

# Supplementary Material of “Sensitivity Analysis under the $f$ -Sensitivity Models: A Distributional Robustness Perspective”

Ying Jin<sup>\*1</sup>, Zhimei Ren<sup>\*1</sup>, and Zhengyuan Zhou<sup>3</sup>

<sup>1</sup>Department of Statistics and Data Science, the Wharton School, University of Pennsylvania

<sup>3</sup>Stern School of Business, New York University

## A Deferred details and discussions

### A.1 Discussion on related and future directions

In this work, we propose a new sensitivity model based on the  $f$ -divergence that characterizes the *average* effect of confounders on selection bias. Under the  $f$ -sensitivity model, we offer a scheme for the estimation and inference on the counterfactual and the ATE. We close the paper by a discussion on possible extensions.

**Sensitivity analysis.** In this paper, we have focused on conducting inference on the counterfactuals and treatment effects under the  $(f, \rho)$ -selection condition, with a prescribed confounding parameter  $\rho$ . Based on this, we can make robust causal conclusions and conduct sensitivity analysis by inverting the confidence intervals as follows. Suppose the goal is to detect if there is a nonzero ATE; we can consider an increasing sequence of  $\rho$ , and construct a level  $1 - \alpha$  confidence interval  $\widehat{C}(\rho)$  for the ATE using the method introduced in this paper at each value of  $\rho$ ; then let  $\widehat{\rho}$  be the smallest  $\rho$  such that  $\widehat{C}(\rho)$  contains zero. We can interpret the results as “either there is a nonzero ATE, or there is a confounder as large as  $\widehat{\rho}$  to explain away the observed treatment effects”.

More rigorously, let  $\rho^*$  denote the true confounding level and suppose the constructed confidence intervals  $\widehat{C}(\rho)$  are nested in  $\rho$ , in the sense that  $\widehat{C}(\rho_1) \subset \widehat{C}(\rho_2)$  for any  $\rho_1 \leq \rho_2$ . We then have

$$\limsup_{n \rightarrow \infty} \mathbb{P}(\text{ATE} = 0, \rho^* < \widehat{\rho}) \leq \limsup_{n \rightarrow \infty} \mathbb{P}(\text{ATE} \notin \widehat{C}(\rho^*)) \leq \alpha,$$

if  $\widehat{C}(\rho^*)$  is an asymptotically valid confidence interval for the ATE. In words, when the ATE is indeed zero,  $\widehat{\rho}$  is an asymptotic level- $(1 - \alpha)$  confidence lower bound for  $\rho^*$ . Similar to the case of [Jin et al. \(2023\)](#), here only point-wise validity is necessary, i.e., we only need our CIs to be asymptotically valid for each fixed ground truth of  $\rho$ . Finally, we note that the monotonicity of the confidence intervals is satisfied with a reasonable estimation procedure; in addition, given a sequence of preliminary confidence intervals  $\widehat{C}(\rho)$  that are not necessarily nested, one can enforce monotonicity while maintaining validity via

$$\widehat{C}^{\text{nest}}(\rho) := \cup_{\rho' \leq \rho} \widehat{C}(\rho').$$

It is straightforward to see that  $\widehat{C}^{\text{nest}}(\rho)$ 's are nested, and they are still asymptotically valid since  $\widehat{C}(\rho) \subseteq \widehat{C}^{\text{nest}}(\rho)$  for every  $\rho \geq 0$ .

---

\*Author names listed alphabetically.

**Calibrating the sensitivity parameter.** It is impossible to estimate the confounding strength  $\rho$  without additional information. However, in practice, a reasonable approach for guessing the magnitude of confounding is to pretend an observed covariate is a missing confounder, and use the leave-one-out approach to estimate the resulting deviation as  $\rho$ . This can be done for both conditional and marginal  $f$ -sensitivity models. To be specific, let  $U$  be the one observed variable that is left out, and  $X$  be the remaining covariates, and suppose data are  $\{X_i, Y_i, T_i\}_{i=1}^n$ . One can fit propensity score models  $\widehat{e}(x, u)$  for  $\mathbb{P}(T = 1 | X = x, U = u)$  and  $\widehat{e}_0(x)$  for  $\mathbb{P}(T = 1 | X = x)$ , and then set

$$\widehat{\text{OR}}(x, u) = \frac{\widehat{e}_0(x)/(1 - \widehat{e}_0(x))}{\widehat{e}(x, u)/(1 - \widehat{e}(x, u))}.$$

Then, for the conditional  $f$ -sensitivity model (consider  $\mathcal{M}_1^{f, \rho}$ ), we can set

$$\widehat{\rho} = \sup_{T_i=1} \widehat{f}(X_i),$$

where  $\widehat{f}(x)$  fits  $\mathbb{E}[f(\widehat{\text{OR}}(x, U)) | X = x, T = 1]$ . For the marginal  $f$ -sensitivity model (consider  $\mathcal{M}_{1, \text{mgn}}^{f, \rho}$ ), we can set

$$\widehat{\rho} = \frac{1}{n_1} \sum_{T_i=1} f(\widehat{\text{OR}}(X_i, U_i)).$$

Note that the accuracy of these estimators depends on the estimation quality of  $\widehat{e}_0$ ,  $\widehat{e}$ , and  $\widehat{f}$ . Analyzing the properties of these estimators are beyond the scope of this paper, and we suggest using these estimators as an exploratory analysis. In addition, since  $\widehat{\rho}$  for the marginal  $f$ -sensitivity model requires less functional estimation, it might be easier to calibrate than the conditional  $f$ -sensitivity model.

**Conditional average treatment effect (CATE).** The methodology proposed in this paper also provides bounds on CATE under the conditional  $(f, \rho)$ -selection condition (however, the situation with the marginal  $(f, \rho)$ -selection condition is more complicated). Concretely, following the proof of Propositions 2 and 3, a lower bound for  $\mathbb{E}[Y(1) | X = x, T = 1]$  is given by the optimal value of

$$\begin{aligned} & \min_{L(x, y) \geq 0 \text{ measurable}} \mathbb{E}[Y(1)L(x, Y(1)) | X = x, T = 1] \\ & \text{s.t. } \mathbb{E}[L(x, Y(1)) | T = 1, X = x], \\ & \mathbb{E}\left[f\left(\frac{L(x, Y(1))}{r_{1,0}(x)}\right) | X = x, T = 1\right] \leq \rho. \end{aligned}$$

The dual form of the above optimization problem is

$$\sup_{\alpha \geq 0, \eta \in \mathbb{R}} -r_{1,0}(x) \cdot \mathbb{E}\left[\alpha f^*\left(-\frac{Y(1) + \eta}{\alpha}\right) + \eta + \alpha \rho \mid X = x, T = 1\right]. \quad (1)$$

Note that the optimizer  $(\alpha^*(x), \eta^*(x))$  defined in (3) is exactly the optimizer of (1). In fact,  $\widehat{\mu}(x) := \widehat{r}^{(j)}(x)\widehat{h}^{(j)}(x)$  where  $\widehat{r}^{(j)}, \widehat{h}^{(j)}$  are defined in Algorithm 1 is an estimator for the optimal objective in (1). These quantities are repeatedly estimated on distinct folds of data as intermediate steps of our procedure. While such sample splitting does not compromise the efficiency of inference due to the final averaging step, how to efficiently estimate these CATE bound functions with statistical guarantees might call for distinct considerations from ours. We leave this for future investigation.

## A.2 Discussions on Assumption 2 on sieve estimation

We provide additional discussion on Assumption 2 for sieve estimators in the context of  $(X, Y(1)) | T = 1$ . In particular, we first justify the smoothness of the optimizers when the conditional distributions are sufficiently smooth. We then verify the technical conditions for two choices of  $f$ -divergences: KL-divergence and  $\chi^2$ -divergence. Then we discuss some considerations of relaxing the conditions with implementations in practice.

**Smoothness of the optimizers.** We first provide some justifications for assuming the optimizers are continuously differentiable. By the strong convexity of  $f$ , its conjugate  $f^*$  is continuous, hence without loss of generality we always assume the differentiation and expectation are exchangeable. We also assume the conjugate  $f^*$  is sufficiently smooth, which is the case for many popular choices of  $f$ -divergence. As we discussed in Proposition 4, under mild conditions, the optimizers  $(\alpha^*(x), \eta^*(x))$  lies in the interior of  $\mathbb{R}^+ \times \mathbb{R}$ . The optimizers are thus the solutions to

$$(\alpha^*(x), \eta^*(x)) = \operatorname{argzero} \nabla_{a,b} \left\{ a \mathbb{E} \left[ f^* \left( \frac{Y(1)+b}{-a} \right) \mid X = x, T = 1 \right] + b + a\rho \right\},$$

where the right-hand side takes the form

$$F(a, b, x) := \mathbb{E} [g(Y(1), a, b) \mid X = x, T = 1] \in \mathbb{R}^2$$

for some differentiable or smooth function  $g$  decided by  $f^*$  and its derivative  $(f^*)'$ . Thus  $F(a, b, x)$  is smooth in  $(a, b)$  when  $f^*$  is sufficiently smooth. Now let us assume the conditional distribution  $\mathbb{P}_{Y(1) \mid X=x, T=1}$  is smooth; for example, for some  $h \in \mathcal{X}$ ,  $\mathbb{P}_{Y(1) \mid X=x+th, T=1} = \mathbb{P}_{Y(1) \mid X=x, T=1} + t \cdot \mathbb{P}_h$  for some measure  $\mathbb{P}_h$  on  $\mathcal{Y}$ ; and similar for higher-order expansions. This is a reasonable assumption if we are willing to assume that the conditional distributions of  $Y(1)$  are close for similar covariates. Concretely, such condition holds when  $Y(1) \mid X = x, T = 1$  is a normal distribution with homoskedastic noise and a smooth mean function, or heteroskedastic noise with a smooth mean function and smooth standard deviation function, etc. When the conditional distributions are smooth in  $x$ , the function  $F(a, b, x)$  is also smooth in  $x$  by the linearity of conditional expectation. Finally, if the derivatives with respect to  $a, b$  is always invertible (which is the case under mild conditions for the examples we discuss shortly) and smooth, invoking the Implicit Function Theorem (Rudin et al., 1976), the minimizer can be smooth in  $x$ .

**KL-divergence.** A popular choice for the function  $f$  is  $f(x) = x \log x$ , which leads to the KL-divergence (Kullback and Leibler, 1951). The dual function in this case is  $f^*(y) = e^{y-1}$ , and the loss function becomes

$$\ell(\theta, x, y) = \alpha(x) e^{\frac{y+\eta(x)}{-\alpha(x)}-1} + \eta(x) + \alpha(x)\rho.$$

The conditional expectation is

$$\mathbb{E} [\ell((a, b), x, Y(1)) \mid X = x, T = 1] = a \mathbb{E} \left[ e^{\frac{Y(1)+b}{-a}-1} \mid X = x, T = 1 \right] + b + a\rho.$$

We first look at the strong convexity assumption. The conditional expectation is twice differentiable, with

$$\begin{aligned} \nabla_a^2 \mathbb{E} [\ell((a, b), x, Y(1)) \mid X = x, T = 1] &= \frac{1}{a^3} \mathbb{E} \left[ (Y(1) + b)^2 e^{\frac{Y(1)+b}{-a}-1} \mid X = x, T = 1 \right], \\ \nabla_b^2 \mathbb{E} [\ell((a, b), x, Y(1)) \mid X = x, T = 1] &= \frac{1}{a} \mathbb{E} \left[ e^{\frac{Y(1)+b}{-a}-1} \mid X = x, T = 1 \right], \\ \nabla_{a,b}^2 \mathbb{E} [\ell((a, b), x, Y(1)) \mid X = x, T = 1] &= -\frac{1}{a^2} \mathbb{E} \left[ (Y(1) + b) e^{\frac{Y(1)+b}{-a}-1} \mid X = x, T = 1 \right]. \end{aligned}$$

Therefore, a simple calculation shows that as long as  $Y(1)$  is not deterministic at  $(\alpha^*(x), \beta^*(x))$ , the Hessian matrix is non-singular. Also, if the underlying distribution  $\mathbb{P}_{Y(1) \mid X=x, T=1}$  is continuous in  $x$ , the above derivatives, hence the eigenvalues of the Hessian matrix is continuous; since  $\mathcal{X}$  is compact, there exists a positive uniform lower bound for the smallest eigenvalue of the Hessian matrix, leading to strong convexity.

