g, \mathbf{Z}_i) &= \mathbb{E}_{D'_n} \{S(g, \mathbf{Z}'_i)\} - 2S(g, \mathbf{Z}_i), \text{ for } g \in \mathcal{G}, \\ \tilde{f}_{g_1, h}(g(\mathbf{x}) - y) &:= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \mathcal{K}_h(g_1(\mathbf{X}_i) - Y_i, g(\mathbf{x}) - y) \right]. \end{aligned}$$

### S.1.2 Conditions

Denote  $f^{(r)}(\cdot)$  to be the  $r$ th derivative of  $f$ , and  $f_{\mathbf{x}}(\cdot)$  to be the density function of covariates  $\mathbf{X}$ , who is supported on a bounded set, and for simplicity, we assume this bounded set to be  $[0, 1]^d$ . In the rest of the paper, the symbol  $c$  denotes a positive constant which may vary across different contexts. The following conditions are required for establishing the rate of the excess risk:

- (C1) Kernel: Let  $U_r = \int K(t)t^r dt$  and  $v_r = \int K^2(t)t^r dt$ . Assume the kernel function  $K(\cdot)$  has a bounded second derivative,  $U_0 = 1$  and  $U_1 = 0$ .
- (C2) Bandwidth:  $h \rightarrow 0$  and  $nh \rightarrow \infty$  as  $n \rightarrow \infty$ .

- (C3) Density function  $f$ : (C3a) For any  $\zeta > 0$ , there exists a  $\beta_\zeta > 0$  such that  $\mathbb{E}(|\log f(\epsilon)I(f(\epsilon) < \beta_\zeta)|) < \zeta$ . For  $\zeta = O(n^{-1} \log n)$ , there exists a  $\beta_\zeta = O(n^{-r})$  for some constant  $r \geq 1$  satisfying the above inequality. (C3b) Assume the density function  $f(\cdot)$  has a continuous first-order derivative and satisfies  $\mathbb{E}(|f^{(1)}/f|) < \mathcal{B}$ . (C3c) Assume the density function  $f(\cdot)$  has a continuous third-order derivative and satisfies  $\mathbb{E}|f^{(r_1)}/f|^{r_2} < \mathcal{B}$  for  $r_1 = 1, 2, 3$ ,  $r_2 = 1, 2$ .
- (C4) Function class for  $g$  and  $g^*$ : For any function  $g \in \mathcal{G}$  and the true function  $g^*$ , we assume  $\|g\|_\infty < \mathcal{B}$  and  $\|g^*\|_\infty < \mathcal{B}$ .

Condition (C1) is a mild condition for the kernel function, which is easy to be satisfied when the kernel function is a symmetric density function. Condition (C2) is the most commonly used assumption for the bandwidth. Condition (C3a) requires a tail condition on the log transformation of the density. Conditions (C3b) and (C3c) require bounded moment for the density and its derivatives to avoid tail-related problems. Clearly, when the errors have bounded supports, conditions (C3b) and (C3c) are automatically satisfied. For errors with unbounded supports, Gaussian errors and several sub-Gaussian errors also meet these conditions. An example is a variable whose characteristic function takes the form:  $\varphi(t) \propto \exp(-\gamma t^2)\psi_b(t)$ , where  $\gamma$  is a positive constant and  $\psi_b(t)$  is a polynomial function of  $t$  with order  $b > 0$ . In particular, a variable with the characteristic function  $\varphi(t) = \exp(-t^2/2)(1 - \alpha t^2 + \beta t^4)$  for  $\alpha \geq \sqrt{2\beta}$  is classified as a strictly sub-Gaussian variable, as defined in Proposition 5.1 in Bobkov et al. (2024) . Condition (C4) is a bounded condition for the function class  $\mathcal{G}$  and the true function  $g^*$ .

## S.2 Propositions and Proofs

### S.2.1 Proposition S.1

**Proposition S.1.** *For any two functions  $g$  and  $g_1$  satisfy the following model:  $Y = g(\mathbf{X}) + \epsilon = g_1(\mathbf{X}) + \epsilon_1$ , where  $\mathbf{X}$  and  $\epsilon$  are independent, and  $\mathbf{X}$  and  $\epsilon_1$  are independent. Denote  $F_\epsilon$  as the distribution function of  $\epsilon$ . Then it follows that for any  $\mathbf{x} \in [0, 1]^d$  and  $m \in R$ ,*

$$g_1(\mathbf{x}) - g(\mathbf{x}) \equiv c, \quad \text{and} \quad F_{\epsilon_1}(m) = F_\epsilon(m - c).$$

for some constant  $c$ .

**Proof of Proposition S.1.** Denote  $\delta(\mathbf{x}) = g_1(\mathbf{x}) - g(\mathbf{x})$ . Then, we have

$$Y = g_1(\mathbf{X}) + \epsilon_1 = g(\mathbf{X}) + (g_1(\mathbf{X}) - g(\mathbf{X}) + \epsilon).$$

Denote  $F_\epsilon(m \mid \mathbf{X} = \mathbf{x}) := P(\epsilon \leq m \mid \mathbf{X} = \mathbf{x})$  and  $F_\epsilon(m) := P(\epsilon \leq m)$ . By the independent assumption on  $\epsilon$  and  $\mathbf{X}$ , it follows that

$$\begin{aligned} F_{\epsilon_1}(m \mid \mathbf{X} = \mathbf{x}) &= \mathbb{E}[I(\epsilon \leq m - g_1(\mathbf{X}) + g(\mathbf{X})) \mid \mathbf{X} = \mathbf{x}] \\ &= F_\epsilon(m - \delta(\mathbf{x}) \mid \mathbf{X} = \mathbf{x}) = F_\epsilon(m - \delta(\mathbf{x})), \end{aligned}$$

and

$$F_{\epsilon_1}(m) = \int F_\epsilon(m - \delta(\mathbf{x}) \mid \mathbf{X} = \mathbf{x}) f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} = \int F_\epsilon(m - \delta(\mathbf{x})) f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x}.$$

By the independent assumption on  $\epsilon_1$  and  $\mathbf{X}$ , it holds that

$$F_{\epsilon_1}(m \mid \mathbf{X} = \mathbf{x}) = F_{\epsilon_1}(m),$$

which leads to

$$F_\epsilon(m - \delta(\mathbf{x})) = \int F_\epsilon(m - \delta(\mathbf{x})) f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x}.$$

Then, it follows that

$$\delta(\mathbf{x}) \equiv c.$$

That is,  $g(\mathbf{x})$  is identifiable up to a constant, and the density of  $\epsilon_1$  is the same as  $\epsilon$  with a constant mean shift.

## S.2.2 Proof of Theorem 1

Let  $g_G^*$  be the estimator in the function class that  $g_G^* = \arg \min_{g \in \mathcal{G}} \mathcal{R}(g)$ . By the definition of the empirical risk minimizer, we have

$$\mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \log f(Y_i - \hat{g}_{oracle}(\mathbf{X}_i)) \right] \geq \mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \log f(Y_i - g_G^*(\mathbf{X}_i)) \right].$$

Then, it follows that

$$\begin{aligned} & \mathbb{E}(\mathcal{R}(\hat{g}_{oracle}) - \mathcal{R}(g^*)) \\ &= \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n [\mathbb{E}_{D'_n}(-\log f(Y'_i - \hat{g}_{oracle}(\mathbf{X}'_i))) - \mathbb{E}_{D'_n}(-\log f(Y'_i - g^*(\mathbf{X}'_i)))] \right\} \\ &= \mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \left\{ \mathbb{E}_{D'_n} \left( n^{-1} \sum_{i=1}^n S(\hat{g}_{oracle}, \mathbf{Z}'_i) \right) - 2S(\hat{g}_{oracle}, \mathbf{Z}_i) \right\} \right] \\ &\quad - 2\mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \{ \log f(Y_i - \hat{g}_{oracle}(\mathbf{X}_i)) - \log f(Y_i - g_G^*(\mathbf{X}_i)) \} \right] \\ &\quad + 2\mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n S(g_G^*, \mathbf{Z}_i) \right] \\ &\leq I_1 + 2(\mathcal{R}(g_G^*) - \mathcal{R}(g^*)). \end{aligned}$$

Next, we will give an upper bound of  $I_1$  and handle it with truncation and classical chaining technique of empirical processes.

**Upper bound for  $I_1$ .**

Write  $\hat{g}_o := \hat{g}_{oracle}$ . Let  $L(g, \mathbf{Z}_i) = \mathbb{E}_{D'_n} \{S(g, \mathbf{Z}'_i)\} - 2S(g, \mathbf{Z}_i)$ . According to the definition of  $S(g, \mathbf{Z}_i)$ , we have

$$\begin{aligned} \left| L(\hat{g}_o, \mathbf{Z}_i) - L(g_{q^*}, \mathbf{Z}_i) \right| &= \left| \mathbb{E}_{D'_n} \{S(\hat{g}_o, \mathbf{Z}'_i)\} - 2S(\hat{g}_o, \mathbf{Z}_i) - \mathbb{E}_{D'_n} \{S(g_{q^*}, \mathbf{Z}'_i)\} + 2S(g_{q^*}, \mathbf{Z}_i) \right|, \\ \left| S(\hat{g}_o, \mathbf{Z}_i) - S(g_{q^*}, \mathbf{Z}_i) \right| &= \left| \log f(\hat{g}_o(\mathbf{X}_i) - Y_i) - \log f(g_{q^*}(\mathbf{X}_i) - Y_i) \right|. \end{aligned}$$

Then given the definition of  $g_{q^*}$ , it follows from

$$\begin{aligned} & \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \{S(\hat{g}_o, \mathbf{Z}_i)\} - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \{S(g_{q^*}, \mathbf{Z}_i)\} \right| \\ & \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} |S(\hat{g}_o, \mathbf{Z}_i) - S(g_{q^*}, \mathbf{Z}_i)| \\ & = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left| \log f(\hat{g}_o(\mathbf{X}_i) - Y_i) - \log f(g_{q^*}(\mathbf{X}_i) - Y_i) \right| \\ & = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left| \frac{\dot{f}(\epsilon)(\hat{g}_o(\mathbf{X}_i) - g_{q^*}(\mathbf{X}_i))}{f(\epsilon)} \right| \\ & \leq \mathcal{B}\delta, \end{aligned}$$

where the last inequality follows from Condition (C3b) and

$$\begin{aligned} & \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n \left| L(\hat{g}_o, \mathbf{Z}_i) - L(g_{q^*}, \mathbf{Z}_i) \right| \right\} \\ & = \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n \left| \mathbb{E}_{D'_n} \{S(\hat{g}_o, \mathbf{Z}'_i)\} - 2S(\hat{g}_o, \mathbf{Z}_i) - \mathbb{E}_{D'_n} \{S(g_{q^*}, \mathbf{Z}'_i)\} + 2S(g_{q^*}, \mathbf{Z}_i) \right| \right\} \\ & \leq \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n \left| \mathbb{E}_{D'_n} \{S(\hat{g}_o, \mathbf{Z}'_i)\} - \mathbb{E}_{D'_n} \{S(g_{q^*}, \mathbf{Z}'_i)\} \right| + 2 \left| S(\hat{g}_o, \mathbf{Z}_i) - S(g_{q^*}, \mathbf{Z}_i) \right| \right\} \\ & \leq 3\mathcal{B}\delta, \end{aligned}$$

that

$$\mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n L(\hat{g}_o, \mathbf{Z}_i) \right\} \leq \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n L(g_{q^*}, \mathbf{Z}_i) \right\} + 3\mathcal{B}\delta. \quad (\text{S.3})$$



Let  $0 < \beta_n$  be a positive number who may depend on the sample size  $n$ . Denote  $T_{\beta_n}f = f$  if  $f \geq \beta_n$  and  $T_{\beta_n}f = \beta_n$  otherwise. Define the function  $g_{\beta_n}^*$  by

$$g_{\beta_n}^*(\mathbf{x}) = \arg \min_{g: |g|_{\infty} < \mathcal{B}} \mathbb{E}(-\log T_{\beta_n}f(Y - g(\mathbf{X})) \mid \mathbf{X} = \mathbf{x}).$$

For any  $g \in \mathcal{G}$ , we let  $S_{\beta_n}(g, \mathbf{Z}_i) = -\log T_{\beta_n}f(g(\mathbf{X}_i) - Y_i) + \log T_{\beta_n}f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i)$ . Then, we have

$$\begin{aligned} \mathbb{E}\{S(g, \mathbf{Z}_i)\} &= \mathbb{E}\{S_{\beta_n}(g, \mathbf{Z}_i)\} + \mathbb{E}\{\log T_{\beta_n}f(g(\mathbf{X}_i) - Y_i) - \log f(g(\mathbf{X}_i) - Y_i)\} \\ &\quad + \mathbb{E}\{\log T_{\beta_n}f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i) - \log T_{\beta_n}f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i)\} \\ &\quad + \mathbb{E}\{\log f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i) - \log T_{\beta_n}f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i)\} \\ &\leq \mathbb{E}\{S_{\beta_n}(g, \mathbf{Z}_i)\} + \mathbb{E}\left|\log T_{\beta_n}f(g(\mathbf{X}_i) - Y_i) - \log f(g(\mathbf{X}_i) - Y_i)\right| \\ &\quad + \mathbb{E}\left|\log f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i) - \log T_{\beta_n}f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i)\right| \\ &\leq \mathbb{E}\{S_{\beta_n}(g, \mathbf{Z}_i)\} + 2\mathbb{E}\left|(\log \beta_n - \log f(\epsilon))I(f(\epsilon) < \beta_n)\right| \\ &\leq \mathbb{E}\{S_{\beta_n}(g, \mathbf{Z}_i)\} + 4\mathbb{E}\left\{\left|\log f(\epsilon)I(f(\epsilon) < \beta_n)\right|\right\}, \\ \mathbb{E}\{S_{\beta_n}(g, \mathbf{Z}_i)\} &= \mathbb{E}\{S(g, \mathbf{Z}_i)\} + \mathbb{E}\{\log f(g(\mathbf{X}_i) - Y_i) - \log T_{\beta_n}f(g(\mathbf{X}_i) - Y_i)\} \\ &\quad + \mathbb{E}\{\log f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i) - \log f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i)\} \\ &\quad + \mathbb{E}\{\log T_{\beta_n}f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i) - \log f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i)\} \\ &\leq \mathbb{E}\{S(g, \mathbf{Z}_i)\} + \mathbb{E}\left|\log T_{\beta_n}f(g(\mathbf{X}_i) - Y_i) - \log f(g(\mathbf{X}_i) - Y_i)\right| \\ &\quad + \mathbb{E}\left|\log f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i) - \log T_{\beta_n}f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i)\right| \\ &\leq \mathbb{E}\{S(g, \mathbf{Z}_i)\} + 4\mathbb{E}\left\{\left|\log f(\epsilon)I(f(\epsilon) < \beta_n)\right|\right\}, \end{aligned}$$

so

$$\left| \frac{1}{n} \sum_{i=1}^n \mathbb{E} \{S(g, \mathbf{Z}_i)\} - \frac{1}{n} \sum_{i=1}^n \mathbb{E} \{S_{\beta_n}(g, \mathbf{Z}_i)\} \right| \leq 4 \mathbb{E} \left\{ \left| \log f(\epsilon) I(f(\epsilon) < \beta_n) \right| \right\},$$

which lead to

$$\left| \mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n (L(g_{q^*}, \mathbf{Z}_i) - L_{\beta_n}(g_{q^*}, \mathbf{Z}_i)) \right] \right| \leq 12 \mathbb{E} \left\{ \left| \log f(\epsilon) I(f(\epsilon) < \beta_n) \right| \right\}. \quad (\text{S.4})$$

On the other hand, for any  $g \in \mathcal{G}$ , we have

$$|S_{\beta_n}(g, \mathbf{Z}_i)| \leq 2 |\log \beta_n|,$$

$$\sigma_S^2(g) := \text{Var}(S_{\beta_n}(g, \mathbf{Z}_i)) \leq \mathbb{E} \{S_{\beta_n}^2(g, \mathbf{Z}_i)\} \leq 2 |\log \beta_n| \mathbb{E}(S_{\beta_n}(g, \mathbf{Z}_i)).$$

Following the Bernstein inequality, for any  $t > 0$ , let  $u = t/2 + \sigma_S^2(g)/4 |\log \beta_n|$ , we have

$$\begin{aligned} & P \left\{ \frac{1}{n} \sum_{i=1}^n L_{\beta_n}(g_q, \mathbf{Z}_i) > t \right\} \\ &= P \left\{ \mathbb{E}_{D'_n} \{S_{\beta_n}(g_q, \mathbf{Z}'_i)\} - \frac{2}{n} \sum_{i=1}^n S_{\beta_n}(g_q, \mathbf{Z}_i) > t \right\} \\ &= P \left\{ \mathbb{E}_{D'_n} \{S_{\beta_n}(g_q, \mathbf{Z}'_i)\} - \frac{1}{n} \sum_{i=1}^n S_{\beta_n}(g_q, \mathbf{Z}_i) > \frac{t}{2} + \frac{1}{2} \mathbb{E}_{D'_n} \{S_{\beta_n}(g_q, \mathbf{Z}'_i)\} \right\} \\ &\leq P \left\{ \mathbb{E}_{D'_n} \{S_{\beta_n}(g_q, \mathbf{Z}'_i)\} - \frac{1}{n} \sum_{i=1}^n S_{\beta_n}(g_q, \mathbf{Z}_i) > \frac{t}{2} + \frac{1}{2} \frac{\sigma_S^2(g)}{2 |\log \beta_n|} \right\} \\ &\leq \exp\left(-\frac{nu^2}{2\sigma_S^2(g) + 8u |\log \beta_n|/3}\right) \\ &\leq \exp\left(-\frac{nu^2}{8u |\log \beta_n| + 8u |\log \beta_n|/3}\right) \\ &\leq \exp\left(-\frac{1}{8 + 8/3} \frac{nu}{|\log \beta_n|}\right) \\ &\leq \exp\left(-\frac{1}{16 + 16/3} \frac{nt}{|\log \beta_n|}\right) \\ &= \exp\left(-\frac{Cnt}{|\log \beta_n|}\right), \end{aligned}$$

This leads to a tail probability bound of  $\frac{1}{n} \sum_{i=1}^n L_{\beta_n}(g_{q^*}, \mathbf{Z}_i)$ , that is

$$P \left\{ \frac{1}{n} \sum_{i=1}^n L_{\beta_n}(g_{q^*}, \mathbf{Z}_i) > t \right\} \leq 2\mathcal{N}_{2n} \exp \left( -\frac{Cnt}{|\log \beta_n|} \right).$$

Then for  $a_n > 0$ ,

$$\begin{aligned} \mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n L_{\beta_n}(g_{q^*}, \mathbf{Z}_i) \right] &\leq a_n + \int_{a_n}^{\infty} P \left\{ \frac{1}{n} \sum_{i=1}^n L_{\beta_n}(g_{q^*}, \mathbf{Z}_i) > t \right\} dt \\ &\leq a_n + \int_{a_n}^{\infty} 2\mathcal{N}_{2n} \exp \left( -\frac{Cnt}{|\log \beta_n|} \right) dt \\ &\leq a_n + 2\mathcal{N}_{2n} \exp \left( -a_n \frac{Cn}{|\log \beta_n|} \right) \frac{|\log \beta_n|}{Cn}. \end{aligned}$$

Choosing  $a_n = \log 2\mathcal{N}_{2n} \frac{|\log \beta_n|}{Cn}$ , the above inequality leads to

$$\mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n L_{\beta_n}(g_{q^*}, \mathbf{Z}_i) \right] \leq \frac{C|\log \beta_n|(\log 2\mathcal{N}_{2n} + 1)}{n}. \quad (\text{S.5})$$

Combining inequalities (S.3), (S.4), (S.5), we have

$$\begin{aligned} I_1 &= \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n L(\hat{g}_o, \mathbf{Z}_i) \right\} \\ &\leq \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n L(g_{q^*}, \mathbf{Z}_i) \right\} + 3\mathcal{B}\delta \\ &\leq \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n L_{\beta_n}(g_{q^*}, \mathbf{Z}_i) \right\} + 12\mathbb{E} \left\{ \left| \log f(\epsilon) I(f(\epsilon) < \beta_n) \right| \right\} + 3\mathcal{B}\delta \\ &\leq \frac{C|\log \beta_n|(\log 2\mathcal{N}_{2n} + 1)}{n} + 12\mathbb{E} \left\{ \left| \log f(\epsilon) I(f(\epsilon) < \beta_n) \right| \right\} + 3\mathcal{B}\delta. \end{aligned} \quad (\text{S.6})$$

Let  $\beta_n = n^{-r}$ . Using the above inequalities and Condition (C3a), we obtain that

$$\mathbb{E}(\mathcal{R}(\hat{g}_{oracle}) - \mathcal{R}(g^*)) \preceq \frac{\log n(\log 2\mathcal{N}_{2n}(n^{-1}, \|\cdot\|_{\infty}, \mathcal{G}|_{\mathbf{x}}) + 2)}{n} + (\mathcal{R}(g_{\mathcal{G}}^*) - \mathcal{R}(g^*)).$$

### S.2.3 Proof of Theorem 2

Recall that  $g_{\mathcal{G}}^*$  is the estimator in the function class that  $g_{\mathcal{G}}^* = \arg \min_{g \in \mathcal{G}} \mathcal{R}(g)$  and

$$\hat{f}_g(z) = \frac{1}{n} \sum_{i=1}^n \mathcal{K}_h(Y_i - g(\mathbf{X}_i), z). \quad (\text{S.7})$$

Since  $\hat{g} \in \arg \min_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \left( -\log \hat{f}_g(Y_i - g(\mathbf{X}_i)) \right)$ , we write  $\hat{g}_c(\cdot) = \hat{g}(\cdot) - \int_{\mathbf{x}} \hat{g}(\mathbf{x}) f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} + \mathbb{E}(Y)$ . It then follows from

$$\hat{f}_g(Y_i - g(\mathbf{X}_i)) = \frac{1}{n} \sum_{j=1}^n \mathcal{K}_h(Y_i - g(\mathbf{X}_i), Y_j - g(\mathbf{X}_j)),$$

and  $\mathcal{K}_h(y, x) = K((y - x)/h)/h$  that

$$\frac{1}{n} \sum_{i=1}^n \left( -\log \hat{f}_{\hat{g}}(Y_i - \hat{g}(\mathbf{X}_i)) \right) \equiv \frac{1}{n} \sum_{i=1}^n \left( -\log \hat{f}_{\hat{g}_c}(Y_i - \hat{g}_c(\mathbf{X}_i)) \right).$$