We then consider the continuity condition  $|\ell(\theta, x, y) - \ell(\theta^*, x, y)| \leq \bar{\ell}(x, y) \|\theta(x) - \theta^*(x)\|_2$  for  $\|\theta(x) - \theta^*(x)\|_2 < \epsilon$  for some sufficiently small  $\epsilon > 0$ , where  $\|\theta(x) - \theta^*(x)\|_2$  is the Euclidean norm, and  $\sup_{x \in \mathcal{X}} \mathbb{E} [\bar{\ell}(x, Y(1))^2 \mid X = x, T = 1] < M$  for some constant  $M > 0$ . By Taylor expansion, we have

$$\ell(\theta, x, y) - \ell(\theta^*, x, y) = \nabla_{\theta} \ell(\tilde{\theta}, x, y) (\theta^*(x) - \theta(x)),$$

where  $\tilde{\theta}(x)$  lies between  $\theta(x)$  and  $\theta^*(x)$ . We note that  $\nabla_{\theta}$  is also a smooth function of  $\theta$ , and the gradient is uniform bounded for  $\theta(x)$  within a neighborhood of  $\theta^*(x)$  in terms of Euclidean  $L_2$ -norm. In particular,

$$\frac{\partial}{\partial a} \ell((a, b), x, y) = \left(1 + \frac{y+b}{a}\right) e^{\frac{y+b}{-a}-1} + \rho, \quad \frac{\partial}{\partial b} \ell((a, b), x, y) = 1 - e^{\frac{y+b}{-a}-1}.$$

For any  $\|(a, b) - \theta^*(x)\|_2 \leq \epsilon$  for sufficiently small  $\epsilon$ , we can take  $\bar{\ell}(x, y)$  as the uniform upper bound of the Euclidean norm of the gradient, which has finite second moment if  $Y(1)$  is not too heavy-tailed.

Finally, the last condition is that there exists a constant  $C_1$  such that  $\mathbb{E}[\ell(\theta, X, Y(1)) - \ell(\theta^*, X, Y(1)) | T = 1] \leq C_1 \|\theta - \theta^*\|_{L_2(\mathbb{P}_{\cdot | T=1})}^2$  when  $\theta \in \Lambda_c^p(\mathcal{X}) \times \Lambda_c^p(\mathcal{X})$  and  $\|\theta - \theta^*\|_{L_2(\mathbb{P}_{\cdot | T=1})}$  is sufficiently small. Similar to arguments in the proof of Theorem 1, sufficiently small  $\|\theta - \theta^*\|_{L_2(\mathbb{P}_{\cdot | T=1})}$  implies sufficiently small  $\|\theta - \theta^*\|_\infty$  for this function class. Therefore, we can consider  $\theta \in \Lambda_c^p(\mathcal{X}) \times \Lambda_c^p(\mathcal{X})$  such that  $\|\theta - \theta^*\|_\infty$  is sufficiently small. With a Taylor expansion of the conditional expectation of the risk at  $(\alpha^*(x), \eta^*(x))$ , we have

$$\begin{aligned} & \mathbb{E}[\ell(\theta, x, Y(1)) | X = x, T = 1] - \mathbb{E}[\ell(\theta^*, x, Y(1)) | X = x, T = 1] \\ &= 1/2 \cdot \nabla_{\tilde{\theta}}^2 \mathbb{E}[\ell(\tilde{\theta}, x, Y(1)) | X = x, T = 1][\theta(x) - \theta^*(x), \theta(x) - \theta^*(x)] \end{aligned}$$

since the gradient is zero, where  $\tilde{\theta}(x)$  lies between  $\theta(x)$  and  $\theta^*(x)$ . Previous derivations have shown that the Hessian is continuous; also, by the compactness of  $\mathcal{X}$  and continuity of  $\theta^*(x)$ , there is a uniform lower bound  $\epsilon > 0$  for  $\alpha^*(x)$ . Thus, when  $\|\theta - \theta^*\|_\infty$  is sufficiently small, the Hessian is also bounded. Again by the compactness of  $\mathcal{X}$ , this bound can be taken to be uniform for  $x \in \mathcal{X}$ , which leads to the desired condition.

**$\chi^2$ -divergence.** Another popular choice is  $f(x) = (x-1)^2$ , so that  $f^*(y) = \frac{1}{4}((y+2)_+^2 - 1)$ . The conjugate function is a quadratic function on  $[-2, \infty)$  and zero on  $(-\infty, -2]$ , with continuous gradient  $(f^*)'(y) = (\frac{y}{2} + 1)_+$ , and second-order derivative  $(f^*)''(y) = \frac{1}{2} \mathbf{1}\{y > -2\}$ ; the latter is almost-everywhere (under Lebesgue measure) except  $y = 2$ . We now proceed to verify the conditions. The loss function is

$$\ell(\theta, x, y) = \frac{\alpha(x)}{4} \left[ \left( \frac{y + \eta(x)}{-\alpha(x)} + 1 \right)_+^2 - 1 \right] + \eta(x) + \alpha(x)\rho.$$

Assuming  $Y(1)$  does not have point measure, the differentiation and expectation are exchangeable, and

$$\begin{aligned} \nabla_a^2 \mathbb{E}[\ell((a, b), x, Y(1)) | X = x, T = 1] &= \frac{1}{2a^3} \mathbb{E} \left[ (Y(1) + b)^2 \mathbf{1}\left\{ \frac{Y(1)+b}{-a} > -2 \right\} \middle| X = x, T = 1 \right], \\ \nabla_b^2 \mathbb{E}[\ell((a, b), x, Y(1)) | X = x, T = 1] &= \frac{1}{2a} \mathbb{E} \left[ \mathbf{1}\left\{ \frac{Y(1)+b}{-a} > -2 \right\} \middle| X = x, T = 1 \right], \\ \nabla_{a,b}^2 \mathbb{E}[\ell((a, b), x, Y(1)) | X = x, T = 1] &= -\frac{1}{2a^2} \mathbb{E} \left[ (Y(1) + b) \mathbf{1}\left\{ \frac{Y(1)+b}{-a} > -2 \right\} \middle| X = x, T = 1 \right]. \end{aligned}$$

Also, the gradient is given by

$$\begin{aligned} \nabla_a \mathbb{E}[\ell((a, b), x, Y(1)) | X = x, T = 1] &= \mathbb{E} \left[ \left( \frac{Y(1)+b}{-2a} + 1 \right)_+^2 - 1 + \frac{Y(1)+b}{a} \left( \frac{Y(1)+b}{-2a} + 1 \right)_+ \middle| X = x, T = 1 \right] + \rho, \\ \nabla_b \mathbb{E}[\ell((a, b), x, Y(1)) | X = x, T = 1] &= -\mathbb{E} \left[ \left( \frac{Y(1)+b}{-2a} + 1 \right)_+ \middle| X = x, T = 1 \right] + 1, \end{aligned}$$

which are both zero at  $(a, b) = (\alpha^*(x), \eta^*(x))$ . The form of the loss function implies that  $\alpha^*(x) > 0$  for almost all  $x$ ; hence there is a uniform lower bound  $\epsilon > 0$  by the compactness of  $\mathcal{X}$ . By Cauchy-Schwarz inequality, the Hessian at  $(a, b) = (\alpha^*(x), \eta^*(x))$  is positive  $\mathbb{P}(\frac{Y(1)+\eta^*(x)}{-\alpha^*(x)} > -2 | X = x, T = 1) = 0$  or  $(Y(1) + \eta^*(x) - c(x)) \mathbf{1}\{\frac{Y(1)+\eta^*(x)}{-a} > -2\} = 0$  almost surely for some  $c(x) \in \mathbb{R}$ . By the optimality condition, the former is impossible, and the latter is also impossible if  $Y(1)$  is not deterministic conditional on  $X = x$ . Thus, as long as  $Y(1) | X = x$  is not deterministic for almost all  $x$ , the Hessian is positive definite for all  $x \in \mathcal{X}$ . By compactness of  $\mathcal{X}$  and the continuity, we know that the minimal eigenvalue of the Hessian is uniformly lower bounded away from zero, hence the strong convexity follows.

The other two conditions are easy to verify in this case: the conjugate function  $f^*$  is a truncation of a quadratic function. Since truncation is a contraction map, these results hold easily by the uniform boundedness of second-order derivatives. We've thus verified the conditions in Assumption 2 for  $\chi^2$ -divergence.

**Practical considerations.** In practice, we might search for  $(\alpha^*(x), \eta^*(x))$  within the function classes with a bounded range of coefficients in the two examples we give, leading to a compact function space. This is typically assumed in the contexts of  $M$ -estimators and sieve estimators (Van der Vaart, 2000; Geer et al., 2000; Chen and Shen, 1998; Chen, 2007). In this case, the regularity conditions are easier to verify given the uniform boundedness. The function space still provides finer and finer approximation to the targets if the bounded range enlarges properly with  $n$ .

### A.3 Intuitions for the debiasing technique

In this part, we explain some high-level intuitions for the debiasing technique introduced in Section 3.3 which enables fast-rate inference in Section 4.2 with slow-rate estimators.

In view of the duality form in Proposition 3, with an estimator  $(\hat{\alpha}(\cdot), \hat{\eta}(\cdot))$  for the optimizer  $(\alpha^*, \eta^*)$ , a natural plug-in estimator for the optimal objective is

$$\hat{\mu}_{1,0}^{+, \text{naive}} = \frac{1}{n_1} \sum_{T_i=1} \hat{r}(X_i) \hat{H}(X_i, Y_i(1)), \quad \hat{H}(x, y) = \hat{\alpha}(x) f^* \left( \frac{y + \hat{\eta}(x)}{-\hat{\alpha}(x)} \right) + \hat{\eta}(x) + \hat{\alpha}(x) \rho.$$

However, the convergence rate for nonparametric estimation of  $\hat{\alpha}(\cdot)$ ,  $\hat{\eta}(\cdot)$ , and  $\hat{r}(\cdot)$  is typically slower than  $n^{-1/2}$ . That is, we usually have  $\|\hat{H} - H\|_{L_2} = O_P(n^{-\beta})$  for  $\beta < 1/2$ , where  $H(x, y)$  is the function with true  $(\alpha^*, \eta^*)$ . Such a bias will lead to a bias in  $\hat{\mu}_{1,0}^{+, \text{naive}}$  which is larger than the typical  $O(n^{-1/2})$  standard deviation, and further violate the Wald-type inference based on  $\hat{\mu}_{1,0}^{+, \text{naive}}$  since it is no longer asymptotically root- $n$  normal.

We will utilize covariates  $X_i$  for  $T_i = 0$  to remove the first-order bias in  $\hat{\mu}_{1,0}^{+, \text{naive}}$ . Assume also that  $\hat{h}$  is an external estimator of  $\bar{h}(x) := \mathbb{E}[\hat{H}(x, Y(1)) | X = x, T = 1]$ . For simplicity, assume  $\hat{r}$ ,  $\hat{\alpha}$  and  $\hat{\eta}$  are given and independent of all the data, and so is  $\hat{H}$ . The bias in  $\hat{\mu}_{1,0}^{+, \text{naive}}$  is

$$\hat{\mu}_{1,0}^{+, \text{naive}} - \hat{\mathbb{E}}_1[rH] = \hat{\mathbb{E}}_1[\hat{r} \cdot (\hat{H} - \bar{h})] + \hat{\mathbb{E}}_1[\hat{r}\bar{h}] - \hat{\mathbb{E}}_1[r_{1,0}H],$$

where  $\hat{\mathbb{E}}_1[\cdot]$  is the empirical mean in the treated group. The first term  $\hat{\mathbb{E}}_1[r(\hat{H} - \bar{h})]$  is mean zero (unbiased) since  $\bar{h}$  is the true conditional mean function of  $\hat{H}(X, Y(1))$ , and thus negligible. Due to the optimality of  $\alpha^*$  and  $\eta^*$ , as well as the convexity of  $H$ , the first-order Taylor expansion of  $H$  implies  $\mathbb{E}_1[r(\hat{H} - H)]$  is of a higher order than the estimation bias, hence also negligible. The remaining bias is thus (we use  $\approx$  to denote equivalence up to higher-order error)

$$\hat{\mu}_{1,0}^{+, \text{naive}} - \hat{\mathbb{E}}_1[r_{1,0}H] \approx \hat{\mathbb{E}}_1[\hat{r}\bar{h}] - \hat{\mathbb{E}}_1[r_{1,0}\hat{H}].$$

Next, the bias  $\hat{\mathbb{E}}_1[\hat{r}\bar{h}] - \hat{\mathbb{E}}_1[r_{1,0}\hat{H}]$  is removed by subtracting  $\hat{\mathbb{E}}_1[\hat{r}\hat{h}] - \hat{\mathbb{E}}_0[\hat{h}]$  from  $\hat{\mu}_{1,0}^{+, \text{naive}}$ , leading to our estimator (ignoring cross-fitting). The idea is that

$$\begin{aligned} \hat{\mathbb{E}}_1[\hat{r}\hat{h}] - \hat{\mathbb{E}}_0[\hat{h}] &\approx \hat{\mathbb{E}}_1[\hat{r}\bar{h}] + \hat{\mathbb{E}}_1[\hat{r}(\hat{h} - \bar{h})] - \hat{\mathbb{E}}_0[\hat{h}] \\ &\approx \hat{\mathbb{E}}_1[\hat{r}\bar{h}] + \hat{\mathbb{E}}_1[r_{1,0}(\hat{h} - \bar{h})] - \hat{\mathbb{E}}_1[r_{1,0}\hat{h}] \\ &\approx \hat{\mathbb{E}}_1[\hat{r}\bar{h}] - \hat{\mathbb{E}}_1[r_{1,0}\bar{h}] \approx \hat{\mathbb{E}}_1[\hat{r}\bar{h}] - \hat{\mathbb{E}}_1[r_{1,0}\hat{H}]. \end{aligned}$$

In the second line, we use the fact that  $\hat{\mathbb{E}}_1[(\hat{r} - r_{1,0})(\hat{h} - \bar{h})]$  is of a higher order. Also,  $r_{1,0}(x)$  is the covariate shift between treated and control groups, so that for any fixed function  $f(x)$ , it holds that  $\mathbb{E}[r_{1,0}(X)f(X) | T = 1] = \mathbb{E}[f(X) | T = 0]$ .

Finally, as all nuisance component estimators have to be computed from data, we use the cross-fitting technique to decouple the estimators and the data they are applied to.

### A.4 Estimators for bounds on counterfactual means

In this section, we summarize the application of the procedure in Section 3.3 to estimate other lower and upper bounds on counterfactual means.

- (a) **Upper bound of  $\mathbb{E}[Y(1) | T = 0]$ :** Let  $-\hat{\mu}_{1,0}^+$  be the estimator obtained from the procedure in Section 3.3 with  $-Y(1)$  replacing  $Y(1)$ . Then  $\sqrt{n}(\hat{\mu}_{1,0}^+ - \mu_{1,0}^+) \rightsquigarrow N(0, \text{Var}(\phi_{1,+}(X, Y, T)))$ , with influence function

$$\phi_{1,+}(X_i, Y_i, T_i) = \frac{T_i}{p_1} r_{1,0}(X_i) [H_{1,+}(X_i, -Y_i(1)) - h_{1,+}(X_i)] + \frac{1 - T_i}{p_0} h_{1,+}(X_i),$$

where  $H_{1,+}(x, y) = \alpha_{1,+}^*(x) f^* \left( \frac{y + \eta_{1,+}^*(x)}{-\alpha_{1,+}^*(x)} \right) + \eta_{1,+}^*(x) + \alpha_{1,+}^*(x) \rho$  with  $(\alpha_{1,+}^*(x), \eta_{1,+}^*(x))$  being the minimizer of  $\mathbb{E}[\alpha f^* \left( \frac{-Y(1) + \eta}{-\alpha} \right) + \eta + \alpha \rho | X = x, T = 1]$ , and  $h_{1,+}(x) = \mathbb{E}[H_{1,+}(X, -Y(1)) | X = x, T = 1]$ .

- (b) **Lower bound of  $\mathbb{E}[Y(0) | T = 1]$ :** Let  $\widehat{\mu}_{0,1}^-$  be the estimator obtained from the procedure in Section 3.3 switching the role of treated and control groups. Then  $\sqrt{n}(\widehat{\mu}_{0,1}^- - \mu_{0,1}^-) \rightsquigarrow N(0, \text{Var}(\phi_{0,-}(X, Y, T)))$  with influence function

$$\phi_{0,-}(X_i, Y_i, T_i) = \frac{1 - T_i}{p_0} r_{0,1}(X_i) [H_{0,-}(X, Y(0)) - h_{0,-}(X_i)] + \frac{T_i}{p_1} h_{0,-}(X_i).$$

Here  $H_{0,-}(x, y) = \alpha_{0,-}^*(x) f^*\left(\frac{y + \eta_{0,-}^*(x)}{-\alpha_{0,-}^*(x)}\right) + \eta_{0,-}^*(x) + \alpha_{0,-}^*(x)\rho$ , and  $(\alpha_{0,-}^*(x), \eta_{0,-}^*(x))$  is the minimizer of  $\mathbb{E}[\alpha f^*\left(\frac{Y(0) + \eta}{-\alpha}\right) + \eta + \alpha\rho | X = x, T = 0]$ , and  $h_{0,-}(x) = \mathbb{E}[H_{0,-}(X, Y(0)) | X = x, T = 0]$ .

- (c) **Upper bound of  $\mathbb{E}[Y(0) | T = 1]$ :** Let  $-\widehat{\mu}_{0,1}^+$  be the estimator obtained from the procedure in Section 3.3 switching the role of treated and control groups and replacing  $Y(0)$  with  $-Y(0)$ . Then  $\sqrt{n}(\widehat{\mu}_{0,1}^+ - \mu_{0,1}^+) \rightsquigarrow N(0, \text{Var}(\phi_{0,+}(X, Y, T)))$  with influence function

$$\phi_{0,+}(X_i, Y_i, T_i) = \frac{1 - T_i}{p_0} r_{0,1}(X_i) [H_{0,+}(X_i, -Y_i(0)) - h_{0,+}(X_i)] + \frac{T_i}{p_1} h_{0,+}(X_i),$$

where  $H_{0,+}(x, y) = \alpha_{0,+}^*(x) f^*\left(\frac{y + \eta_{0,+}^*(x)}{-\alpha_{0,+}^*(x)}\right) + \eta_{0,+}^*(x) + \alpha_{0,+}^*(x)\rho$  with  $(\alpha_{0,+}^*(x), \eta_{0,+}^*(x))$  being the minimizer of  $\mathbb{E}[\alpha f^*\left(\frac{-Y(0) + \eta}{-\alpha}\right) + \eta + \alpha\rho | X = x, T = 0]$ , and  $h_{0,+}(x) = \mathbb{E}[H_{0,+}(X, -Y(0)) | X = x, T = 0]$ .

## B Sensitivity analysis under the marginal $f$ -sensitivity model

This section collects the identification, estimation and inference for treatment effects under the marginal  $f$ -sensitivity model introduced in Section 2.4.

### B.1 Identification results

Parallel to Proposition 1, below we show that the marginal  $(f, \rho)$ -selection condition implies a bounded average  $f$ -divergence between  $\mathbb{P}_{X, Y(1) | T=0}$  and  $\mathbb{P}_{X, Y(1) | T=1}$ , averaged over  $\mathbb{P}_{X | T=1}$ . Switching the role of treated and control groups leads to a similar result for distribution shifts under  $\mathcal{M}_{0, \text{mgn}}^{f, \rho}$ . Its proof is also similar to that of Proposition 1, and we defer it to Appendix D.1.

**Proposition B.1.** *Let  $\mathbb{P}$  be the true unknown joint distribution over  $(X, U, T, Y(0), Y(1))$ , which induces  $\mathbb{P}_{X, Y, T}$ , the joint distribution of  $(X, Y, T)$  where  $Y = Y(T) = TY(1) + (1 - T)Y(0)$ . Let  $\mathcal{P}$  be the set of all distributions over  $(X, Y(0), Y(1), U, T)$  and  $\mathcal{P}_{X, Y(1)}$  be the set of all distributions over  $(X, Y(1))$ . Define  $\mathcal{Q}_{1,0}^m$  as the ambiguity set of all counterfactual distributions that agree with the observables and satisfy marginal  $(f, \rho)$ -selection for the treated:*

$$\mathcal{Q}_{1,0}^m = \{Q_{X, Y(1) | T=0} : Q \in \mathcal{P}, Q_{X, Y, T} = \mathbb{P}_{X, Y, T}, Q \in \mathcal{M}_{1, \text{mgn}}^{f, \rho}\}$$

Then  $\mathbb{P}_{X, Y(1) | T=0} \in \mathcal{Q}_{1,0}^m$ , and a sharp characterization of  $\mathcal{Q}_{1,0}^m$  is given by

$$\mathcal{Q}_{1,0}^m = \left\{ Q : \frac{dQ_X}{d\mathbb{P}_{X | T=1}^{\text{obs}}}(x) = r_{1,0}(x), \int_x D_f(Q_{Y | X=x} \| \mathbb{P}_{Y | X=x, T=1}^{\text{obs}}) d\mathbb{P}^{\text{obs}}(x | T = 1) \leq \rho \right\}.$$

### B.2 Convex optimization characterization

Parallel to Proposition 2, below we characterize the partial identification bounds for the counterfactual mean  $\mathbb{E}[Y(1) | T = 0]$  under  $\mathcal{M}_{1, \text{mgn}}^{f, \rho}$  via an optimization program.

**Proposition B.2.** *Let  $\mu_{1,0}^{m+}$  (resp.  $\mu_{0,1}^{m-}$ ) be the optimal value of the convex optimization problem:*

$$\begin{aligned} & \min(\text{resp. max})_{L(x,y) \text{ measurable}} \mathbb{E}[Y(1)L(X, Y(1)) | T = 1] \\ & \text{s.t. } \mathbb{E}[L(X, Y(t)) | X = x, T = 1] = r_{1,0}(x), \quad \forall x \\ & \mathbb{E}[f(L(X, Y(1))/r_{1,0}(X)) | T = 1] \leq \rho, \end{aligned} \tag{2}$$

Then  $\mu_{1,0}^{m,-} \leq \mathbb{E}[Y(1) | T = 0] \leq \mu_{1,0}^{m,+}$  if  $\mathbb{P} \in \mathcal{M}_{1,\text{mgn}}^{f,\rho}$ .

The dual formulation is provided in the following proposition as preparation for our estimation procedure, whose proof is in Appendix D.2.

**Proposition B.3.** *The solution to (2) is given by*

$$\mu_{1,0}^{m,-} = - \inf_{\alpha \geq 0, \eta(\cdot)} \mathbb{E} \left[ \alpha f^* \left( \frac{r_{1,0}(X)}{-\alpha} (Y(1) + \eta(X)) \right) + \eta(X) r_{1,0}(X) + \alpha \rho \mid T = 1 \right], \quad (3)$$

where  $f^*(s) = \sup_{t \geq 0} \{st - f(t)\}$  is the conjugate function of  $f$ .

### B.3 Sensitivity analysis with stable risk minimizer

We then discuss estimation and inference strategies under the marginal  $f$ -sensitivity model. At a high level, we will still follow a plug-in-then-debias strategy. However, the specific debiasing technique and the assumptions required for consistent estimation and valid inference will differ from previous ones. We first introduce the procedure.

#### B.3.1 The procedure

We start by dividing both the treated group and the control group into four equal-sized folds and use the notation  $\mathcal{I}_1^{(j)}$  and  $\mathcal{I}_0^{(j)}$  to denote the resulting folds,  $j = 1, \dots, 4$ . In view of (3), for any  $a \in \mathbb{R}^+$ ,  $b, r \in \mathbb{R}$ , we denote the generic risk function as

$$m(x, y, a, b, r) = a f^* \left( \frac{r}{-a} (y + \eta) \right) + br + a\rho.$$

For each fold  $j = 1, \dots, 4$ , we use  $\mathcal{I}_1^{(j+1)}$  and  $\mathcal{I}_0^{(j+1)}$  to obtain an estimator  $\hat{r}^{(j)}$  for  $r_{1,0}$ ; as  $r_{1,0}$  is decided by the propensity score  $e(x) = \mathbb{P}(T = 1 | X = x)$ , there are many methods in the literature that can be applied. We then use  $\mathcal{I}_1^{(j+2)}$  to obtain the empirical risk minimizer

$$(\hat{\alpha}^{(j)}, \hat{\eta}^{(j)}) = \arg \min_{\alpha \in \mathbb{R}, \eta \in \Theta_n} \frac{1}{|\mathcal{I}_1^{(j+2)}|} \sum_{i \in \mathcal{I}_1^{(j+2)}} m(X_i, Y_i, \alpha, \eta(X_i), \hat{r}^{(j)}(X_i)),$$

where  $\Theta_n$  is some function class that might grow with  $n$ . This estimation procedure is the same for the preceding naive estimator, whose convergence has been discussed in Section B.3.2.

We assume throughout that  $f^*$  is differentiable, and define

$$H^{(j)}(x, y) = \frac{\partial f^*}{\partial r} (x, y, \hat{\alpha}^{(j)}, \hat{\eta}^{(j)}(x), \hat{r}^{(j)}(x)).$$

We then use the fold  $\mathcal{I}_1^{(j+3)}$  to obtain an estimator  $\hat{h}^{(j)}(\cdot)$  for the conditional mean function

$$\bar{h}^{(j)}(x) = \mathbb{E}[H^{(j)}(x, Y(1)) | X = x, T = 1],$$

viewing  $H^{(j)}$  as fixed. Finally, we define the adjusted estimator on folds  $\mathcal{I}_1^{(j)}$  and  $\mathcal{I}_0^{(j)}$  as

$$\hat{\mu}_{1,0}^{m,(j)} = \frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} \left\{ m(X_i, Y_i, \hat{\alpha}^{(j)}, \hat{\eta}^{(j)}(X_i), \hat{r}^{(j)}(X_i)) - \hat{\alpha}^{(j)} \hat{r}^{(j)}(X_i) \hat{h}^{(j)}(X_i) \right\} + \frac{\hat{\alpha}^{(j)}}{|\mathcal{I}_0^{(j)}|} \sum_{i \in \mathcal{I}_0^{(j)}} \hat{h}^{(j)}(X_i).$$

The final estimator is the average of the estimator obtained from each fold:

$$\hat{\mu}_{1,0}^{m,-} = -\frac{1}{4} \sum_{j=1}^4 \hat{\mu}_{1,0}^{m,(j)}. \quad (4)$$

### B.3.2 Estimation of $\hat{\alpha}$ and $\hat{\eta}$

The estimation of  $\hat{\alpha} \in \mathbb{R}^+$  and  $\hat{\eta}: \mathcal{X} \rightarrow \mathbb{R}$  can be done in a similar fashion as Section 3.4. To be specific, we still employ the sieve estimator for (2), and let  $\Theta_n = \mathbb{R}^+ \times \Theta_n^\eta$  where  $\Theta_n^\eta$  is the function class in Examples 2 or 3. We define the population risk minimizer given the estimator  $\hat{r}$  as

$$(\bar{\alpha}, \bar{\eta}) = \arg \min_{\alpha \in \mathbb{R}, \eta(\cdot)} \mathbb{E} \left[ m(X, Y(1), \alpha, \eta(X), \hat{r}(X)) \mid T = 1, \mathcal{I}_1^{(1)} \cup \mathcal{I}_0^{(1)} \right].$$

We impose the following assumption on regularity conditions of the risk function under  $\hat{r}$ . Here we view  $\bar{m}(x, y, a, b) = m(x, y, a, b, \hat{r}(x))$  as a fixed function; for simplicity, we write  $\bar{m}(x, y, \theta) = \bar{m}(x, y, \alpha, \eta(x))$  for  $\theta = (\alpha, \eta(\cdot))$  and denote  $\bar{\theta} = (\bar{\alpha}, \bar{\eta}(\cdot))$ ,  $\hat{\theta} = (\hat{\alpha}, \hat{\eta}(\cdot))$ , and  $\theta(x) = (\alpha, \eta(x))$ . We define  $\bar{I}(x, a, b) = \mathbb{E}[\bar{m}(x, Y(1), a, b) \mid X = x, T = 1]$ , which is convex in  $(a, b) \in \mathbb{R}^+ \times \mathbb{R}$ .