- We first show that

$$\begin{aligned} \mathbb{E}(\mathcal{R}(\hat{g}_c) - \mathcal{R}(g^*)) &\preceq n^{-1} \log n \log \mathcal{N}_{2n}(n^{-1}, \|\cdot\|_{\infty}, \mathcal{G}|_{\mathbf{x}}) \\ &\quad + (\mathcal{R}(g_{\mathcal{G}}^*) - \mathcal{R}(g^*)) + (\|g_{\mathcal{G}}^* - g^*\|_{\infty}^2 + h^2). \end{aligned} \quad (\text{S.8})$$

By the definition of the empirical risk minimizer, we have

$$\mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \log \hat{f}_{\hat{g}_c}(Y_i - \hat{g}_c(\mathbf{X}_i)) \right] \geq \mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \log \hat{f}_{g_{\mathcal{G}}^*}(Y_i - g_{\mathcal{G}}^*(\mathbf{X}_i)) \right]. \quad (\text{S.9})$$

Then, it follows that

$$\begin{aligned}
& \mathbb{E}(\mathcal{R}(\hat{g}_c) - \mathcal{R}(g^*)) \\
&= \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n [\mathbb{E}_{D'_n}(-\log f(Y'_i - \hat{g}_c(\mathbf{X}'_i))) - \mathbb{E}_{D'_n}(-\log f(Y'_i - g^*(\mathbf{X}'_i)))] \right\} \\
&= \mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \left\{ \mathbb{E}_{D'_n} \left( n^{-1} \sum_{i=1}^n S(\hat{g}_c, \mathbf{Z}'_i) \right) - 2S(\hat{g}_c, \mathbf{Z}_i) \right\} \right] \\
&\quad - 2\mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \left\{ \log f(Y_i - \hat{g}_c(\mathbf{X}_i)) - \log \left( \frac{1}{n} \sum_{j=1}^n \mathcal{K}_h(\hat{g}_c(\mathbf{X}_j) - Y_j, \hat{g}_c(\mathbf{X}_i) - Y_i) \right) \right\} \right] \\
&\quad - 2\mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \left\{ \log \hat{f}_{\hat{g}_c}(Y_i - \hat{g}_c(\mathbf{X}_i)) - \log \hat{f}_{g_{\mathcal{G}}^*}(Y_i - g_{\mathcal{G}}^*(\mathbf{X}_i)) \right\} \right] \\
&\quad - 2\mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \left\{ \log \left( \frac{1}{n} \sum_{j=1}^n \mathcal{K}_h(g_{\mathcal{G}}^*(\mathbf{X}_j) - Y_j, g_{\mathcal{G}}^*(\mathbf{X}_i) - Y_i) \right) - \log f(Y_i - g_{\mathcal{G}}^*(\mathbf{X}_i)) \right\} \right] \\
&\quad + 2\mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n S(g_{\mathcal{G}}^*, \mathbf{Z}_i) \right] \\
&\leq \mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \left\{ \mathbb{E}_{D'_n} \left( n^{-1} \sum_{i=1}^n S(\hat{g}_c, \mathbf{Z}'_i) \right) - 2S(\hat{g}_c, \mathbf{Z}_i) \right\} \right] \\
&\quad - 2\mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \left\{ \log f(Y_i - \hat{g}_c(\mathbf{X}_i)) - \log \left( \frac{1}{n} \sum_{j=1}^n \mathcal{K}_h(\hat{g}_c(\mathbf{X}_j) - Y_j, \hat{g}_c(\mathbf{X}_i) - Y_i) \right) \right\} \right] \\
&\quad - 2\mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \left\{ \log \left( \frac{1}{n} \sum_{j=1}^n \mathcal{K}_h(g_{\mathcal{G}}^*(\mathbf{X}_j) - Y_j, g_{\mathcal{G}}^*(\mathbf{X}_i) - Y_i) \right) - \log f(Y_i - g_{\mathcal{G}}^*(\mathbf{X}_i)) \right\} \right] \\
&\quad + 2(\mathcal{R}(g_{\mathcal{G}}^*) - \mathcal{R}(g^*)) \\
&= I_1 + I_2 + I_3 + I_4.
\end{aligned}$$

Next, we will give an upper bound of  $I_r$  for  $r = 1, 2, 3$  and handle it with truncation and classical chaining technique of empirical processes.

**Upper bound for  $I_1$ .** According to the definition of  $L(g, \mathbf{Z}_i)$  and  $S(g, \mathbf{Z}_i)$ , we have

$$\begin{aligned} \left| L(\hat{g}_c, \mathbf{Z}_i) - L(g_{q^*}, \mathbf{Z}_i) \right| &= \left| \mathbb{E}_{D'_n} \{S(\hat{g}_c, \mathbf{Z}'_i)\} - 2S(\hat{g}_c, \mathbf{Z}_i) - \mathbb{E}_{D'_n} \{S(g_{q^*}, \mathbf{Z}'_i)\} + 2S(g_{q^*}, \mathbf{Z}_i) \right|, \\ \left| S(\hat{g}_c, \mathbf{Z}_i) - S(g_{q^*}, \mathbf{Z}_i) \right| &= \left| \log f(\hat{g}_c(\mathbf{X}_i) - Y_i) - \log f(g_{q^*}(\mathbf{X}_i) - Y_i) \right|. \end{aligned}$$

Then given the definition of  $g_{q^*}$ , it follows from

$$\begin{aligned} & \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \{S(\hat{g}_c, \mathbf{Z}_i)\} - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \{S(g_{q^*}, \mathbf{Z}_i)\} \right| \\ & \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} |S(\hat{g}_c, \mathbf{Z}_i) - S(g_{q^*}, \mathbf{Z}_i)| \\ & = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} |\log f(\hat{g}_c(\mathbf{X}_i) - Y_i) - \log f(g_{q^*}(\mathbf{X}_i) - Y_i)| \\ & = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left| \frac{f^{(1)}(\epsilon)((\hat{g}_c(\mathbf{X}_i) - Y_i) - (g_{q^*}(\mathbf{X}_i) - Y_i))}{f(\epsilon)} \right| \\ & \leq \mathcal{B}\delta, \end{aligned}$$

where the last inequality follows from Condition (C3c), and

$$\begin{aligned} & \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n \left| L(\hat{g}_c, \mathbf{Z}_i) - L(g_{q^*}, \mathbf{Z}_i) \right| \right\} \\ & = \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n \left| \mathbb{E}_{D'_n} \{S(\hat{g}_c, \mathbf{Z}'_i)\} - 2S(\hat{g}_c, \mathbf{Z}_i) - \mathbb{E}_{D'_n} \{S(g_{q^*}, \mathbf{Z}'_i)\} + 2S(g_{q^*}, \mathbf{Z}_i) \right| \right\} \\ & \leq \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n \left| \mathbb{E}_{D'_n} \{S(\hat{g}_c, \mathbf{Z}'_i)\} - \mathbb{E}_{D'_n} \{S(g_{q^*}, \mathbf{Z}'_i)\} \right| + 2 \left| S(\hat{g}_c, \mathbf{Z}_i) - S(g_{q^*}, \mathbf{Z}_i) \right| \right\} \\ & \leq 3\mathcal{B}\delta, \end{aligned}$$

that

$$\mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n L(\hat{g}_c, \mathbf{Z}_i) \right\} \leq \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n L(g_{q^*}, \mathbf{Z}_i) \right\} + 3\mathcal{B}\delta. \quad (\text{S.10})$$

Let  $0 < \beta_n$  be a positive number who may depend on the sample size  $n$ . Denote  $T_{\beta_n}f = f$  if  $f \geq \beta_n$  and  $T_{\beta_n}f = \beta_n$  otherwise. Define the function  $g_{\beta_n}^*$  by

$$g_{\beta_n}^*(\mathbf{x}) = \arg \min_{g: |g|_{\infty} < \mathcal{B}} \mathbb{E}(-\log T_{\beta_n}f(Y - g(\mathbf{X})) \mid \mathbf{X} = \mathbf{x}).$$

For any  $g \in \mathcal{G}$ , we let  $S_{\beta_n}(g, \mathbf{Z}_i) = -\log T_{\beta_n}f(g(\mathbf{X}_i) - Y_i) + \log T_{\beta_n}f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i)$ . Then, we have

$$\begin{aligned} \mathbb{E}\{S(g, \mathbf{Z}_i)\} &= \mathbb{E}\{S_{\beta_n}(g, \mathbf{Z}_i)\} + \mathbb{E}\{\log T_{\beta_n}f(g(\mathbf{X}_i) - Y_i) - \log f(g(\mathbf{X}_i) - Y_i)\} \\ &\quad + \mathbb{E}\{\log T_{\beta_n}f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i) - \log T_{\beta_n}f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i)\} \\ &\quad + \mathbb{E}\{\log f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i) - \log T_{\beta_n}f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i)\} \\ &\leq \mathbb{E}\{S_{\beta_n}(g, \mathbf{Z}_i)\} + \mathbb{E}\left|\log T_{\beta_n}f(g(\mathbf{X}_i) - Y_i) - \log f(g(\mathbf{X}_i) - Y_i)\right| \\ &\quad + \mathbb{E}\left|\log f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i) - \log T_{\beta_n}f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i)\right| \\ &\leq \mathbb{E}\{S_{\beta_n}(g, \mathbf{Z}_i)\} + 2\mathbb{E}\left|(\log \beta_n - \log f(\epsilon))I(f(\epsilon) < \beta_n)\right| \\ &\leq \mathbb{E}\{S_{\beta_n}(g, \mathbf{Z}_i)\} + 4\mathbb{E}\left\{\left|\log f(\epsilon)I(f(\epsilon) < \beta_n)\right|\right\}, \\ \mathbb{E}\{S_{\beta_n}(g, \mathbf{Z}_i)\} &= \mathbb{E}\{S(g, \mathbf{Z}_i)\} + \mathbb{E}\{\log f(g(\mathbf{X}_i) - Y_i) - \log T_{\beta_n}f(g(\mathbf{X}_i) - Y_i)\} \\ &\quad + \mathbb{E}\{\log f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i) - \log f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i)\} \\ &\quad + \mathbb{E}\{\log T_{\beta_n}f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i) - \log f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i)\} \\ &\leq \mathbb{E}\{S(g, \mathbf{Z}_i)\} + \mathbb{E}\left|\log T_{\beta_n}f(g(\mathbf{X}_i) - Y_i) - \log f(g(\mathbf{X}_i) - Y_i)\right| \\ &\quad + \mathbb{E}\left|\log f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i) - \log T_{\beta_n}f(g_{\beta_n}^*(\mathbf{X}_i) - Y_i)\right| \\ &\leq \mathbb{E}\{S(g, \mathbf{Z}_i)\} + 4\mathbb{E}\left\{\left|\log f(\epsilon)I(f(\epsilon) < \beta_n)\right|\right\}, \end{aligned}$$

so

$$\left| \frac{1}{n} \sum_{i=1}^n \mathbb{E} \{S(g, \mathbf{Z}_i)\} - \frac{1}{n} \sum_{i=1}^n \mathbb{E} \{S_{\beta_n}(g, \mathbf{Z}_i)\} \right| \leq 4\mathbb{E} \left\{ \left| \log f(\epsilon) I(f(\epsilon) < \beta_n) \right| \right\},$$

which lead to

$$\left| \mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n (L(g_{q^*}, \mathbf{Z}_i) - L_{\beta_n}(g_{q^*}, \mathbf{Z}_i)) \right] \right| \leq 12\mathbb{E} \left\{ \left| \log f(\epsilon) I(f(\epsilon) < \beta_n) \right| \right\}. \quad (\text{S.11})$$