**Assumption B.1.** *Suppose  $\mathcal{X}$  is a compact set, and  $\mathbb{P}_{X \mid T=1}$  has positive density on  $\mathcal{X}$ . Suppose  $\bar{\eta} \in \Theta = \mathbb{R}^+ \times \Lambda_c^p(\mathcal{X})$  for some  $c > 0$ . Assume  $|\bar{m}(x, y, a, b) - \bar{m}(x, y, \bar{\alpha}, \bar{\eta}(x)) - \nabla_{a,b} \bar{m}(x, y, \bar{\alpha}, \bar{\eta}(x))[a - \bar{\alpha}, b - \bar{\eta}(x)]| \leq g(x, y) \|(a - \bar{\alpha}, b - \bar{\eta}(x))\|_2^2$  for  $\|(a - \bar{\alpha}, b - \bar{\eta}(x))\|_2 < \epsilon$  for sufficiently small  $\epsilon > 0$ , where  $\|\cdot\|_2$  is the Euclidean norm, and  $\mathbb{E}[g(x, Y(1)) \mid X = x, T = 1] < M$  for some constant  $M > 0$ . We assume that  $|\bar{I}(x, a, b) - \bar{I}(x, \bar{\alpha}^{(j)}, \bar{\eta}(x)) - \nabla_{a,b} \bar{I}(x, \bar{\alpha}, \bar{\eta}(x))[a - \bar{\alpha}, b - \bar{\eta}(x)]| \geq \lambda \|(a - \bar{\alpha}, b - \bar{\eta}(x))\|_2^2$  for  $\|(a - \bar{\alpha}, b - \bar{\eta}(x))\|_2 < \epsilon$  for sufficiently small  $\epsilon > 0$ , where  $\lambda > 0$  is a constant. Also, assume  $|\bar{m}(\theta, x, y) - \bar{m}(\bar{\theta}^{(j)}, x, y)| \leq \bar{g}(x, y) \|\theta(x) - \bar{\theta}(x)\|_2$  for  $\|\theta(x) - \bar{\theta}(x)\|_2 < \epsilon$  for sufficiently small  $\epsilon > 0$ , and  $\sup_{x \in \mathcal{X}} \mathbb{E}[\bar{g}(x, Y(1))^2 \mid X = x, T = 1] < M$  for some constant  $M > 0$ . Furthermore, we assume the exchangeability of differentiation and expectation for  $\mathbb{E}[\nabla_{a,b} \bar{m}(x, Y(1), \bar{\alpha}, \bar{\eta}(x)) \mid X = x, T = 1]$  and  $\mathbb{E}[\nabla_{a,b} \bar{m}(X, Y(1), \bar{\alpha}, \bar{\eta}(X)) \mid T = 1]$ .*

In Assumption B.1, we assume the smoothness of the underlying optimizer, which is reasonable if  $\hat{r}$  is smooth and the distribution of  $Y(1) \mid X = x$  is smooth in  $x$  (c.f. the discussion in Appendix A.2 for the conditional  $f$ -sensitivity model). The stable expansion condition is also mild, which is satisfied if the risk function is second-order smooth or have a Lipschitz first-order derivative. The condition on the convex function  $\bar{I}$  can be satisfied if its Hessian is positive definite (or strongly convex) at  $(\bar{\alpha}, \bar{\eta}(x))$  for all  $x$  in the compact set  $\mathcal{X}$ . The first-order expansions of  $\bar{m}$  are similar to what we require in Assumption 2; concrete examples are similar to those discussed in Appendix A.2.

The following proposition provides the convergence rate of  $(\hat{\alpha}, \hat{\eta}(\cdot))$  to  $(\bar{\alpha}, \bar{\eta}(\cdot))$  with sieve estimation under the above regularity conditions. Its proof is in Appendix D.3.

**Proposition B.4.** *Suppose Assumption B.1 holds. We set  $J_n = (\frac{\log n}{n})^{1/(2p+d)}$  for the sieve estimators in Examples 2 or 3. Denote the empirical risk for sieve estimation with parameters  $\theta$  as  $\hat{\mathbb{E}}_n[\bar{m}(\theta, X, Y(1))]$ , and suppose  $\hat{\theta}$  satisfies  $\hat{\mathbb{E}}_n[\bar{m}(\hat{\theta}, X, Y(1))] \leq \inf_{\theta \in \Theta_n} \hat{\mathbb{E}}_n[\bar{m}(\theta, X, Y(1))] - O_P((\frac{\log n}{n})^{\frac{2p}{2p+d}})$ . Then  $\|\hat{\theta} - \bar{\theta}\|_{L_2(\mathbb{P}_{\cdot \mid T=1})} = O_P((\frac{\log n}{n})^{\frac{p}{2p+d}})$ ,  $\|\hat{\theta} - \bar{\theta}\|_\infty = O_P((\frac{\log n}{n})^{\frac{2p^2}{(2p+d)^2}})$ .*

### B.3.3 Theoretical guarantees

Our key additional assumption for Wald-type inference is the stability of population risk minimizers. It is needed for our analysis because estimating  $\hat{\alpha}$  and  $\hat{\eta}$  is based on an estimated  $\hat{r}$ , thereby leading to optimization error compared with estimating them based on the true  $r_{1,0}(x)$ . In a nutshell, it ensures that the estimation error of  $\hat{\alpha}$  and  $\hat{\eta}$  is of the same order as that of  $\hat{r}$ .

**Assumption B.2.** *For each  $j$ , it holds that  $\|(\bar{\alpha}^{(j)}, \bar{\eta}^{(j)}(\cdot)) - (\alpha^*, \eta^*(\cdot))\|_{L_2(\mathbb{P}_{\cdot \mid T=1})} = O(\|\hat{r}^{(j)} - r_{1,0}\|_{L_2(\mathbb{P}_{\cdot \mid T=1})})$ , where the  $L_2$ -norm is defined viewing  $\alpha$  and  $\bar{\alpha}^{(j)}$  as constant functions.*

Under the above assumption, the correction term using  $\hat{h}^{(j)}$  helps to remove the bias incurred by  $\hat{r}^{(j)}$  in each fold, so that the errors caused by nuisance estimation is negligible compared to  $n^{-1/2}$ . Before presenting our inferential guarantee, we state estimation rates and some mild regularity conditions. The functions  $\hat{r}^{(j)}$ ,  $\bar{\alpha}^{(j)}$  and  $\bar{\eta}^{(j)}$  are all viewed as fixed functions when taking expectations. We also denote  $\hat{\theta}^{(j)} = (\hat{\alpha}^{(j)}, \hat{\eta}^{(j)})$  and  $\bar{\theta}^{(j)} = (\bar{\alpha}^{(j)}, \bar{\eta}^{(j)})$ , on which  $L_2(\mathbb{P}_{X \mid T=1})$  is defined by viewing the first argument as a constant function.

**Assumption B.3.** Assume  $\|\widehat{r}^{(j)} - r_{1,0}\|_{L_2(\mathbb{P}_{\cdot|T=1})} = o_P(n^{-1/4})$ ,  $\|\widehat{\theta}^{(j)} - \bar{\theta}^{(j)}\|_{L_2(\mathbb{P}_{\cdot|T=1})} = o_P(n^{-1/4})$  and  $\|\widehat{h}^{(j)} - \bar{h}^{(j)}\|_{L_2(\mathbb{P}_{\cdot|T=1})} = o_P(n^{-1/4})$  for all  $j$ .

For each  $j$ , assume that  $|m(x, y, a, b, \widehat{r}^{(j)}(x)) - m(x, y, \bar{\alpha}^{(j)}, \bar{\eta}^{(j)}(x), \widehat{r}^{(j)}(x)) - \nabla_{a,b} m(x, y, \bar{\alpha}^{(j)}, \bar{\eta}^{(j)}(x), \widehat{r}^{(j)}(x)) [a - \bar{\alpha}^{(j)}, b - \bar{\eta}^{(j)}(x)]| \leq g(x, y) \|(a, b) - (\bar{\alpha}^{(j)}, \bar{\eta}^{(j)}(x))\|_2^2$  for all  $(a, b)$  in a neighborhood of  $(\bar{\alpha}^{(j)}, \bar{\eta}^{(j)}(x))$  when  $\|\widehat{r}^{(j)} - r_{1,0}\|_{L_2(\mathbb{P}_{X|T=1})}$  is sufficiently small, where  $\mathbb{E}[g(x, Y(1))^2 | X = x, T = 1] \leq M$  for some constant  $M > 0$ . We assume  $\mathbb{E}[(\nabla_b m(x, Y(1), \bar{\alpha}^{(j)}, \bar{\eta}^{(j)}(x), \widehat{r}^{(j)}(x)))^2 | X = x, T = 1] \leq M$ ,  $\mathbb{E}[(\nabla_a m(X, Y(1), \bar{\alpha}^{(j)}, \bar{\eta}^{(j)}(X), \widehat{r}^{(j)}(X)))^2 | T = 1] \leq M$  for some constant  $M > 0$  and  $\mathbb{P}_{X|T=1}$ -almost all  $x$ , and the differentiation and expectation is exchangeable in  $\mathbb{E}[\nabla_b m(x, Y(1), \bar{\alpha}^{(j)}, \bar{\eta}^{(j)}(x), \widehat{r}^{(j)}(x)) | X = x, T = 1]$  and  $\mathbb{E}[\nabla_a m(X, Y(1), \bar{\alpha}^{(j)}, \bar{\eta}^{(j)}(X), \widehat{r}^{(j)}(X)) | T = 1]$ . Also assume that  $|f^*(\frac{r}{-\bar{\alpha}}(y + \bar{\eta}(x))) - f^*(\frac{r_{1,0}(x)}{-\bar{\alpha}}(y + \bar{\eta}(x))) - \frac{\partial f^*}{\partial r}(\frac{r_{1,0}(x)}{-\bar{\alpha}}(y + \bar{\eta}(x)))(r - r_{1,0}(x))| \leq \bar{g}(x, y)(r - r_{1,0}(x))^2$  for  $|r - r_{1,0}(x)| < \epsilon$  for  $\mathbb{P}_{X|T=1}$ -almost all  $x$ , where  $\epsilon > 0$  is a sufficiently small constant, and  $\mathbb{E}[\bar{g}(x, Y(1))^2 | X = x, T = 1] < M$  for some constant  $M > 0$ .

In Assumption B.3, we require a slow rate  $o(n^{-1/4})$  of convergence for the regression problem of  $\widehat{r}^{(j)}$  and  $\widehat{h}^{(j)}$ , which are achievable by many methods under smoothness of the underlying regression functions. The estimation of  $\widehat{\theta}^{(j)}$  is the same as we discussed in Section B.3.2.

The regularity conditions we require on the smoothness of  $m$  are mild and similar to preceding ones. For example, the second order expansion can be satisfied if the first-order derivatives are locally Lipschitz or smooth, for which the uniform upper bound of  $\mathbb{E}[g(x, Y(1))^2 | X = x, T = 1]$  exists as long as the per- $x$  expectations are finite within the compact set  $\mathcal{X}$ . Other stability conditions can be similarly satisfied.

We then have the following inferential guarantee for  $\widehat{\mu}_{1,0}^{m,-}$ . The proof of Theorem B.1 is in Appendix D.4.