On the other hand, for any  $g \in \mathcal{G}$ , we have

$$\begin{aligned} |S_{\beta_n}(g, \mathbf{Z}_i)| &\leq 2|\log \beta_n|, \\ \sigma_S^2(g) &:= \text{Var}(S_{\beta_n}(g, \mathbf{Z}_i)) \leq \mathbb{E} \{S_{\beta_n}^2(g, \mathbf{Z}_i)\} \leq 2|\log \beta_n| \mathbb{E}(S_{\beta_n}(g, \mathbf{Z}_i)). \end{aligned}$$

Following the Bernstein inequality, for any  $t > 0$ , let  $u = t/2 + \sigma_S^2(g)/4|\log \beta_n|$ , we



have

$$\begin{aligned}
& P \left\{ \frac{1}{n} \sum_{i=1}^n L_{\beta_n}(g_q, \mathbf{Z}_i) > t \right\} \\
&= P \left\{ \mathbb{E}_{D'_n} \{S_{\beta_n}(g_q, \mathbf{Z}'_i)\} - \frac{2}{n} \sum_{i=1}^n S_{\beta_n}(g_q, \mathbf{Z}_i) > t \right\} \\
&= P \left\{ \mathbb{E}_{D'_n} \{S_{\beta_n}(g_q, \mathbf{Z}'_i)\} - \frac{1}{n} \sum_{i=1}^n S_{\beta_n}(g_q, \mathbf{Z}_i) > \frac{t}{2} + \frac{1}{2} \mathbb{E}_{D'_n} \{S_{\beta_n}(g_q, \mathbf{Z}'_i)\} \right\} \\
&\leq P \left\{ \mathbb{E}_{D'_n} \{S_{\beta_n}(g_q, \mathbf{Z}'_i)\} - \frac{1}{n} \sum_{i=1}^n S_{\beta_n}(g_q, \mathbf{Z}_i) > \frac{t}{2} + \frac{1}{2} \frac{\sigma_S^2(g)}{2|\log \beta_n|} \right\} \\
&\leq \exp\left(-\frac{nu^2}{2\sigma_S^2(g) + 8u|\log \beta_n|/3}\right) \\
&\leq \exp\left(-\frac{nu^2}{8u|\log \beta_n| + 8u|\log \beta_n|/3}\right) \\
&\leq \exp\left(-\frac{1}{8 + 8/3} \frac{nu}{|\log \beta_n|}\right) \\
&\leq \exp\left(-\frac{1}{16 + 16/3} \frac{nt}{|\log \beta_n|}\right) \\
&= \exp\left(-\frac{Cnt}{|\log \beta_n|}\right),
\end{aligned}$$

This leads to a tail probability bound of  $\frac{1}{n} \sum_{i=1}^n L_{\beta_n}(g_{q^*}, \mathbf{Z}_i)$ , that is

$$P \left\{ \frac{1}{n} \sum_{i=1}^n L_{\beta_n}(g_{q^*}, \mathbf{Z}_i) > t \right\} \leq 2\mathcal{N}_{2n} \exp\left(-\frac{Cnt}{|\log \beta_n|}\right).$$

Then for  $a_n > 0$ ,

$$\begin{aligned}
\mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n L_{\beta_n}(g_{q^*}, \mathbf{Z}_i) \right] &\leq a_n + \int_{a_n}^{\infty} P \left\{ \frac{1}{n} \sum_{i=1}^n L_{\beta_n}(g_{q^*}, \mathbf{Z}_i) > t \right\} dt \\
&\leq a_n + \int_{a_n}^{\infty} 2\mathcal{N}_{2n} \exp\left(-\frac{Cnt}{|\log \beta_n|}\right) dt \\
&\leq a_n + 2\mathcal{N}_{2n} \exp\left(-a_n \frac{Cn}{|\log \beta_n|}\right) \frac{|\log \beta_n|}{Cn}.
\end{aligned}$$

Choosing  $a_n = \log 2\mathcal{N}_{2n} \frac{|\log \beta_n|}{C_n}$ , the above inequality leads to

$$\mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n L_{\beta_n}(g_{q^*}, \mathbf{Z}_i) \right] \leq \frac{C |\log \beta_n| (\log 2\mathcal{N}_{2n} + 1)}{n}. \quad (\text{S.12})$$

Combining inequalities (S.10), (S.11), (S.12), we have

$$\begin{aligned} I_1 &= \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n L(\hat{g}_c, \mathbf{Z}_i) \right\} \\ &\leq \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n L(g_{q^*}, \mathbf{Z}_i) \right\} + 3\mathcal{B}\delta \\ &\leq \mathbb{E}_{D_n} \left\{ \frac{1}{n} \sum_{i=1}^n L_{\beta_n}(g_{q^*}, \mathbf{Z}_i) \right\} + 12\mathbb{E} \left\{ \left| \log f(\epsilon) I(f(\epsilon) < \beta_n) \right| \right\} + 3\mathcal{B}\delta \\ &\leq \frac{C |\log \beta_n| (\log 2\mathcal{N}_{2n} + 1)}{n} + 12\mathbb{E} \left\{ \left| \log f(\epsilon) I(f(\epsilon) < \beta_n) \right| \right\} + 3\mathcal{B}\delta. \end{aligned}$$

**Upper bound for  $I_2$  and  $I_3$ .** Recall that  $\tilde{f}_{g_1, h} = \mathbb{E}(\mathcal{K}_h(Y_i - g_1(\mathbf{X}_i), x))$ . For  $I_2$ , we have

$$\begin{aligned} I_2 &= \mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \left\{ \log f(Y_i - \hat{g}_c(\mathbf{X}_i)) - \log \tilde{f}_{\hat{g}_c, h}(Y_i - \hat{g}_c(\mathbf{X}_i)) \right\} \right] \\ &\quad + \mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \left\{ \log \tilde{f}_{\hat{g}_c, h}(Y_i - \hat{g}_c(\mathbf{X}_i)) - \log \left( \frac{1}{n_1} \sum_{j=1}^n \mathcal{K}_h(\hat{g}_c(\mathbf{X}_j) - Y_j, \hat{g}_c(\mathbf{X}_i) - Y_i) \right) \right\} \right] \\ &= I_{2,1} + I_{2,2}. \end{aligned}$$

Denote  $f^{(r)}(\epsilon)$  as the  $r$ th derivative of  $f(\epsilon)$ . We first show that for  $\|g_1 - g^*\|_\infty = o_p(1)$  and  $h \rightarrow 0$ , we have

$$\begin{aligned} \tilde{f}_{g_1, h}(z) &= U_0 f(z) + f^{(1)}(z) [U_0 L_1(g_1 - g^*) + U_1 h] \\ &\quad + 0.5 f^{(2)}(z) [U_0 L_2(g_1 - g^*) + 2U_1 h L_1(g_1 - g^*) + U_2 h^2] (1 + o(1)) \end{aligned} \quad (\text{S.13})$$

Note that

$$\begin{aligned}
& \mathbb{E} \left\{ \frac{1}{n} \sum_{j=1}^n \mathcal{K}_h(Y_j - g_1(\mathbf{X}_j), y - g(\mathbf{x})) \right\} \\
&= \int \frac{1}{h} K \left( \frac{Y - g_1(\mathbf{X}) - y + g(\mathbf{x})}{h} \right) f(Y - g^*(\mathbf{X})) f_{\mathbf{x}}(\mathbf{X}) dY d\mathbf{X} \\
&= \int K(t) f(g_1(\mathbf{X}) - g^*(\mathbf{X}) + y - g(\mathbf{x}) + th) f_{\mathbf{x}}(\mathbf{X}) dt d\mathbf{X} \\
&= \int K(t) \left\{ f(g_1(\mathbf{X}) - g^*(\mathbf{X}) + y - g(\mathbf{x})) + f^{(1)}(g_1(\mathbf{X}) - g^*(\mathbf{X}) + y - g(\mathbf{x})) th \right. \\
&\quad \left. + f^{(2)}(g_1(\mathbf{X}) - g^*(\mathbf{X}) + y - g(\mathbf{x})) \frac{1}{2} t^2 h^2 + o(h^2) \right\} f_{\mathbf{x}}(\mathbf{X}) dt d\mathbf{X} \\
&= U_0 \left( f(y - g(\mathbf{x})) + f^{(1)}(y - g(\mathbf{x})) L_1(g_1 - g^*) \right. \\
&\quad \left. + f^{(2)}(y - g(\mathbf{x})) \frac{1}{2} L_2(g_1 - g^*) + O(L_3(g_1 - g^*)) \right) \\
&\quad + U_1 h \left( f^{(1)}(y - g(\mathbf{x})) + f^{(2)}(y - g(\mathbf{x})) L_1(g_1 - g^*) + O(L_2(g_1 - g^*)) \right) \\
&\quad + \frac{U_2}{2} h^2 \left( f^{(2)}(y - g(\mathbf{x})) + O(L_1(g_1 - g^*)) \right) (1 + o(1)),
\end{aligned}$$

Then, using  $U_0 = 0$  and  $U_1 = 1$ , we obtain that

$$\begin{aligned}
|I_{2,1}| &= \left| \mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \left\{ \log f(Y_i - \hat{g}_c(\mathbf{X}_i)) - \log \tilde{f}_{\hat{g}_c, h}(Y_i - \hat{g}_c(\mathbf{X}_i)) \right\} \right] \right| \tag{S.14} \\
&= \left| \mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \left\{ \log f(Y_i - \hat{g}_c(\mathbf{X}_i)) - \log(f(Y_i - \hat{g}_c(\mathbf{X}_i) + f^{(1)}(Y_i - \hat{g}_c(\mathbf{X}_i)) L_1(\hat{g}_c - g^*)) \right. \right. \right. \\
&\quad \left. \left. + f^{(2)}(Y_i - \hat{g}_c(\mathbf{X}_i)) 0.5 \{L_2(\hat{g}_c - g^*) + U_2 h^2\} \} \right\} \right] \right| \\
&= \left| \mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \left\{ \log \left( 1 + \frac{f^{(1)}(Y_i - \hat{g}_c(\mathbf{X}_i))}{f(Y_i - \hat{g}_c(\mathbf{X}_i))} L_1(\hat{g}_c - g^*) \right) \right. \right. \right. \\
&\quad \left. \left. + \frac{f^{(2)}(Y_i - \hat{g}_c(\mathbf{X}_i))}{f(Y_i - \hat{g}_c(\mathbf{X}_i))} 0.5 \{L_2(\hat{g}_c - g^*) + U_2 h^2\} \right\} \right] \right|.
\end{aligned}$$

Let

$$\begin{aligned} t_1 &= \mathbb{E} \left| \frac{f^{(1)}(\epsilon)}{f(\epsilon)} \right|, \quad t_4 = \mathbb{E} \left| \frac{f^{(2)}(\epsilon)}{f(\epsilon)} \right|, \\ t_2 &= \mathbb{E} \left| \frac{f^{(2)}(\epsilon)f(\epsilon) - (f^{(1)}(\epsilon))^2}{f^2(\epsilon)} \right|, \quad t_3 = \mathbb{E} \left| \frac{f^{(3)}(\epsilon)f(\epsilon) - f^{(2)}(\epsilon)f^{(1)}(\epsilon)}{f^2(\epsilon)} \right|. \end{aligned}$$