**Theorem B.1.** Suppose Assumptions B.2 and B.3 hold. Then the estimator (4) satisfies  $\sqrt{n}(\widehat{\mu}_{1,0}^{m,-} - \mu_{1,0}^{m,-}) \rightsquigarrow N(0, \text{Var}(\phi_{1,-}^m(X, Y, T)))$ , where  $\mu_{1,0}^{m,-}$  is the lower bound of  $\mathbb{E}[Y(1) | T = 0]$  under the marginal  $(f, \rho)$ -selection condition in Proposition 2, and

$$\phi_{1,-}^m(X, Y, T) = \frac{T}{p_1} \left\{ m(X, Y, \alpha^*, \eta^*(X), r_{0,1}(X)) - \alpha^* r_{0,1}(X) h(X) \right\} + \frac{1-T}{p_0} \alpha^* h(X),$$

with  $p_1 = \mathbb{P}(T = 1) = 1 - p_0$ . Here we define  $h(x) = \mathbb{E}[H(x, Y(1)) | X = x, T = 1]$ , where  $H(x, y) = H(x, y) = \frac{\partial f^*}{\partial r}(x, y, \alpha^*, \eta^*(x), r_{0,1}(x))$ , and  $\alpha^*, \eta^*(\cdot)$  is the population optimizer in (3). Define

$$\widehat{\sigma}_{1,0,m,-}^2 := \frac{1}{\widehat{p}_1} \left( \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} d_{1,i}^2 - \left( \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} d_{1,i} \right)^2 \right) + \frac{1}{\widehat{p}_0} \left( \frac{1}{n_0} \sum_{i \in \mathcal{I}_0} d_{0,i}^2 - \left( \frac{1}{n_0} \sum_{i \in \mathcal{I}_0} d_{0,i} \right)^2 \right)$$

where  $\widehat{p}_1 = |\mathcal{I}_1|/n$ ,  $\widehat{p}_0 = |\mathcal{I}_0|/n$ ,  $d_{1,i} = m(X_i, Y_i, \widehat{\alpha}^{j[i]}, \widehat{\eta}^{j[i]}(X_i), \widehat{r}^{j[i]}(X_i)) - \widehat{\alpha}^{j[i]} \widehat{r}^{j[i]}(X_i) \widehat{h}^{j[i]}(X_i)$ ,  $d_{0,i} = \widehat{\alpha}^{j[i]} \widehat{h}^{j[i]}(X_i)$ , and  $j[i]$  is the fold that sample  $i$  lies in. Then  $\sqrt{n}(\widehat{\mu}_{1,0}^{m,-} - \mu_{1,0}^{m,-}) / \widehat{\sigma}_{1,0,m,-} \rightsquigarrow N(0, 1)$ .

Finally, similar to Section 4.4, the consistent naive estimators and inference procedures we introduce in this section can be applied to other quantities and combined to form bounds on average treatment effects.

## C Technical proofs for conditional $f$ -sensitivity models

### C.1 Proof of Lemma 1

*Proof.* Proof of Lemma 1 Suppose a distribution  $\mathbb{P}$  over  $(X, U, T, Y(0), Y(1))$  satisfies the  $(f, \rho)$ -selection condition. We note that the likelihood ratio can be decomposed as

$$\begin{aligned} \frac{d\mathbb{P}_{Y(1),U|X=x,T=0}}{d\mathbb{P}_{Y(1),U|X=x,T=1}} &= \frac{d\mathbb{P}_{Y(1)|U,X=x,T=0}}{d\mathbb{P}_{Y(1)|U,X=x,T=1}} \cdot \frac{d\mathbb{P}_{U|X=x,T=0}}{d\mathbb{P}_{U|X=x,T=1}} \\ &\stackrel{(a)}{=} \frac{d\mathbb{P}_{U|X=x,T=0}}{d\mathbb{P}_{U|X=x,T=1}} = \frac{\mathbb{P}(T=0|X=x,U)}{\mathbb{P}(T=1|X=x,U)} \cdot \frac{\mathbb{P}(T=1|X=x)}{\mathbb{P}(T=0|X=x)}, \end{aligned}$$

where step (a) is due to condition (1). Marginalizing over the conditional distribution of  $Y(1), U$  given  $X = x, T = 1$  yields the first result in Lemma 1.

In addition, by condition (1) and the data-processing inequality,

$$\begin{aligned} D_f(\mathbb{P}_{Y(1)|X=x,T=0} \parallel \mathbb{P}_{Y(1)|X=x,T=1}) &\leq D_f(\mathbb{P}_{Y(1),U|X=x,T=0} \parallel \mathbb{P}_{Y(1),U|X=x,T=1}) \\ &= \mathbb{E}_{Y(1),U|X=x,T=1} \left[ f \left( \frac{d\mathbb{P}_{Y(1),U|X=x,T=0}}{d\mathbb{P}_{Y(1),U|X=x,T=1}} \right) \right]. \end{aligned}$$

Combining the above two facts yields

$$\begin{aligned} D_f(\mathbb{P}_{Y(1)|X,T=0} \parallel \mathbb{P}_{Y(1)|X,T=1}) &\leq \mathbb{E}_{Y(1),U|X,T=1} \left[ f \left( \frac{\mathbb{P}(T=0|X,U)}{\mathbb{P}(T=1|X,U)} \cdot \frac{\mathbb{P}(T=1|X)}{\mathbb{P}(T=0|X)} \right) \right] \\ &= \mathbb{E}_{U|X,T=1} \left[ f \left( \frac{\mathbb{P}(T=0|X,U)}{\mathbb{P}(T=1|X,U)} \cdot \frac{\mathbb{P}(T=1|X)}{\mathbb{P}(T=0|X)} \right) \right] \leq \rho \end{aligned}$$

almost surely, where the last inequality is due to the  $(f, \rho)$ -selection condition.  $\square$

## C.2 Proof of Proposition 1

*Proof.* Proof of Proposition 1 We prove the claim by showing that  $\mathcal{Q}_{1,0}$  and the set on the right-handed side (RHS) of Proposition 1 contain each other.

**Step 1:  $\mathcal{Q}_{1,0} \subseteq \text{RHS}$ .**

By definition, any member of  $\mathcal{Q}_{1,0}$  is the induced distribution over  $(X, Y(1))$  given  $T = 0$  of some super-population  $Q \in \mathcal{M}_1^{f,\rho}$  with  $\mathbb{Q}_{X,Y,T} = \mathbb{P}_{X,Y,T}$ . Thus, we write it as  $\mathbb{Q}_{X,Y(1)|T=0} \in \mathcal{Q}_{1,0}$ . Since  $\mathbb{Q}_{X,Y,T} = \mathbb{P}_{X,Y,T}$ , we know that

$$\frac{d\mathbb{Q}_{X|T=0}}{d\mathbb{P}_{X|T=1}} = \frac{d\mathbb{P}_{X|T=0}}{d\mathbb{P}_{X|T=1}} = r_{1,0}(x).$$

By Lemma 1,  $D_f(\mathbb{Q}_{X,Y(1)|T=0} \parallel \mathbb{P}_{X,Y(1)|T=1}) \leq \rho$ . Also, we note that  $\mathbb{Q}_{Y(1)|X,T=1} = \mathbb{Q}_{Y|X,T=1} = \mathbb{P}_{Y|X,T=1}$  since  $Y = TY(1) + (1-T)Y(0)$  and  $\mathbb{Q}_{X,Y,T} = \mathbb{P}_{X,Y,T}$ . Combining the above, we establish the “ $\supseteq$ ” direction. It remains to prove the reverse.

**Step 2:  $\text{RHS} \subseteq \mathcal{Q}_{1,0}$ .**

Given any distribution  $Q$  over  $(X, Y(1))$  in RHS, we aim to find a distribution (super-population)  $\mathbb{Q}$  over  $(X, Y(0), Y(1), U, T)$  such that

- $(Y(1), Y(0)) \perp\!\!\!\perp T | X, U$  under  $\mathbb{Q}$ ;
- $\mathbb{Q}_{X,T,Y} = \mathbb{P}_{X,T,Y}$ ;
- $\mathbb{Q} \in \mathcal{M}_1^{f,\rho}$ ;
- $\mathbb{Q}_{X,Y(1)|T=0}(x, y) = Q(x, y)$ .

These conditions imply that  $Q \in \mathcal{Q}_{1,0}$ , and thus  $\text{RHS} \subseteq \mathcal{Q}_{1,0}$ . To construct  $\mathbb{Q}$ , we first set  $\mathbb{Q}_{X,T} = \mathbb{P}_{X,T}$ . Then we specify the distribution of  $Y(1)|X, T$  via

$$\mathbb{Q}_{Y(1)|X,T=1} = \mathbb{P}_{Y|X,T=1}, \quad \mathbb{Q}_{Y(1)|X,T=0} = \mathbb{Q}_{Y|X}.$$

So far,  $\mathbb{Q}_{X,T,Y(1)}$  has been determined. We let  $U = Y(1)$  be the unobserved confounder. Finally, the distribution of  $Y(0)|X, T, Y(1), U$  is specified via

$$\mathbb{Q}_{Y(0)|X,T,Y(1),U} = \mathbb{Q}_{Y(0)|X} = \mathbb{P}_{Y|X,T=0}.$$

Having constructed  $\mathbb{Q}$ , we proceed to check that it satisfies the conditions listed at the beginning of Step 2. Since  $Y(1) = U$  and  $\mathbb{Q}_{Y(0)|X,T,Y(1),U} = \mathbb{Q}_{Y(0)|X}$  by construction, it is easy to check that  $(Y(1), Y(0)) \perp\!\!\!\perp T | X, U$ . By construction, it also clearly holds that  $\mathbb{Q}_{X,T,Y} = \mathbb{P}_{X,T,Y}$ .

We now check that  $\mathbb{Q} \in \mathcal{M}_1^{f,\rho}$ . For any  $x$ , again by the construction of  $\mathbb{Q}$ ,

$$\begin{aligned} \mathbb{E}_{\mathbb{Q}_{Y(1)|X=x,T=1}} \left[ f \left( \frac{d\mathbb{Q}_{Y(1)|X=x,T=0}}{d\mathbb{Q}_{Y(1)|X=x,T=1}} \right) \right] &= \mathbb{E}_{\mathbb{P}_{Y|X=x,T=1}} \left[ f \left( \frac{dQ_{Y|X=x}}{d\mathbb{P}_{Y|X=x,T=1}} \right) \right] \\ &= D_f(Q_{Y|X=x} \parallel \mathbb{P}_{Y|X=x,T=1}) \leq \rho, \end{aligned}$$

where the last inequality is due to the fact that  $Q$  is in the set on the RHS. Therefore, we verify that  $\mathbb{Q} \in \mathcal{M}_1^{f,\rho}$ . By construction,  $\mathbb{Q}_{Y(1)|X,T=0} = Q_{Y|X}$ . It remains to show that  $\mathbb{Q}_{X|T=0} = Q_X$ . For any measurable set  $A$ ,

$$\begin{aligned} \mathbb{Q}(X \in A | T=0) &= \mathbb{E}_{\mathbb{Q}} \left[ \frac{d\mathbb{Q}_{X|T=0}}{d\mathbb{Q}_{X|T=1}} \cdot \mathbf{1}\{X \in A\} \mid T=1 \right] = \mathbb{E}_{\mathbb{Q}} \left[ \frac{dQ_{X|T=0}}{d\mathbb{Q}_{X|T=1}} \cdot \mathbf{1}\{X \in A\} \mid T=1 \right] \\ &= \mathbb{E}_{\mathbb{Q}} \left[ r_{1,0}(X) \cdot \mathbf{1}\{X \in A\} \mid T=1 \right] = \mathbb{E}_{\mathbb{Q}} \left[ \frac{dQ_X}{d\mathbb{Q}_{X|T=1}} \cdot \mathbf{1}\{X \in A\} \mid T=1 \right] = Q(X \in A). \end{aligned}$$

Since the above holds for any measurable set  $A$ , we know that  $\mathbb{Q}_{X|T=0} = Q_X$ . Finally, switching the role of the treated and control groups yields similar results for the counterfactual distribution of  $(X, Y(0))$  given  $T=1$  under  $\mathcal{M}_0^{f,\rho}$ . This completes the proof of Proposition 1.  $\square$

### C.3 Proof of Proposition 3

*Proof.* Proof of Proposition 3 We first claim that solving (5) amounts to solving the following problem for each  $x$ :

$$\begin{aligned} \min_{L(x) \text{ measurable}} \quad & \mathbb{E}[Y(1)L(x) | X=x, T=1] \\ \text{s.t.} \quad & \mathbb{E}[L(x) | X=x, T=1] = r_{1,0}(x) \\ & \mathbb{E}[f(L(x)/r_{1,0}(x)) | X=x, T=1] \leq \rho. \end{aligned} \tag{5}$$

To be specific, denoting the optimal objective of (5) as  $\mu(x)$  and that of (5) as  $\mu_{1,0}^-$ , we are to show that  $\mu_{1,0}^- = \mathbb{E}[\mu(X) | T=1]$ . To see why it is the case, suppose  $L^*$  is the optimizer of (5), then it is measurable with respect to  $X$  and  $Y(1)$  and satisfies the constraints of (5). Then  $L(x)(\cdot) := L^*(x, \cdot)$  is measurable with respect to  $Y(1)$ , and satisfy the constraints of (5). As a result, we have  $\mathbb{E}[Y(1)L^*(x, Y(1)) | X=x, T=1] \geq \mu(x)$ . Marginalizing over  $X$  yields  $\mu_{1,0}^- = \mathbb{E}[Y(1)L(X, Y(1)) | T=1] \geq \mathbb{E}[\mu(X) | T=1]$ . On the other hand, suppose  $L^*(x)(\cdot)$  is measurable with respect to  $Y(1)$  and is the minimizer for (5) for  $\mathbb{P}_{X|T=1}$ -almost all  $x$ . We let  $L(x, y) = L^*(x)(y)$ , so that it is measurable with respect to  $(X, Y(1))$  and satisfy the constraints of (5). Thus we have  $\mathbb{E}[L(X, Y(1))Y(1) | T=1] = \mathbb{E}[\mu(X) | T=1] \geq \mu_{1,0}^-$ . Combining the two directions leads to the equivalence.