Note that

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left\{ \frac{f^{(1)}(Y_i - \hat{g}_c(\mathbf{X}_i))}{f(Y_i - \hat{g}_c(\mathbf{X}_i))} L_1(\hat{g}_c - g^*) \right\} \\ & \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left\{ \frac{f^{(1)}(Y_i - g_{q^*}(\mathbf{X}_i))}{f(Y_i - g_{q^*}(\mathbf{X}_i))} L_1(g_{q^*} - g^*) \right\} \\ & \quad + \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left| \frac{f^{(1)}(Y_i - \hat{g}_c(\mathbf{X}_i))}{f(Y_i - \hat{g}_c(\mathbf{X}_i))} L_1(\hat{g}_c - g^*) - \frac{f^{(1)}(Y_i - g_{q^*}(\mathbf{X}_i))}{f(Y_i - g_{q^*}(\mathbf{X}_i))} L_1(g_{q^*} - g^*) \right| \\ & \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left\{ \frac{f^{(1)}(Y_i - g_{q^*}(\mathbf{X}_i))}{f(Y_i - g_{q^*}(\mathbf{X}_i))} L_1(g_{q^*} - g^*) \right\} \\ & \quad + \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left| \frac{f^{(2)}(\epsilon)f(\epsilon) - (f^{(1)}(\epsilon))^2}{f^2(\epsilon)} \right| \|\hat{g}_c(\mathbf{X}_i) - g_{q^*}(\mathbf{X}_i)\|_{\infty} L_1(|\hat{g}_c - g^*|) \\ & \leq \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left\{ \frac{f^{(1)}(Y_i - g_{q^*}(\mathbf{X}_i))}{f(Y_i - g_{q^*}(\mathbf{X}_i))} L_1(g_{q^*} - g^*) \right\} + t_2 \mathcal{B} \delta, \end{aligned}$$

and

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left\{ \frac{f^{(1)}(Y_i - g_{q^*}(\mathbf{X}_i))}{f(Y_i - g_{q^*}(\mathbf{X}_i))} L_1(g_{q^*} - g^*) \right\} \\
\leq & \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left\{ \frac{f^{(1)}(Y_i - \hat{g}_c(\mathbf{X}_i))}{f(Y_i - \hat{g}_c(\mathbf{X}_i))} L_1(\hat{g}_c - g^*) \right\} + t_2 \mathcal{B} \delta, \\
& \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left[ \frac{f^{(2)}(Y_i - \hat{g}_c(\mathbf{X}_i))}{f(Y_i - \hat{g}_c(\mathbf{X}_i))} 0.5 \{L_2(\hat{g}_c - g^*) + U_2 h^2\} \right] \right. \\
& \quad \left. - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left[ \frac{f^{(2)}(Y_i - g_{q^*}(\mathbf{X}_i))}{f(Y_i - g_{q^*}(\mathbf{X}_i))} 0.5 \{L_2(g_{q^*} - g^*) + U_2 h^2\} \right] \right| \\
\leq & (t_3 \mathcal{B}^2 + t_3 h^2) \delta.
\end{aligned}$$

Then, it follows from expression (S.14) and the above inequalities and Condition

(C3c), we have

$$\begin{aligned}
|I_{2,1}| &\leq \left| \mathbb{E}_{D_n} \left[ \frac{1}{n} \sum_{i=1}^n \frac{f^{(1)}(Y_i - \hat{g}_c(\mathbf{X}_i))}{f(Y_i - \hat{g}_c(\mathbf{X}_i))} L_1(\hat{g}_c - g^*) \right. \right. \\
&\quad \left. \left. + \frac{1}{n} \sum_{i=1}^n \frac{f^{(2)}(Y_i - \hat{g}_c(\mathbf{X}_i))}{f(Y_i - \hat{g}_c(\mathbf{X}_i))} 0.5 \{L_2(\hat{g}_c - g^*) + U_2 h^2\} (1 + o_p(1)) \right] \right| \\
&\leq \left| \frac{1}{n_1} \sum_{i=1}^{n_1} \mathbb{E}_{D_n} \left\{ \frac{f^{(1)}(Y_i - g_{q^*}(\mathbf{X}_i))}{f(Y_i - g_{q^*}(\mathbf{X}_i))} L_1(g_{q^*} - g^*) \right\} \right| \\
&\quad + \left| \frac{1}{n_1} \sum_{i=1}^{n_1} \mathbb{E}_{D_n} \left[ \frac{f^{(2)}(Y_i - g_{q^*}(\mathbf{X}_i))}{f(Y_i - g_{q^*}(\mathbf{X}_i))} 0.5 \{L_2(g_{q^*} - g^*) + U_2 h^2\} \right] \right| (1 + o_p(1)) \\
&\quad + (t_2 \mathcal{B} + t_3 \mathcal{B}^2 + t_3 h^2) \delta \\
&\leq \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left\{ \frac{f^{(1)}(Y_i - g^*(\mathbf{X}_i))}{f(Y_i - g^*(\mathbf{X}_i))} \right\} \mathbb{E}_{D_n} \{L_1(g_{q^*} - g^*)\} \right| + \|g_{q^*} - g^*\|_\infty^2 \\
&\quad + \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left[ \frac{f^{(2)}(Y_i - g_{q^*}(\mathbf{X}_i))}{f(Y_i - g_{q^*}(\mathbf{X}_i))} 0.5 \{L_2(g_{q^*} - g^*) + U_2 h^2\} \right] \right| (1 + o_p(1)) \\
&\quad + (t_2 \mathcal{B} + t_3 \mathcal{B}^2 + t_3 h^2) \delta \\
&\preceq \|g_{\mathcal{G}}^* - g^*\|_\infty^2 + h^2 + \delta, \tag{S.15}
\end{aligned}$$

where the last inequality follows from the definition of  $g_{q^*}$ , i.e.,  $\|g_{q^*} - \hat{g}_c\|_\infty < \delta$  and  $g_{q^*} \in \mathcal{G}$ , and the unbiasedness of the score function, i.e.,  $\mathbb{E}_{D_n} \left\{ \frac{f^{(1)}(Y_i - g^*(\mathbf{X}_i))}{f(Y_i - g^*(\mathbf{X}_i))} \right\} = 0$ .

For  $I_{2,2}$ , let  $L_{\mathcal{K}}(g_1, g, \mathbf{Z}_i) = \log \left( \frac{1}{n_1} \sum_{j=1}^{n_1} \mathcal{K}_h(g_1(\mathbf{X}_j) - Y_j, g(\mathbf{X}_i) - Y_i) \right) - \log \tilde{f}_{g_1, h}(Y_i -$

$g(\mathbf{X}_i)$ ). Then, we have that

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} (L_{\mathcal{K}}(\hat{g}_c, \hat{g}_c, \mathbf{Z}_i)) \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left[ \log \left( \frac{1}{n} \sum_{j=1}^n \mathcal{K}_h(\hat{g}_c(\mathbf{X}_j) - Y_j, \hat{g}_c(\mathbf{X}_i) - Y_i) \right) - \log \tilde{f}_{\hat{g}_c, h}(Y_i - \hat{g}_c(\mathbf{X}_i)) \right] \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left[ \log \left( 1 + \frac{\frac{1}{n} \sum_{j=1}^n \mathcal{K}_h(\hat{g}_c(\mathbf{X}_j) - Y_j, \hat{g}_c(\mathbf{X}_i) - Y_i) - \tilde{f}_{\hat{g}_c, h}(Y_i - \hat{g}_c(\mathbf{X}_i))}{\tilde{f}_{\hat{g}_c, h}(Y_i - \hat{g}_c(\mathbf{X}_i))} \right) \right] \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \left[ \frac{\frac{1}{n} \sum_{j=1}^n \mathcal{K}_h(\hat{g}_c(\mathbf{X}_j) - Y_j, \hat{g}_c(\mathbf{X}_i) - Y_i) - \tilde{f}_{\hat{g}_c, h}(Y_i - \hat{g}_c(\mathbf{X}_i))}{\tilde{f}_{\hat{g}_c, h}(Y_i - \hat{g}_c(\mathbf{X}_i))} (1 + o_p(1)) \right] \\
&\leq \left( \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \frac{\frac{1}{n} \sum_{j=1}^n \mathcal{K}_h(g^*(\mathbf{X}_j) - Y_j, g^*(\mathbf{X}_i) - Y_i) - \tilde{f}_{g^*, h}(Y_i - g^*(\mathbf{X}_i))}{\tilde{f}_{g^*, h}(Y_i - g^*(\mathbf{X}_i))} \right| \right. \\
&\quad \left. + \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \frac{-\frac{1}{n} \sum_{j=1}^n \dot{\mathcal{K}}_h(g^*(\mathbf{X}_j) - Y_j, g^*(\mathbf{X}_i) - Y_i) - \dot{\tilde{f}}_{g^*, h}(Y_i - g^*(\mathbf{X}_i))}{\dot{\tilde{f}}_{g^*, h}(Y_i - g^*(\mathbf{X}_i))} \right| \|\hat{g}_c - g^*\|_{\infty} \right) \\
&\quad \times (1 + o_p(1)).
\end{aligned}$$

Write  $\mathbf{Z}_i = Y_i - g^*(\mathbf{X}_i)$ . Let  $U(\mathbf{Z}_j, \mathbf{Z}_i) = \frac{\mathcal{K}_h(\mathbf{Z}_j, \mathbf{Z}_i) - \tilde{f}_{g^*, h}(\mathbf{Z}_i)}{\tilde{f}_{g^*, h}(\mathbf{Z}_i)} \times 0.5 + \frac{\mathcal{K}_h(\mathbf{Z}_i, \mathbf{Z}_j) - \tilde{f}_{g^*, h}(\mathbf{Z}_j)}{\tilde{f}_{g^*, h}(\mathbf{Z}_j)} \times 0.5$ . Clearly, using the U-statistics theory (Theorems 1 and 3 of Chapter 1.3 in Lee

2019), we have

$$\begin{aligned}
\tilde{U}(\mathbf{Z}_j) &= \mathbb{E}_{\mathbf{Z}_i} [U(\mathbf{Z}_j, \mathbf{Z}_i | \mathbf{Z}_j)] \\
&= \mathbb{E}_{\mathbf{Z}_i} \left[ \frac{\mathcal{K}_h(\mathbf{Z}_i, \mathbf{Z}_j) - \tilde{f}_{g^*,h}(\mathbf{Z}_i)}{\tilde{f}_{g^*,h}(\mathbf{Z}_i)} \times 0.5 \right] \\
&= 0.5 \times \int \frac{\mathcal{K}_h(\mathbf{Z}_i, \mathbf{Z}_j) - \tilde{f}_{g^*,h}(\mathbf{Z}_i)}{\tilde{f}_{g^*,h}(\mathbf{Z}_i)} f(\mathbf{Z}_i) d\mathbf{Z}_i \\
&= 0.5 \times \int \frac{\mathcal{K}_h(\mathbf{Z}_i, \mathbf{Z}_j)}{\tilde{f}_{g^*,h}(\mathbf{Z}_i)} f(\mathbf{Z}_i) d\mathbf{Z}_i - 0.5 \\
&= 0.5 \times \int \frac{1}{h} \frac{\mathcal{K}(\frac{\mathbf{Z}_i - \mathbf{Z}_j}{h})}{\tilde{f}_{g^*,h}(\mathbf{Z}_i)} f(\mathbf{Z}_i) d\mathbf{Z}_i - 0.5 \\
&= 0.5 \times \int \frac{\mathcal{K}(t)}{\tilde{f}_{g^*,h}(th + \mathbf{Z}_j)} f(th + \mathbf{Z}_j) dt - 0.5 \\
&= 0.5 \times \int \mathcal{K}(t) \frac{f(th + \mathbf{Z}_j)}{\tilde{f}_{g^*,h}(th + \mathbf{Z}_j)} dt - 0.5 \\
&= 0.5 \times \int \mathcal{K}(t) \left( \tilde{f}_{g^*,h}(th + \mathbf{Z}_j) - f^{(1)}(th + \mathbf{Z}_j) L_1(g^* - g^*) \right. \\
&\quad \left. - f^{(2)}(th + \mathbf{Z}_j) 0.5 \{ L_2(g^* - g^*) + U_2 h^2 \} \right) / \left( \tilde{f}_{g^*,h}(th + \mathbf{Z}_j) \right) dt - 0.5 \\
&= 0.5 \times \int \mathcal{K}(t) \left( 1 - \frac{f^{(2)}}{\tilde{f}_{g^*,h}}(th + \mathbf{Z}_j) 0.5 U_2 h^2 \right) dt - 0.5 \\
&= - \left( \frac{f^{(2)}}{\tilde{f}_{g^*,h}}(\mathbf{Z}_j) 0.5 U_2 h^2 (1 + o(1)) \right),
\end{aligned}$$

which leads to

$$\text{Var}(\tilde{U}(\mathbf{Z}_j)) \leq h^4.$$

Then, we have

$$\mathbb{E}_{D_n} \left| \frac{1}{n} \sum_{i=1}^n \frac{1}{n} \sum_{j=1}^n U(\mathbf{Z}_j, \mathbf{Z}_i) \right| \leq n^{-1/2} h^2.$$



Note that

$$\begin{aligned}
& \mathbb{E} \left\{ \frac{1}{n} \sum_{j=1}^n \mathcal{K}_h^{(1)}(Y_j - g_1(\mathbf{X}_j), y - g(\mathbf{x})) \right\} \\
&= \int \frac{1}{h^2} K^{(1)} \left( \frac{Y - g_1(\mathbf{X}) - y + g(\mathbf{x})}{h} \right) f(Y - g^*(\mathbf{X})) f_{\mathbf{x}}(\mathbf{X}) dY d\mathbf{X} \\
&= \int \frac{1}{h} K^{(1)}(t) f(g_1(\mathbf{X}) - g^*(\mathbf{X}) + y - g(\mathbf{x}) + th) f_{\mathbf{x}}(\mathbf{X}) dt d\mathbf{X} \\
&= \int K^{(1)}(t) \left\{ h^{-1} f(g_1(\mathbf{X}) - g^*(\mathbf{X}) + y - g(\mathbf{x})) + f^{(1)}(g_1(\mathbf{X}) - g^*(\mathbf{X}) + y - g(\mathbf{x})) t \right. \\
&\quad \left. + f^{(2)}(g_1(\mathbf{X}) - g^*(\mathbf{X}) + y - g(\mathbf{x})) \frac{1}{2} t^2 h + o(h) \right\} f_{\mathbf{x}}(\mathbf{X}) dt d\mathbf{X} \\
&= h^{-1} \int K^{(1)}(t) dt \left( f(y - g(\mathbf{x})) + f^{(1)}(y - g(\mathbf{x})) L_1(g_1 - g^*) \right. \\
&\quad \left. + f^{(2)}(y - g(\mathbf{x})) \frac{1}{2} L_2(g_1 - g^*) + O(L_3(g_1 - g^*)) \right) \\
&+ \int K^{(1)}(t) t dt \left( f^{(1)}(y - g(\mathbf{x})) + f^{(2)}(y - g(\mathbf{x})) L_1(g_1 - g^*) + O(L_2(g_1 - g^*)) \right) \\
&+ \frac{\int K^{(1)}(t) t^2 dt}{2} h \left( f^{(2)}(y - g(\mathbf{x})) + O(L_1(g_1 - g^*)) \right) (1 + o(1)).
\end{aligned}$$

Then, after similar calculations, we have that

$$\begin{aligned}
& \left| \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{D_n} \frac{-\frac{1}{n} \sum_{j=1}^n \dot{\mathcal{K}}_h(g^*(\mathbf{X}_j) - Y_j, g^*(\mathbf{X}_i) - Y_i) - \dot{f}_{g^*,h}(Y_i - g^*(\mathbf{X}_i))}{\dot{f}_{g^*,h}(Y_i - g^*(\mathbf{X}_i))} \right| \\
&\quad \times \|\hat{g}_c - g^*\|_{\infty} \leq n^{-1/2} h \|\hat{g}_c - g^*\|_{\infty}.
\end{aligned}$$

Thus, we can obtain that

$$|I_{2,2}| \leq \frac{1}{\sqrt{n}} (h^2 + h \|\hat{g}_c - g^*\|_{\infty}). \quad (\text{S.16})$$

Combining inequalities (S.15) and (S.16), it follows from  $h \rightarrow 0$ ,

$$|I_2| \leq h^2 + n^{-1/2} h^2 + \delta + \|g_{\mathcal{G}}^* - g^*\|_{\infty}^2.$$

Similarly, we can obtain that  $|I_3| = O(|I_2|)$ . Thus, equation (S.8) holds from condition (C3a).

- We then show that

$$\|\tilde{g}(\mathbf{x}) - \hat{g}_c(\mathbf{x})\|_\infty = O_p(n^{-1/2}).$$

In particular,

–

$$\begin{aligned} \tilde{g}(\cdot) &= \hat{g}_c(\cdot) + \int \hat{g}(\mathbf{x}) f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} - \mathbb{E}(Y) - \frac{1}{n} \sum_{i=1}^n \hat{g}(\mathbf{X}_i) + \frac{1}{n} \sum_{i=1}^n Y_i \\ &= \hat{g}_c(\cdot) + \left( \frac{1}{n} \sum_{i=1}^n Y_i - \mathbb{E}(Y) \right) + \left( \int \hat{g}(\mathbf{x}) f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} - \frac{1}{n} \sum_{i=1}^n \hat{g}(\mathbf{X}_i) \right) \\ &:= \hat{g}_c(\cdot) + I_1 + I_2. \end{aligned}$$

- Given  $\mathbb{E}Y^2 < \infty$ , it follows from central limit theorem that

$$I_1 := \frac{1}{n} \sum_{i=1}^n Y_i - \mathbb{E}(Y) = O_p(n^{-1/2}). \quad (\text{S.17})$$

- Define  $F_{n,\mathbf{x}}(\mathbf{x}) = n^{-1} \sum_{i=1}^n I(\mathbf{X}_i \leq \mathbf{x})$ . Let  $F_{\mathbf{x}}$  be the distribution function of  $\mathbf{x}$ . Based on the empirical measure, we have

$$\begin{aligned} |I_2| &:= \left| \frac{1}{n} \sum_{i=1}^n \hat{g}(\mathbf{X}_i) - \int \hat{g}(\mathbf{x}) f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \right| \\ &= \left| \int \hat{g}(\mathbf{x}) dF_n(\mathbf{x}) - \int \hat{g}(\mathbf{x}) dF(\mathbf{x}) \right| \\ &= \left| \int \hat{g}(\mathbf{x}) d(F_{n,\mathbf{x}}(\mathbf{x}) - F_{\mathbf{x}}(\mathbf{x})) \right| \\ &\leq \frac{\mathcal{B}}{n^{1/2}}, \end{aligned} \quad (\text{S.18})$$

where the last inequality follows from the upper bounded condition of  $\hat{g}$  and the  $n^{1/2}$  convergence rate of empirical distribution. Similar conclusions can also be found in Bickel and Ritov (2003). They claimed that functionals of a nonparametric regression function is able to be estimated efficiently, i.e., at  $n^{1/2}$  convergence rate (Last paragraph on Page 1050).

- Last, we show that

$$\mathbb{E}(\mathcal{R}(\tilde{g}) - \mathcal{R}(g^*)) \simeq \mathbb{E}(\mathcal{R}(\hat{g}_c) - \mathcal{R}(g^*)) + O(n^{-1}).$$

Particularly,

$$\begin{aligned} \mathbb{E}(\mathcal{R}(\tilde{g}) - \mathcal{R}(g^*)) &= \mathbb{E}(\mathcal{R}(\hat{g}_c + I_1 + I_2) - \mathcal{R}(g^*)) \\ &= \mathbb{E}(\mathcal{R}(\hat{g}_c) - \mathcal{R}(g^*)) + \mathbb{E}(\mathcal{R}(\hat{g}_c + I_1 + I_2) - \mathcal{R}(\hat{g}_c)) \\ &\simeq \mathbb{E}(\mathcal{R}(\hat{g}_c) - \mathcal{R}(g^*)) + O_p(n^{-1}), \end{aligned}$$

where the last approximation follows from

$$\mathbb{E}(\mathcal{R}(g_1) - \mathcal{R}(g_2)) \simeq \|g_1 - g_2\|_\infty^2,$$

and  $\|I_1 + I_2\|_\infty^2 = O_p(n^{-1})$  by (S.17) and (S.18).

Then, Theorem 2 follows.

## S.2.4 Proof of Proposition 1

According to Theorem 2, the variance of  $\tilde{g}$  equals to the variance of  $\hat{g}_{oracle}$  (Substituting  $f$  with  $\hat{f}$  and rescaling does not introduce additional variance). Thus, it suffices to show that for any asymptotically unbiased  $\tilde{g}$ ,

$$\text{Var}(\tilde{g}) \geq \text{Var}(\hat{g}_{oracle}) = \text{Var}(\tilde{g})(1 + o(1)).$$

According to Stoica and Marzetta (2001), when  $J$  is singular, any biased estimator  $\check{g}(\mathbf{x}) := g(\mathbf{x}; \check{\boldsymbol{\theta}})$  with finite variance  $C$  must satisfy (inequality (16) in Stoica and Marzetta (2001))

$$C := \text{Var}(\check{g}(\mathbf{x})) \geq n H J^\dagger H^\top, \quad (\text{S.19})$$

where

$$H = \frac{1}{n} \left( \frac{\partial (\mathbb{E}\{\check{g}(\mathbf{x})\} - g(\mathbf{x}; \boldsymbol{\theta}))}{\partial \boldsymbol{\theta}^\top} + \frac{\partial g(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} + \mathbb{E}(\check{g}(\mathbf{x})) o(1) \right) |_{\boldsymbol{\theta}=\boldsymbol{\theta}^*}$$

and  $B^\dagger$  is the Moore-Penrose generalized inverse of a matrix  $B$ .