In the following, we solve (5) and write  $\mathbb{E}_x$  in place of  $\mathbb{E}[\cdot | X=x, T=1]$  for simplicity. Invoking [Luenberger \(1997, Theorem 8.6.1\)](#) to this convex problem, we have

$$\min_{\substack{\mathbb{E}_x[L]=r_{1,0}(x), \\ \mathbb{E}_x[f(L/r_{1,0}(X))] - \rho \leq 0}} \mathbb{E}_x[Y(1)L(x)] = \max_{\alpha \geq 0, \eta \in \mathbb{R}} \varphi(\alpha, \eta, x),$$

where the Slater's condition is satisfied and strong duality holds, and

$$\begin{aligned} \varphi(\alpha, \eta, x) &= \inf_{L \geq 0 \text{ measurable}} \mathcal{L}(\alpha, \eta, L, x), \\ \mathcal{L}(\alpha, \eta, L, x) &= \mathbb{E}_x[Y(1)L(x)] + \eta \mathbb{E}_x[L - r_{1,0}(x)] + \alpha (\mathbb{E}_x[f(L/r_{1,0}(x))] - \rho). \end{aligned}$$

The minimum of  $\mathcal{L}(\alpha, \eta, L, x)$  is thus given by

$$\begin{aligned} \varphi(\alpha, \eta, x) &= \mathbb{E}_x \left[ \min_{z \geq 0} \{Y(1)z + \eta z - \eta r_{1,0}(x) + \alpha f(z/r_{1,0}(x)) - \alpha \rho\} \right] \\ &= \mathbb{E}_x \left[ -\alpha f^* \left( \frac{r_{1,0}(x)}{-\alpha} (Y(1) + \eta) \right) - \eta r_{1,0}(x) - \alpha \rho \right]. \end{aligned}$$

Now we write  $\alpha(x)$  and  $\eta(x)$  to emphasize its dependency on  $x$ . Therefore, by the equivalence discussed in the beginning, we have

$$\begin{aligned}\mu_{1,0}^- &= \mathbb{E} \left[ \max_{\alpha(X) \geq 0, \eta(X) \in \mathbb{R}} \varphi(\alpha(X), \eta(X), X) \mid T = 1 \right] \\ &= \mathbb{E} \left[ \varphi(\alpha^*(X), \eta^*(X), X) \mid T = 1 \right],\end{aligned}$$

where for  $\mathbb{P}_{X|T=1}$ -almost all  $x$ ,

$$(\alpha^*(x), \eta^*(x)) \in \arg \max_{\alpha \geq 0, \eta \in \mathbb{R}} \mathbb{E} \left[ -\alpha f^* \left( \frac{r_{1,0}(x)}{-\alpha} (Y(1) + \eta) \right) - \eta r_{1,0}(x) - \alpha \rho \mid X = x, T = 1 \right].$$

With a change-of-variable from  $\alpha(x)$  to  $\alpha(x)r_{1,0}(x)$ , we have

$$(\alpha^*(x)/r_{1,0}(x), \eta^*(x)) \in \arg \max_{\alpha \geq 0, \eta \in \mathbb{R}} \mathbb{E} \left[ -\alpha r_{1,0}(x) f^* \left( \frac{Y(1) + \eta}{-\alpha} \right) - \eta r_{1,0}(x) - \alpha r_{1,0}(x) \rho \mid X = x, T = 1 \right].$$

The minimum of (5) can thus be written as

$$\mu_{1,0}^- = -\mathbb{E} \left[ r_{1,0}(x) \left\{ \alpha^*(X) f^* \left( \frac{Y(1) + \eta^*(X)}{-\alpha^*(X)} \right) + \eta^*(X) + \alpha^*(X) \rho \right\} \mid T = 1 \right],$$

where for  $\mathbb{P}_{X|T=1}$ -almost all  $x$ , it holds that

$$(\alpha^*(x), \eta^*(x)) \in \arg \min_{\alpha \geq 0, \eta \in \mathbb{R}} \mathbb{E} \left[ \alpha f^* \left( \frac{Y(1) + \eta}{-\alpha} \right) + \eta + \alpha \rho \mid X = x, T = 1 \right].$$

Therefore, we complete the proof of Proposition 3.  $\square$

## C.4 Proof of convergence of sieve estimator

*Proof.* Proof of Theorem 1 We analyze the behavior of  $\widehat{\theta}^{(j)}$  for each fold  $j$ . As  $|\mathcal{I}_1^{(j)}| \asymp n$ , we take the generic notation of  $\widehat{\theta}$  and sample size  $n$ , so that

$$\widehat{\mathbb{E}}_n[\ell(\widehat{\theta}, X, Y(1))] \geq \inf_{\theta \in \Theta_n} \widehat{\mathbb{E}}_n[\ell(\theta, X, Y(1))] - O_P\left(\left(\frac{\log n}{n}\right)^{2p/(2p+d)}\right),$$

where  $(X_i, Y_i) \sim \mathbb{P}_{X, Y(1) | T=1}$  are i.i.d. data. For some fixed  $b > 0$ , we denote the sequence

$$\delta_n := \inf \left\{ \delta \in (0, 1) : \frac{1}{\sqrt{n}\delta^2} \int_{b\delta^2}^{\delta} \sqrt{\log N(\epsilon^{1+d/2p}, \Theta_n, \|\cdot\|_{L_2(\mathbb{P}_{\cdot|T=1})})} d\epsilon \leq 1 \right\},$$

where  $N(\epsilon, \Theta_n, \|\cdot\|_{L_2(\mathbb{P}_{\cdot|T=1})})$  is the  $\epsilon$ -covering number of  $\Theta_n$  in the  $L_2$ -norm under  $\mathbb{P}_{\cdot|T=1}$ . We employ the established convergence results for sieve estimators adapted from Chen (2007, Theorem 3.2) and Yadlowsky et al. (2018, Lemma B.3), stated in Lemma C.1.

**Lemma C.1.** *Let  $\theta^* \in \Theta$  be a population risk minimizer. Suppose there exists constants  $c_1, c_2 > 0$  such that  $c_1 \mathbb{E}[\ell(\theta, X, Y) - \ell(\theta^*, X, Y)] \leq d(\theta, \theta^*)^2 \leq c_2 \mathbb{E}[\ell(\theta, X, Y) - \ell(\theta^*, X, Y)]$  for  $\theta$  in a neighborhood of  $\theta^*$ . Suppose the following conditions hold:*

- (i) *For sufficiently small  $\epsilon > 0$ ,  $\text{Var}(\ell(\theta, X, Y) - \ell(\theta^*, X, Y)) \leq C_1 \epsilon^2$  for all  $\theta \in \Theta_n$  such that  $d(\theta, \theta^*) \leq \epsilon$ .*
- (ii) *For any  $\delta > 0$ , there exists a constant  $s \in (0, 2)$  and a measurable function  $U_n(\cdot)$  such that  $\sup_n \mathbb{E}[U_n(X, Y)^2] \leq C_3$  and  $\sup_{\theta \in \Theta_n : d(\theta, \theta^*) \leq \delta} |\ell(\theta, X, Y) - \ell(\theta^*, X, Y)| \leq \delta^s U_n(X, Y)$  for constant  $C_3 > 0$ .*

Then  $d(\widehat{\theta}_n - \theta^*) = O_P(\max\{\delta_n, \inf_{\theta' \in \Theta_n} d(\theta', \theta^*)\})$ .

We define the distance as  $L_2$ -norm  $d(\theta, \theta') = \|\theta - \theta'\|_{L_2(\mathbb{P})}$ , and verify the conditions in Lemma C.1. We define  $\Theta = \Lambda_c^p(\mathcal{X}) \times \Lambda_c^p(\mathcal{X})$  without truncation. The upper bound  $\mathbb{E}[\ell(\theta, X, Y) - \ell(\theta^*, X, Y) | T = 1]$  is directly implied by Assumption 2. By the  $\lambda$ -strong convexity of  $\mathbb{E}[\ell((a, b), x, Y) | X = x]$  is at  $(a, b) = \theta^*(x)$ ,

$$\mathbb{E}[\ell(\theta(x), x, Y) | X = x] - \mathbb{E}[\ell(\theta^*(x), x, Y) | X = x] \geq \lambda(\theta(x) - \theta^*(x))^2.$$

Integrating over  $X$  yields  $\mathbb{E}[\ell(\theta(X), X, Y)] - [\ell(\theta^*(X), X, Y)] \geq c''d(\theta, \theta^*)$  for some constant  $c'' > 0$ .

We then check condition (i). By the positive density condition, we have  $\|\cdot\|_{L_2(\lambda)} \asymp \|\cdot\|_{L_2(\mathbb{P})}$ . Hence  $\|\theta - \theta^*\|_\infty = o(1)$  once  $\|\theta - \theta^*\|_{L_2(\mathbb{P})} = o(1)$ . By Lemma 2 of Chen and Shen (1998), we have  $\|\theta\|_\infty \lesssim \|\theta\|_{L_2(\lambda)}^{2p/(2p+d)}$  for any  $\theta \in \Theta$ , where  $\lambda$  is the Lebesgue measure. Therefore, sufficiently small  $\|\theta - \theta^*\|_{L_2(\mathbb{P})}$  implies sufficiently small  $\|\theta - \theta^*\|_\infty$ . Since for  $\|\theta - \theta^*\|_\infty$  sufficiently small,  $|\ell(\theta, x, y) - \ell(\theta^*, x, y)| \leq \bar{\ell}(x, y)(\theta(x) - \theta^*(x))$  where  $\mathbb{E}[\bar{\ell}(x, Y)^2 | X = x] \leq M$  for all  $x$ , we have

$$\text{Var}(\ell(\theta, X, Y) - \ell(\theta^*, X, Y)) \leq \mathbb{E}[|\ell(\theta, x, y) - \ell(\theta^*, x, y)|^2] \leq M\mathbb{E}[(\theta(X) - \theta^*(X))^2] \leq M\epsilon^2$$

for all  $\theta \in \Theta_n$  such that  $d(\theta, \theta^*) \leq \epsilon$  for sufficiently small  $\epsilon > 0$ . Condition (ii) follows from the same argument by taking  $U_n(x, y) = \bar{\ell}(x, y)$ . Therefore, applying Lemma C.1 we have  $\|\hat{\theta}_n - \theta^*\|_{L_2(\mathbb{P})} = O_P(\max\{\delta_n, \inf_{\theta' \in \Theta_n} d(\theta', \theta^*)\})$ . Here according to Chen and Shen (1998) and Geer et al. (2000), we have

$$\log N(\epsilon, \Theta_n^\eta, \|\cdot\|_{2, \mathbb{P}}) \lesssim \dim(\Theta_n^\eta) \log \frac{1}{\epsilon},$$

where  $\dim(\Theta_n^\eta) = J_n^p$ . Since truncation is a contraction map, the covering number of  $\Theta_n^\alpha$  is upper bounded by the above quantity. As a result, we have

$$\log N(\epsilon, \Theta_n, \|\cdot\|_{2, \mathbb{P}}) \lesssim J_n^p \log \frac{1}{\epsilon}.$$

Similar to the results in Yadlowsky et al. (2018), we have

$$\delta_n \asymp \sqrt{\frac{J_n^d \log n}{n}}.$$

We finally bound the approximation error using  $\Theta_n$ . Note that we take  $\Theta_n$  to be truncated at  $\epsilon$ . However, since the population minimizer  $\theta^*$  is uniformly bounded above  $\epsilon$ , since truncation is a contraction map, we have  $\inf_{\theta \in \Theta_n} \|\theta - \theta^*\|_{L_2(\mathbb{P})} \leq \inf_{\theta \in \Theta_n^\eta \times \Theta_n^\eta} \|\theta - \theta^*\|_{L_2(\mathbb{P})} \leq O(J_n^p)$ , where the last inequality is a well-established result, see, e.g., Timan (2014). We now set  $J_n = (n/\log n)^{1/(2p+d)}$ , so that  $\|\hat{\theta} - \theta^*\|_{L_2(\mathbb{P})} = O_P((\log n/n)^{p/(2p+d)})$ . This completes our proof.  $\square$

## C.5 Proof of Proposition 4

Given  $X = x$ , suppose instead  $\alpha^*(x) = 0$ . We consider the following two cases:

- If  $\eta^*(x) < -y(x)$ , then

$$\begin{aligned} & \liminf_{\alpha \rightarrow 0} \mathbb{E} \left[ \alpha f^* \left( \frac{Y(1) + \eta^*(x)}{-\alpha} \right) + \eta^*(x) + \alpha \rho \mid X = x, T = 1 \right] \\ &= \liminf_{\alpha \rightarrow 0} \mathbb{E} \left[ \alpha f^* \left( \frac{Y(1) + \eta^*(x)}{-\alpha} \right) \mathbf{1}\{Y(1) \leq -\eta^*(x)\} \right. \\ & \quad \left. + \alpha f^* \left( \frac{Y(1) + \eta^*(x)}{-\alpha} \right) \mathbf{1}\{Y(1) > -\eta^*(x)\} \mid X = x, T = 1 \right] + \eta^*(x) + \alpha \rho \\ &\stackrel{(a)}{\geq} \liminf_{\alpha \rightarrow 0} \mathbb{E} \left[ \alpha f^* \left( \frac{Y(1) + \eta^*(x)}{-\alpha} \right) \mathbf{1}\{Y(1) \leq -\eta^*(x)\} - \alpha L \mid X = x, T = 1 \right] \\ & \quad + \liminf_{\alpha \rightarrow 0} \left[ \alpha f^* \left( \frac{Y(1) + \eta^*(x)}{-\alpha} \right) \mathbf{1}\{Y(1) > -\eta^*(x)\} - \alpha L \mid X = x, T = 1 \right] + \eta^*(x) \\ &\stackrel{(b)}{\geq} +\infty, \end{aligned}$$

where step (a) uses the fact that  $\liminf_{n \rightarrow \infty} a_n + b_n \geq \liminf_{n \rightarrow \infty} a_n + \liminf_{n \rightarrow \infty} b_n$  and step (b) follows from Fatou's lemma and the condition that  $f^*(x)/x \rightarrow \infty$  when  $x \rightarrow \infty$ .