Let's consider the asymptotically unbiased estimator  $g(x, \check{\boldsymbol{\theta}})$  obtained by minimizing a loss function  $\ell$  under the framework of FNN. That is, minimizing the loss  $\frac{1}{n} \sum_{i=1}^n \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}))$ . Using  $\frac{1}{n} \sum_{i=1}^n \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} |_{\boldsymbol{\theta}=\check{\boldsymbol{\theta}}} = 0$ , we can obtain that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}))}{\partial g(\mathbf{x}; \boldsymbol{\theta})} |_{\boldsymbol{\theta}=\check{\boldsymbol{\theta}}} &= \frac{1}{n} \sum_{i=1}^n \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \frac{\partial \boldsymbol{\theta}}{\partial g(\mathbf{x}; \boldsymbol{\theta})} |_{\boldsymbol{\theta}=\check{\boldsymbol{\theta}}} \\ &= \left( \frac{1}{n} \sum_{i=1}^n \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right) \frac{\partial \boldsymbol{\theta}}{\partial g(\mathbf{x}; \boldsymbol{\theta})} |_{\boldsymbol{\theta}=\check{\boldsymbol{\theta}}} \\ &= 0. \end{aligned}$$

Then it follows that

$$\begin{aligned} 0 &= \frac{1}{n} \sum_{i=1}^n \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}))}{\partial g(\mathbf{x}; \boldsymbol{\theta})} |_{\boldsymbol{\theta}=\check{\boldsymbol{\theta}}} \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)} + \frac{1}{n} \sum_{i=1}^n \frac{\partial^2 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^2} (g(\mathbf{x}; \check{\boldsymbol{\theta}}) - g(\mathbf{x}; \boldsymbol{\theta}^*)) \\ &\quad + \frac{1}{2} \frac{1}{n} \sum_{i=1}^n \frac{\partial^3 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^3} (g(\mathbf{x}; \check{\boldsymbol{\theta}}) - g(\mathbf{x}; \boldsymbol{\theta}^*))^2 (1 + o_p(1)), \end{aligned} \quad (\text{S.20})$$

which leads to

$$\begin{aligned}
& g(\mathbf{x}; \check{\boldsymbol{\theta}}) \\
&= g(\mathbf{x}; \boldsymbol{\theta}^*) - \left( \frac{1}{n} \sum_{i=1}^n \frac{\partial^2 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^2} \right)^{-1} \frac{1}{n} \sum_{i=1}^n \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)} \\
&- \left\{ \frac{1}{n} \sum_{i=1}^n \frac{\partial^2 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^2} \right\}^{-1} \frac{1}{2} \frac{1}{n} \sum_{i=1}^n \frac{\partial^3 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^3} (g(\mathbf{x}; \check{\boldsymbol{\theta}}) - g(\mathbf{x}; \boldsymbol{\theta}^*))^2 (1 + o_p(1)).
\end{aligned} \tag{S.21}$$

Note that the approximation bias of  $g(\mathbf{x}; \boldsymbol{\theta})$  to  $g(\mathbf{x})$ , expressed by  $g(\mathbf{x}; \boldsymbol{\theta}^*) - g(\mathbf{x})$ , is determined by the function class  $g(\mathbf{x})$  belongs to and the framework of FNN, and tends to zero if  $g(\cdot; \check{\boldsymbol{\theta}})$  achieves nearly minimax optimal rate (Jiao et al. 2023). Therefore, an asymptotically unbiased FNN estimator achieving nearly minimax optimal rate actually satisfies  $Eg(\mathbf{x}; \check{\boldsymbol{\theta}}) - g(\mathbf{x}, \boldsymbol{\theta}^*) = o(1)$ . Then it follows that  $\mathbb{E} \left( \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)} \right) = o(1)$  from (S.21). Furthermore, from (S.21), we have

$$\begin{aligned}
& \mathbb{E} (g(\mathbf{x}; \check{\boldsymbol{\theta}}) - g(\mathbf{x}; \boldsymbol{\theta}^*))^2 \\
&= \left\{ \mathbb{E} \left( \frac{\partial^2 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^2} \right) \right\}^{-1} \left\{ \mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)} \right)^2 \right\} \\
&\times \left\{ \mathbb{E} \left( \frac{\partial^2 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^2} \right) \right\}^{-1} (1 + o(1)) \\
&= \left[ \mathbb{E} \left( \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)} \right) \left\{ \mathbb{E} \left( \frac{\partial^2 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^2} \right) \right\}^{-1} \right]^2 \\
&+ \frac{1}{n} \left\{ \mathbb{E} \left( \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)} \right) - \mathbb{E} \left( \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)} \right) \right\}^2 \\
&\times \left\{ \mathbb{E} \left( \frac{\partial^2 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^2} \right) \right\}^{-2} (1 + o(1)).
\end{aligned} \tag{S.22}$$

Substituting (S.22) into (S.21), it then follows from  $\mathbb{E} \left( \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)} \right) = o(1)$  that

$$\begin{aligned}
Bias &:= \mathbb{E}\{g(\mathbf{x}; \check{\boldsymbol{\theta}})\} - g(\mathbf{x}; \boldsymbol{\theta}^*) \\
&= -\mathbb{E} \left( \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)} \right) \left\{ \mathbb{E} \left( \frac{\partial^2 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^2} \right) \right\}^{-1} \\
&\quad - \left\{ \mathbb{E} \left( \frac{\partial^2 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^2} \right) \right\}^{-1} \frac{1}{2} \mathbb{E} \left( \frac{\partial^3 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^3} \right) \\
&\quad \times \mathbb{E} (g(\mathbf{x}; \check{\boldsymbol{\theta}}) - g(\mathbf{x}; \boldsymbol{\theta}^*))^2 (1 + o(1)) \\
&= -\mathbb{E} \left( \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)} \right) \left\{ \mathbb{E} \left( \frac{\partial^2 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^2} \right) \right\}^{-1} \\
&\quad - \frac{1}{2n} \left\{ \mathbb{E} \left( \frac{\partial^2 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^2} \right) \right\}^{-1} \mathbb{E} \left( \frac{\partial^3 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^3} \right) \\
&\quad \times \left\{ \mathbb{E} \left( \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)} \right) - \mathbb{E} \left( \frac{\partial \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)} \right) \right\}^2 \\
&\quad \times \left\{ \mathbb{E} \left( \frac{\partial^2 \ell(Y_i, \mathbf{X}_i; g(\cdot; \boldsymbol{\theta}^*))}{\partial g(\mathbf{x}; \boldsymbol{\theta}^*)^2} \right) \right\}^{-2} (1 + o(1)).
\end{aligned}$$

Thus, if the loss function has bounded third order derivatives with respect to  $g$ , it holds that

$$\frac{\partial Bias}{\partial \boldsymbol{\theta}^*} = o(1).$$

Then it follows from (S.19) that

$$C = \text{Var}(\check{g}(\mathbf{x})) \geq n \left( \frac{1}{n} \frac{\partial g(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \right) J^\dagger \left( \frac{1}{n} \frac{\partial g(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right) |_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} := V.$$

We next show that  $V > 0$ . Utilizing the eigenvector/eigenvalue representation of  $J$ , we have

$$J = U \Lambda U^\top = \begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{bmatrix} \Lambda_1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} U_1^\top \\ U_2^\top \end{bmatrix},$$

where  $U$  is orthonormal,  $\Lambda_1 \in \mathbb{R}^{r \times r}$  is diagonal and positive definite. Then, it follows from the definition of Moore-Penrose inverse  $J^\dagger$  that

$$J^\dagger = U_1 \Lambda_1^{-1} U_1^\top,$$

Thus, combining with  $\partial g(\mathbf{x}; \boldsymbol{\theta}) / \partial \boldsymbol{\theta} |_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \neq 0$ , we have

$$V = n \left( \frac{1}{n} \frac{\partial g(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right) U_1 \Lambda_1^{-1} U_1^\top \left( \frac{1}{n} \frac{\partial g(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right) > 0. \quad (\text{S.23})$$

On the other hand, note that, the oracle estimator

$$\hat{g}_{oracle} = \arg \min_{g \in \mathcal{G}} \left\{ \frac{1}{n} \sum_{i=1}^n (-\log f(Y_i - g(\mathbf{X}_i))) \right\}.$$

Then, it follows from

$$\begin{aligned} 0 &\equiv \left( \frac{1}{n} \sum_{i=1}^n \frac{\partial \log f(Y_i - g(\mathbf{X}_i; \boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right) \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_{oracle}} \\ &= \left( \frac{1}{n} \sum_{i=1}^n \frac{\partial \log f(Y_i - g(\mathbf{X}_i; \boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \\ &\quad + \left[ \left( \frac{1}{n} \sum_{i=1}^n \frac{\partial^2 \log f(Y_i - g(\mathbf{X}_i; \boldsymbol{\theta}))}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right] (\hat{\boldsymbol{\theta}}_{oracle} - \boldsymbol{\theta}^*) + O_p(\|\hat{\boldsymbol{\theta}}_{oracle} - \boldsymbol{\theta}^*\|_2^2), \end{aligned}$$

that

$$\begin{aligned} &\left( \frac{1}{n} \frac{\partial g(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right) J^\dagger \left( \frac{1}{n} \sum_{i=1}^n \frac{\partial \log f(Y_i - g(\mathbf{X}_i; \boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \\ &= - \left( \frac{1}{n} \frac{\partial g(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right) J^\dagger \left[ \left( \frac{1}{n} \sum_{i=1}^n \frac{\partial^2 \log f(Y_i - g(\mathbf{X}_i; \boldsymbol{\theta}))}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right] (\hat{\boldsymbol{\theta}}_{oracle} - \boldsymbol{\theta}^*) (1 + o_p(1)) \\ &= \left( \frac{1}{n} \frac{\partial g(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right) J^\dagger J (\hat{\boldsymbol{\theta}}_{oracle} - \boldsymbol{\theta}^*) (1 + o_p(1)) \\ &= \left( \frac{1}{n} \frac{\partial g(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right) (\hat{\boldsymbol{\theta}}_{oracle} - \boldsymbol{\theta}^*) (1 + o_p(1)), \end{aligned}$$

where the last equality follows from equation (18) in Stoica and Marzetta (2001). According to Theorem 1, the oracle estimator  $\|\hat{g}_{oracle} - g(\mathbf{x}; \boldsymbol{\theta}^*)\|_\infty^2 \rightarrow 0$  under some conditions on the function class and network width and depth (Corollary 1). Thus, we can obtain that

$$\begin{aligned} & g(\mathbf{x}; \hat{\boldsymbol{\theta}}_{oracle}) \\ = & g(\mathbf{x}; \boldsymbol{\theta}^*) + \frac{\partial g(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} (\hat{\boldsymbol{\theta}}_{oracle} - \boldsymbol{\theta}^*) + O_p(\|\hat{\boldsymbol{\theta}}_{oracle} - \boldsymbol{\theta}^*\|_2^2) \\ = & g(\mathbf{x}; \boldsymbol{\theta}^*) + n \left( \frac{1}{n} \frac{\partial g(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}^\top} \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} \right) J^\dagger \left( \frac{1}{n} \sum_{i=1}^n \frac{\partial \log f(Y_i - g(\mathbf{X}_i; \boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} \right) \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*} + O_p(\|\hat{\boldsymbol{\theta}}_{oracle} - \boldsymbol{\theta}^*\|_2^2). \end{aligned}$$

It then follows that the variance of  $\hat{g}_{oracle}$  achieves the bound  $V$ , i.e.,

$$\text{Var}(\hat{g}_{oracle}) = V(1 + o(1)).$$

Thus, we can obtain that for any asymptotically unbiased FNN based estimator,

$$\text{Var}(\tilde{g}) \geq \text{Var}(\hat{g}_{oracle}).$$

Then Proposition 1 follows.