- If  $\eta^*(x) \geq -\underline{y}(x)$ , then

$$\begin{aligned} & \lim_{\alpha \rightarrow 0} \mathbb{E} \left[ \alpha f^* \left( \frac{Y(1) + \eta^*(x)}{-\alpha} \right) + \eta^*(x) + \alpha \rho \mid X = x, T = 1 \right] \\ &= \lim_{\alpha \rightarrow 0} \mathbb{E} \left[ \alpha f^* \left( \frac{Y(1) + \eta^*(x)}{-\alpha} \right) \mathbb{1}\{Y(1) \geq \text{essinf } Y(1)\} + \alpha \rho \mid X = x, T = 1 \right] + \eta^*(x) \\ &\stackrel{(a)}{=} \mathbb{E} \left[ \lim_{\alpha \rightarrow 0} \alpha f^* \left( \frac{Y(1) + \eta^*(x)}{-\alpha} \right) \mathbb{1}\{Y(1) \geq \text{essinf } Y(1)\} + \alpha \rho \mid X = x, T = 1 \right] + \eta^*(x) \\ &\stackrel{(b)}{=} \eta^*(x). \end{aligned}$$

Above, step (a) is due to the fact that  $f^*(x)$  is bounded when  $x \leq 0$  and the dominated convergence theorem; step (b) is because  $f^*(x)/x \rightarrow 0$  as  $x \rightarrow -\infty$ .

Combining the two cases above, we conclude that  $\eta^*(x) = -\underline{y}(x)$  and the optimal value of the dual problem is  $-\underline{y}(x)$ . By the strong duality, the optimal value of the primal objective function is  $\mathbb{E}[r_{1,0}(X)\underline{y}(X) \mid T = 1]$ . As an implication, there exists a feasible  $L(x, y)$  such that  $\mathbb{E}[Y(1)L(X, Y(1)) \mid T = 1] = \mathbb{E}[r_{1,0}(X)\underline{y}(X) \mid T = 1]$ . Let  $\mathbb{Q}_{Y \mid X=x}$  denote the measure induced by  $L(x, y)$ :

$$\frac{d\mathbb{Q}_{Y \mid X=x}}{d\mathbb{P}_{Y(1) \mid X=x, T=1}}(y) = \frac{L(x, y)}{r_{1,0}(x)}.$$

This is a valid transformation of measure because  $L(x, y)$  is feasible. Then  $Y(1) = \underline{y}(X)$  a.s. under  $\mathbb{Q}_{Y \mid X}$ . Consequently,

$$\begin{aligned} 1 &= \mathbb{Q}(Y(1) = \underline{y}(X) \mid X = x) = \mathbb{E} \left[ \frac{L(x, \underline{y}(x))}{r_{1,0}(x)} \cdot \mathbb{1}\{Y(1) = \underline{y}(x)\} \mid X = x, T = 1 \right] = \frac{L(x, \underline{y}(x))}{r_{1,0}(x)} \cdot \bar{p}(x), \\ 0 &= L(x, Y(1)) \cdot \mathbb{1}\{Y(1) > \underline{y}(x)\}, \text{ a.s. under } \mathbb{P}_{Y(1) \mid X=x, T=1}. \end{aligned}$$

Again since  $L$  is feasible,

$$\rho \geq \mathbb{E} \left[ f \left( \frac{L(x, Y(1))}{r_{1,0}(x)} \right) \mid X = x, T = 1 \right] = \bar{p}(x) \cdot f \left( \frac{1}{\bar{p}(x)} \right) + (1 - \bar{p}(x)) \cdot f(0).$$

This is a contradiction to the condition. Therefore, we have  $\alpha^*(x) > 0$ .

## C.6 Proof of Theorem 2

*Proof.* Proof of Theorem 2 We consider the general scenario where  $(\hat{\alpha}^{(j)}, \hat{\eta}^{(j)})$  converges in sup-norm to some fixed  $(\alpha^\diamond, \eta^\diamond)$ , and show that  $-\hat{\mu}_{1,0}^{(j)} \xrightarrow{P} \mathbb{E}[r(X)\ell(\theta^\diamond(X), X, Y(1)) \mid T = 1]$  for any fixed  $j$ , where the risk function  $\ell$  is defined in Proposition 3. In the following, we drop the dependency on  $j$  for notational convenience; we are to show that with estimators  $\hat{r}$ ,  $\hat{H}$  and  $\hat{h}$  that are independent of  $\mathcal{I}_1$  and  $\mathcal{I}_0$ ,

$$\begin{aligned} \hat{\mu} &:= \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} \hat{r}(X_i) (\hat{H}(X_i, Y_i) - \hat{h}(X_i)) + \frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} \hat{h}(X_i) \\ &\xrightarrow{P} \mathbb{E}[r(X)\ell(\theta^\diamond(X), X, Y(1)) \mid T = 1]. \end{aligned} \tag{6}$$

Therefore, if  $\theta^\diamond = \theta^*$ , we have  $\hat{\mu}_{1,0} = \mu_{1,0}^- + o_P(1)$  since  $\mu_{1,0}^- = -\mathbb{E}[r(X)\ell(\theta^*(X), X, Y(1)) \mid T = 1]$  by Proposition 3. Otherwise, since  $\mathbb{E}[r(X)\ell(\theta^*(X), X, Y(1)) \mid T = 1] \leq \mathbb{E}[r(X)\ell(\theta^\diamond(X), X, Y(1)) \mid T = 1]$ , we have the one-sided validity that  $\hat{\mu}_{1,0} \xrightarrow{P} -\mathbb{E}[r(X)\ell(\theta^\diamond(X), X, Y(1)) \mid T = 1] \leq \mu_{1,0}^-$ , i.e., our estimator converges to a valid lower bound.

It thus remains to show (6). We prove the results when either  $\hat{r}$  or  $\hat{h}$  is consistent.

**Consistent  $\hat{r}$ .** We first show the case where  $\hat{r}$  is consistent for  $r_{1,0}$ , but not necessarily the regression function  $\hat{h}$ . Recall that  $\hat{H}(x, y) = \ell(\hat{\theta}(x), x, y)$ . Note that

$$\hat{\mu} = \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r(X_i)) (\hat{H}(X_i, Y_i) - \hat{h}(X_i)) + \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} r(X_i) (\hat{H}(X_i, Y_i) - \hat{h}(X_i)) + \frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} \hat{h}(X_i).$$

The first summation can be controlled as (where the expectation is implicitly conditional on other folds except  $\mathcal{I}_0^{(j)} \cup \mathcal{I}_1^{(j)}$ )

$$\begin{aligned} & \mathbb{E} \left[ \left( \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r_{1,0}(X_i)) (\hat{H}(X_i, Y_i) - \hat{h}(X_i)) \right)^2 \right] \\ & \leq \mathbb{E} \left[ (\hat{r}(X_i) - r_{1,0}(X_i))^2 \mathbb{E} [ (\hat{H}(X_i, Y_i) - \hat{h}(X_i))^2 \mid X_i ] \right] \leq M \cdot \|\hat{r} - r_{1,0}\|_{L_2(\mathbb{P}_{X \mid T=1})}^2 = o_P(1). \end{aligned}$$

Invoking Lemma E.1, we can drop the conditioning and the first summation is  $o_P(1)$ . On the other hand, since the covariate shift between  $\mathbb{P}_{X \mid T=1}$  and  $\mathbb{P}_{X \mid T=0}$  is exactly  $r_{1,0}$ , we know that  $\mathbb{E}[\hat{h}(X)r(X) \mid T=1] = \mathbb{E}[\hat{h}(X) \mid T=0]$ , where we still implicitly condition on other folds. As a result,

$$\begin{aligned} & - \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} r_{1,0}(X_i) \hat{h}(X_i) + \frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} \hat{h}(X_i) \\ & = - \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} (r_{1,0}(X_i) \hat{h}(X_i) - \mathbb{E}[\hat{h}(X)r(X) \mid T=1]) + \frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} (\hat{h}(X_i) - \mathbb{E}[r(X) \mid T=0]), \end{aligned}$$

where both terms are unbiased. Thus by Cauchy-Schwarz inequality,

$$\begin{aligned} & \mathbb{E} \left[ \left( - \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} r_{1,0}(X_i) \hat{h}(X_i) + \frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} \hat{h}(X_i) \right)^2 \right] \\ & \leq \frac{2}{|\mathcal{I}_1|} \text{Var}(r_{1,0}(X) \hat{h}(X) \mid T=1) + \frac{2}{|\mathcal{I}_0|} \text{Var}(\hat{h}(X) \mid T=0) = o_P(1) \end{aligned}$$

invoking the assumption that  $\hat{h}$  has finite second moment. Drop the conditioning by Lemma E.1, we know that this summation is also  $o_P(1)$ , hence

$$\begin{aligned} \hat{\mu} & = \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} r(X_i) \hat{H}(X_i, Y_i) + o_P(1) \\ & = \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} r(X_i) \ell(\theta^\circ(X_i), X_i, Y_i) + \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} r(X_i) \{ \ell(\hat{\theta}(X_i), X_i, Y_i) - \ell(\theta^\circ(X_i), X_i, Y_i) \} + o_P(1) \\ & = \mathbb{E}[r(X_i) \ell(\theta^\circ(X), X, Y(1)) \mid T=1] + \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} r(X_i) \{ \ell(\hat{\theta}(X_i), X_i, Y_i) - \ell(\theta^\circ(X_i), X_i, Y_i) \} + o_P(1). \end{aligned}$$

Finally, once  $\|\hat{\theta} - \theta^\circ\|_{\infty, \mathbb{P}_{X \mid T=1}} = o_P(1)$ , by the local expansion around  $\theta^\circ(x)$ , we have

$$\left| \ell(\hat{\theta}(X_i), X_i, Y_i) - \ell(\theta^\circ(X_i), X_i, Y_i) \right| \leq M(X_i, Y_i) \|\hat{\theta}(X_i) - \theta^\circ(X_i)\|_2,$$

hence (implicitly conditioning on other folds) we have

$$\begin{aligned} & \mathbb{E} \left[ \left( \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} r(X_i) \{ \ell(\hat{\theta}(X_i), X_i, Y_i) - \ell(\theta^\circ(X_i), X_i, Y_i) \} \right)^2 \right] \\ & \leq \mathbb{E} \left[ r(X)^2 M(X, Y(1))^2 \|\hat{\theta}(X_i) - \theta^\circ(X_i)\|_2^2 \right] = o_P(1) \end{aligned}$$

since  $\mathbb{E}[M(X, Y(1))^2 \mid T=1] \leq M$  for some constant  $M > 0$ . We've thus completed the proof of (6).