### S.3 Lemmas

**Lemma S.1.** (*Approximation error, (Theorem 3.3 in Jiao et al. 2023)*)

Given Hölder smooth functions  $g^* \in \mathcal{H}_\beta([0, 1]^d, B_0)$ , for any  $D \in \mathbb{N}^+$  and  $W \in \mathbb{N}^+$ , there exists a function  $g_{\mathcal{G}}^*$  implemented by a ReLU feedforward neural network with width  $\mathcal{W} = 38(\lfloor \beta \rfloor + 1)^2 d^{\lfloor \beta \rfloor + 1} W \lceil \log_2(8W) \rceil$  and depth  $\mathcal{D} = 21(\lfloor \beta \rfloor + 1)^2 D \lceil \log_2(8D) \rceil$  such that

$$|g^* - g_{\mathcal{G}}^*| \leq 18B_0(\lfloor \beta \rfloor + 1)^2 d^{\lfloor \beta \rfloor + \max\{\beta, 1\}/2} (WD)^{-2\beta/d},$$

for all  $x \in [0, 1]^d \setminus \Omega([0, 1]^d, K, \delta)$  where

$$\Omega([0, 1]^d, K, \delta) = \bigcup_{i=1}^d \{x = [x_1, \dots, x_d]^T : x_i \in \bigcup_{k=1}^{K-1} (k/K - \delta, k/K)\},$$



with  $K = \lfloor WD \rfloor$  and  $\delta$  is an arbitrary number in  $(0, 1/3K]$ .

**Lemma S.2.** (Bounding the covering number, (Theorem 12.2 in Anthony et al. 1999) and (Theorems 3 and 7 in Bartlett et al. 2019))

Let ReLU feedforward neural network  $\mathcal{G}$  be a set of real functions from a domain  $\mathcal{X}$  to the bounded interval  $[0, \mathcal{B}]$ . There exists a universal constant  $C$  such that the following holds. Given any  $\mathcal{D}, \mathcal{S}$  with  $\mathcal{S} > C\mathcal{D} > C^2$ , there exists network class  $\mathcal{G}$  with  $\leq \mathcal{D}$  layers and  $\leq \mathcal{S}$  parameters with VC-dimension  $\geq \mathcal{S}\mathcal{D} \log(\mathcal{S}/\mathcal{D})/C$  and given  $\delta > 0$

$$\mathcal{N}_{2n}(\delta, \|\cdot\|_\infty, \mathcal{G}|_{\mathcal{X}}) \leq \sum_{i=1}^{\mathcal{S}\mathcal{D} \log(\mathcal{S}/\mathcal{D})/C} \binom{2n}{i} \left(\frac{\mathcal{B}}{\delta}\right) = O(\mathcal{S}\mathcal{D} \log(\mathcal{S})).$$

**Lemma S.3.** (Intrinsic dimensionality, Shen (2020))

Let  $f$  be a continuous functions on  $[0, 1]^d$  and  $\mathcal{M} \subseteq [0, 1]^d$  be a compact  $d_{\mathcal{M}}$ -dimensional Riemannian submanifold. For any  $D \in \mathbb{N}^+$ ,  $W \in \mathbb{N}^+$ ,  $\epsilon \in (0, 1)$  and  $\delta \in (0, 1)$ , there exists a function  $g_\phi^*$  implemented by a ReLU feedforward neural network with width  $\mathcal{W} = 3^{d_\delta} \max\{d_\delta \lfloor W^{1/d_\delta} \rfloor, W + 1\}$  and depth  $\mathcal{D} = 12D + 14 + 2d_\delta$  such that

$$\|g^*(x) - g_\phi^*(x)\| \leq 2\omega_g\left(\frac{2\epsilon}{1-\delta} \sqrt{\frac{d}{d_\delta}} + 2\epsilon\right) + 19\sqrt{d}\omega_g\left(\frac{2\sqrt{d}}{(1-\delta)\sqrt{d_\delta}} W^{-2/d_\delta} D^{-2/d_\delta}\right)$$

for any  $x \in \mathcal{M}_\epsilon$ , where  $\mathcal{M} := \{x \in [0, 1]^d : \inf\{\|x - y\|_2 : y \in \mathcal{M}\} \leq \epsilon\}$ , for  $\epsilon \in (0, 1)$  and  $d_\delta$  is an integer such that  $d_{\mathcal{M}} \leq d_\delta \leq d$ .

## S.4 Results in numerical studies

		LS	LAD	Huber	Cauchy	Turkey	EML
$n = 256, d = 100$							
Normal	PE	<b>1.1169</b>	1.3431	1.2402	1.3904	1.3421	1.1181
	SD	<b>0.0473</b>	0.0591	0.0518	0.0653	0.0583	0.0479
Mixture Gaussian	PE	10.2935	9.8321	10.1538	9.5132	9.3231	<b>8.2394</b>
	SD	0.8363	0.7035	0.8111	0.6683	0.6051	<b>0.5582</b>
Student-t	PE	10.7079	9.1759	10.2351	8.8393	9.9832	<b>8.0394</b>
	SD	3.9652	2.7836	2.8532	2.7563	2.8165	<b>2.3942</b>
Heteroscedasticit	PE	18.6772	17.5086	18.4008	17.1440	16.6032	<b>14.5063</b>
	SD	1.4171	1.2793	1.2759	1.1432	1.1234	<b>0.9011</b>
$n = 256, d = 500$							
Normal	PE	<b>1.1175</b>	1.3967	1.3531	1.4973	1.4196	1.1193
	SD	<b>0.0487</b>	0.0610	0.0523	0.0671	0.0610	0.0491
Mixture Gaussian	PE	10.5336	9.8851	10.4314	9.6366	9.9842	<b>8.3159</b>
	SD	0.8630	0.7246	0.8616	0.7502	0.8166	<b>0.6142</b>
Student-t	PE	11.3261	9.5716	10.8394	9.3742	10.0432	<b>8.4261</b>
	SD	4.1326	2.8154	2.9862	2.7738	2.9143	<b>2.5142</b>
Heteroscedasticit	PE	18.7571	17.6223	18.4314	17.3645	16.7491	<b>14.5154</b>
	SD	1.4224	1.2831	1.3238	1.2243	1.2032	<b>0.9132</b>
$n = 1024, d = 500$							
Normal	PE	1.1112	1.2548	1.1942	1.3613	1.2332	<b>1.0972</b>
	SD	<b>0.0461</b>	0.0568	0.0503	0.0609	0.0515	0.0466
Mixture Gaussian	PE	10.1931	9.806	9.9978	9.3856	9.3031	<b>8.2185</b>
	SD	0.7147	0.6712	0.6990	0.5343	0.5343	<b>0.5154</b>
Student-t	PE	9.5032	8.1356	8.7142	8.0797	8.3921	<b>7.1264</b>
	SD	3.924	2.5163	2.7168	2.4766	2.6032	<b>2.3012</b>
Heteroscedasticit	PE	17.9693	17.1379	17.6407	16.8798	16.5712	<b>14.4745</b>
	SD	1.2401	1.1957	1.2394	1.1396	1.1073	<b>0.8879</b>

Table S1: The mean and standard deviation of prediction error of  $g_5$  when using six methods for four error distributions with sample sizes  $n = 256, 1024$ , testing sample size  $t = 2048$  and input dimensions  $d = 100, 500$ .

		LS	LAD	Huber	Cauchy	Turkey	EML
$n = 256, d = 200$							
Normal	PE	<b>1.1603</b>	1.3703	1.1909	1.3904	1.3450	1.1632
	SD	0.0592	0.0911	0.0685	0.0879	0.0905	<b>0.0464</b>
Mixture Gaussian	PE	10.4945	10.1405	10.4869	9.7103	9.5147	<b>8.3187</b>
	SD	0.9168	0.8779	0.9067	0.8695	0.8453	<b>0.7168</b>
Student-t	PE	13.9393	11.7430	13.3536	11.7075	11.273	<b>10.0289</b>
	SD	10.1406	9.9717	10.0178	9.1084	9.0983	<b>8.0537</b>
Heteroscedasticit	PE	18.8234	17.8664	18.6526	16.8545	16.6812	<b>14.4935</b>
	SD	1.2655	1.13	1.1842	1.0594	1.0032	<b>0.8785</b>
$n = 256, d = 600$							
Normal	PE	<b>1.1835</b>	1.5035	1.2898	1.4856	1.5304	1.1912
	SD	0.0651	0.1508	0.0732	0.0988	0.1768	<b>0.0537</b>
Mixture Gaussian	PE	11.4058	10.2008	10.534	9.7404	9.9544	<b>8.3311</b>
	SD	1.1526	0.9003	0.9527	0.8813	0.8651	<b>0.7556</b>
Student-t	PE	14.6317	12.6555	13.7298	12.3261	13.3261	<b>11.1408</b>
	SD	11.9669	10.8725	11.8775	10.852	11.8499	<b>9.6821</b>
Heteroscedasticit	PE	20.1975	18.3254	19.0039	17.2722	16.7207	<b>14.5214</b>
	SD	1.4803	1.2841	1.5119	1.1595	1.0528	<b>0.8937</b>
$n = 1024, d = 600$							
Normal	PE	1.1593	1.3111	1.2323	1.3772	1.3347	<b>1.1082</b>
	SD	0.0531	0.0842	0.0548	0.0907	0.0895	<b>0.0416</b>
Mixture Gaussian	PE	10.3174	9.5994	10.2685	9.4418	9.3905	<b>8.2696</b>
	SD	0.8825	0.863	0.8779	0.8407	0.7982	<b>0.6617</b>
Student-t	PE	12.8441	11.1071	12.6051	11.0201	12.4215	<b>9.4071</b>
	SD	9.8297	8.8399	9.7424	8.8095	9.7735	<b>7.7424</b>
Heteroscedasticit	PE	18.1690	17.4382	18.2861	16.6768	16.5519	<b>14.3442</b>
	SD	1.1637	1.0451	1.0982	1.0028	0.8494	<b>0.78</b>

Table S2: The mean and standard deviation of prediction error of  $g_{10}$  when using six methods for four error distributions with sample sizes  $n = 256, 1024$ , testing sample size  $t = 2048$  and input dimensions  $d = 200, 600$ .

## References

- Anthony, M., Bartlett, P. L., Bartlett, P. L., et al. (1999). *Neural network learning: Theoretical foundations*, volume 9. cambridge university press Cambridge.
- Bartlett, P. L., Harvey, N., Liaw, C., and Mehrabian, A. (2019). Nearly-tight vc-dimension and pseudodimension bounds for piecewise linear neural networks. *The Journal of Machine Learning Research*, 20(1):2285–2301.
- Bickel, P. J. and Ritov, Y. (2003). Nonparametric estimators which can be” plugged-in”. *The Annals of Statistics*, 31(4):1033–1053.
- Bobkov, S., Chistyakov, G., and Götze, F. (2024). Strictly subgaussian probability distributions. *Electronic Journal of Probability*, 29:1–28.
- Jiao, Y., Shen, G., Lin, Y., and Huang, J. (2023). Deep nonparametric regression on approximate manifolds: Nonasymptotic error bounds with polynomial prefactors. *The Annals of Statistics*, 51(2):691–716.
- Lee, A. J. (2019). *U-statistics: Theory and Practice*. Routledge.
- Shen, Z. (2020). Deep network approximation characterized by number of neurons. *Communications in Computational Physics*, 28(5).
- Stoica, P. and Marzetta, T. L. (2001). Parameter estimation problems with singular information matrices. *IEEE Transactions on Signal Processing*, 49(1):87–90.