**Consistent  $\hat{h}$ .** We then show the results when  $\hat{h}$  is consistent, but not necessarily  $\hat{r}$ . In this case,  $\|\hat{h} - \bar{h}\|_{L_2(\mathbb{P}_{X|T=1})} = o_P(1)$ , where  $\bar{h} = \mathbb{E}[\hat{H}(X, Y(1)) | X = x, T = 1]$  viewing  $\hat{H}$  as fixed. Note that

$$\hat{\mu} = \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} \hat{r}(X_i) (\hat{H}(X_i, Y_i) - \bar{h}(X_i)) + \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} \hat{r}(X_i) (\bar{h}(X_i) - \hat{h}(X_i)) + \frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} \hat{h}(X_i).$$

The first summation is unbiased conditional on other folds, hence

$$\mathbb{E} \left[ \left( \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} \hat{r}(X_i) (\hat{H}(X_i, Y_i) - \bar{h}(X_i)) \right)^2 \right] = \frac{1}{|\mathcal{I}_1|} \text{Var}(\hat{r}(X) \{ \hat{H}(X, Y(1)) - \bar{h}(X) \} | T = 1) = o_P(1)$$

due to the finite second moments. By Cauchy-Schwarz inequality, the second summation satisfies

$$\mathbb{E} \left[ \left( \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} (\bar{h}(X_i) - \hat{h}(X_i)) \right)^2 \right] \leq \|\hat{r} \cdot (\hat{h} - \bar{h})\|_{L_2(\mathbb{P}_{X|T=1})}^2 = o_P(1)$$

due to the boundedness of  $\hat{r}$ . Similarly, we know that  $\frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} \bar{h}(X_i) - \hat{h}(X_i) = o_P(1)$ . Consequently, invoking Lemma E.1 we drop the implicit conditioning and arrive at

$$\hat{\mu} = \frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} \bar{h}(X_i) + o_P(1) = \mathbb{E}[\bar{h}(X) | T = 0] + o_P(1)$$

further using the finite second moment of  $\bar{h}$  (or  $\hat{H}$ ), where we implicitly condition on other folds and view  $\hat{h}$  as fixed. Finally, denoting  $h^\diamond(x) = \mathbb{E}[\ell(\theta^\diamond(x), x, Y(1)) | X = x, T = 1]$ , we note that by Jensen's inequality,

$$\begin{aligned} \left| \mathbb{E}[\bar{h}(X) | T = 0] - \mathbb{E}[h^\diamond(X) | T = 0] \right|^2 &\leq \|\bar{h} - h^\diamond\|_{L_2(\mathbb{P}_{X|T=0})}^2 \\ &\leq \|\hat{H}(X, Y(1)) - \ell(\theta^\diamond(X), X, Y(1))\|_{L_2(\mathbb{P}_{\cdot|T=0})}^2. \end{aligned}$$

By the same argument as the previous case and due to the uniform boundedness of the covariate shift  $r_{1,0}(\cdot)$ , the above term is  $o_P(1)$ . Therefore, by the change-of-measure with  $r_{1,0}$ , we have

$$\hat{\mu} = \mathbb{E}[h^\diamond(X) | T = 0] + o_P(1) = \mathbb{E}[r(X)h^\diamond(X) | T = 1] + o_P(1) = \mathbb{E}[r(X)\ell(\theta^\diamond(X), X, Y(1)) | T = 1] + o_P(1)$$

by the tower property of conditional expectations. We thus complete the proof of two cases and conclude the proof of Theorem 2.  $\square$

### C.7 Proof of Theorem 3

*Proof.* Proof of Theorem 3 We show that for each  $j$ , we have  $\hat{\mu}_{1,0}^{(j)} = \hat{\mu}_{1,0}^{*,(j)} + o_P(1/\sqrt{n})$ , where

$$\hat{\mu}_{1,0}^{*(j)} = \frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} r_{1,0}(X_i) (H(X_i, Y_i) - h(X_i)) + \frac{1}{|\mathcal{I}_0^{(j)}|} \sum_{i \in \mathcal{I}_0^{(j)}} h(X_i),$$

and we define

$$H(x, y) = \alpha^*(x) f^* \left( \frac{y + \eta^*(x)}{-\alpha^*(x)} \right) + \eta^*(x) + \alpha^*(x) \rho, \quad h(x) = \mathbb{E}[H(X, Y(1)) | X = x, T = 1].$$

We show this result for any  $j$ ; we implicitly condition on all the remaining folds other than  $\mathcal{I}_1^{(j)}$  and  $\mathcal{I}_0^{(j)}$ , so that all nuisance components are viewed as fixed. To simplify notations, we write  $\mathcal{I}_1 := \mathcal{I}_1^{(j)}$ ,  $\mathcal{I}_0 := \mathcal{I}_0^{(j)}$  and  $r := r_{1,0}$ ,  $\hat{r} := \hat{r}^{(j)}$ ,  $\hat{h} := \hat{h}^{(j)}$ ,  $\hat{H} := \hat{H}^{(j)}$ ,  $\bar{h} := \bar{h}^{(j)}$ . We also represent the parameters (functionals) with

$$\hat{\theta}(\cdot) = (\hat{\alpha}(\cdot), \hat{\eta}(\cdot)) := (\hat{\alpha}^{(j)}(\cdot), \hat{\eta}^{(j)}(\cdot)), \quad \theta^*(\cdot) := (\alpha^*(\cdot), \eta^*(\cdot)),$$

and recall the generic function (where  $\theta = (\alpha(\cdot), \eta(\cdot))$ )

$$\ell(\theta, x, y) = \alpha(x) f^* \left( \frac{y + \eta(x)}{-\alpha(x)} \right) + \eta(x) + \alpha(x) \rho,$$

so that  $H(x, y) = \ell(\theta^*, x, y)$  and  $\hat{H}(x, y) = \ell(\hat{\theta}, x, y)$ . By definition, we have the decomposition

$$\begin{aligned} \hat{\mu}_{1,0}^{(j)} - \hat{\mu}_{1,0}^{*(j)} &= \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} \left[ \hat{r}(X_i) (\hat{H}(X_i, Y_i) - \hat{h}(X_i)) - r(X_i) (H(X_i, Y_i) - h(X_i)) \right] + \frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} (\hat{h}(X_i) - h(X_i)) \\ &= \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} r(X_i) (\hat{H}(X_i, Y_i) - H(X_i, Y_i)) - \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r(X_i)) (\hat{h}(X_i) - \bar{h}(X_i)) \\ &\quad + \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r(X_i)) (\hat{H}(X_i, Y_i) - \bar{h}(X_i)) \\ &\quad - \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} r(X_i) (\hat{h}(X_i) - h(X_i)) + \frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} (\hat{h}(X_i) - h(X_i)). \end{aligned}$$

In the following, we are to bound the several summations separately. Firstly, by Cauchy-Schwarz inequality,

$$\begin{aligned} \left| \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r(X_i)) (\hat{h}(X_i) - \bar{h}(X_i)) \right| &\leq \sqrt{\frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r(X_i))^2} \sqrt{\frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} (\hat{h}(X_i) - \bar{h}(X_i))^2} \\ &= O_P(\|\hat{r} - r\|_{L_2(\mathbb{P}_{X|T=1})} \cdot \|\hat{h} - h\|_{L_2(\mathbb{P}_{X|T=1})}) = o_P(1/\sqrt{n}) \end{aligned}$$

under the given convergence rate of the product. Since  $\bar{h}(x) = \mathbb{E}[\hat{H}(X, Y(1)) | X = x, T = 1]$  for the fixed function  $\hat{H}$ , the term  $(\hat{r}(X_i) - r(X_i)) (\hat{H}(X_i, Y_i) - \bar{h}(X_i))$  has mean zero, hence by Markov's inequality,

$$\frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r(X_i)) (\hat{h}(X_i) - \bar{h}(X_i)) = O_P(\sqrt{\text{Var}(\hat{r}(X_i) - r(X_i)) (\hat{H}(X_i, Y_i) - \bar{h}(X_i))} / \sqrt{n}),$$

where, by the consistency of  $\hat{r}$ , this term is  $o_P(1/\sqrt{n})$ . Furthermore, note that

$$\begin{aligned} &\frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} r(X_i) (\hat{h}(X_i) - h(X_i)) - \frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} (\hat{h}(X_i) - h(X_i)) \\ &= \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} \left( r(X_i) (\hat{h}(X_i) - h(X_i)) - \mathbb{E}[r(X) (\hat{h}(X) - h(X)) | T = 1] \right) \\ &\quad - \frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} (\hat{h}(X_i) - h(X_i) - \mathbb{E}[\hat{h}(X_i) - h(X_i) | T = 0]), \end{aligned} \tag{7}$$

where we use the equivalence of the two expectations: this is because there is a covariate shift  $r(X)$  from  $\mathbb{P}_{X|T=1}$  to  $\mathbb{P}_{X|T=0}$ , hence  $\mathbb{E}[\phi(X)r(X) | T = 1] = \mathbb{E}[\phi(X) | X = 0]$  for any integrable function  $\phi: \mathcal{X} \rightarrow \mathbb{R}$ . The two summations in (7) is thus both unbiased, indicating

$$\begin{aligned} &\frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} r(X_i) (\hat{h}(X_i) - h(X_i)) - \frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} (\hat{h}(X_i) - h(X_i)) \\ &= O_P(\sqrt{\text{Var}(r(X) (\hat{h}(X) - h(X)) | T = 1) / n} + \sqrt{\text{Var}(r(X) (\hat{h}(X) - h(X)) | T = 0) / n}) \\ &= O_P(\|r(X) (\hat{h}(X) - h(X))\|_{L_2(\mathbb{P}_{X|T=1})} / \sqrt{n} + \|\hat{h} - h\|_{L_2(\mathbb{P}_{X|T=0})} / \sqrt{n}) = o_P(1/\sqrt{n}), \end{aligned}$$

where the last equality follows from the  $L_2$ -consistency of  $\hat{h}$  to  $\bar{h}$  and the fact that  $\|\bar{h} - h\|_{L_2(\mathbb{P}_{X|T=1})}$  by the stability of the conditional expectations induced by the stability of  $g$  in Assumption 4. Finally, we turn to

$$\frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} r(X_i) (\hat{H}(X_i, Y_i) - H(X_i, Y_i)).$$

Since  $\ell(\theta, x, y)$  is a convex function in  $\theta$  for any  $(x, y)$ , for any ( $\mathbb{P}_{X|T=1}$ -almost all)  $x$ ,

$$\mathbb{E}[\ell(\theta, x, y) | X = x, T = 1] = \alpha(x)\mathbb{E}\left[f^*\left(\frac{Y(1) + \eta(x)}{-\alpha(x)}\right) \middle| X = x, T = 1\right] + \eta(x) + \alpha(x)\rho$$

is also convex and differentiable by the given regularity condition. In particular, by the optimality of  $(\alpha^*(x), \eta^*(x))$  for the per- $x$  minimization problem and the exchangeability of differentiation and expectation,

$$\nabla_{\theta}\mathbb{E}[\ell(\theta, x, Y(1)) | X = x, T = 1] \Big|_{\theta=(\alpha^*(x), \eta^*(x))} = \mathbb{E}[\nabla_{\theta}\ell(\theta^*(x), x, Y(1)) | X = x, T = 1] = 0.$$

Multiplying  $r(X)$  and integrating over  $X | T = 1$ , we know that

$$\mathbb{E}[r(X)\nabla_{\theta}\ell(\theta^*(X), X, Y(1))[\widehat{\theta}(X) - \theta^*(X)] | T = 1] = 0. \quad (8)$$

By Lemma 2 of [Chen and Shen \(1998\)](#), when both  $\widehat{\theta}$  and  $\theta^*$  is smooth enough, sufficiently small  $\|\widehat{\theta} - \theta^*\|_{L_2(\mathbb{P})}$  implies sufficiently small  $\|\widehat{\theta} - \theta^*\|_{\infty}$ . As a result, when  $\|\widehat{\theta}(x) - \theta^*(x)\|_{L_2(\mathbb{P}_{\cdot|T=1})}$  is sufficiently small, by the condition that  $|\ell(\widehat{\theta}, x, y) - \ell(\theta^*, x, y) - \nabla_{\theta}\ell(\theta^*(x), x, y)[\theta^*(x) - \widehat{\theta}(x)]| \leq M(x, y)\|\widehat{\theta}(x) - \theta^*(x)\|_2^2$  as well as Jensen's inequality, we have

$$\begin{aligned} & \left| \mathbb{E}[\widehat{H}(X, Y(1)) - H(X, Y(1)) | T = 1] \right| \\ &= \left| \mathbb{E}\left[r(X)\left\{\ell(\widehat{\theta}, X, Y(1)) - \ell(\theta^*, X, Y(1)) - \nabla_{\theta}\ell(\theta^*(X), X, Y(1))[\widehat{\theta}(X) - \theta^*(X)]\right\} \middle| T = 1\right] \right| \\ &\leq \mathbb{E}\left[r(X)\left|\ell(\widehat{\theta}, X, Y(1)) - \ell(\theta^*, X, Y(1)) - \nabla_{\theta}\ell(\theta^*(X), X, Y(1))[\widehat{\theta}(X) - \theta^*(X)]\right| \middle| T = 1\right] \\ &\leq \mathbb{E}[r(X)M(X, Y(1))\|\widehat{\theta}(X) - \theta^*(X)\|_2^2 | T = 1] = M \cdot \|r(\widehat{\theta} - \theta)\|_{L_2(\mathbb{P}_{X|T=0})}^2. \end{aligned}$$

Returning to our problem, we note that due to unbiasedness,

$$\begin{aligned} & \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} \left[ r(X_i)(\widehat{H}(X_i, Y_i) - H(X_i, Y_i)) - \mathbb{E}[\widehat{H}(X, Y(1)) - H(X, Y(1)) | T = 1] \right] \\ &= o_P(\|r(X)\widehat{H}(X, Y(1)) - H(X, Y(1))\|_{L_2(\mathbb{P}_{X|T=1})}/\sqrt{n}), \end{aligned}$$

where by the given conditions, we have

$$\|r(X)\widehat{H}(X, Y(1)) - H(X, Y(1))\|_{L_2(\mathbb{P}_{X|T=1})} = O(\|\widehat{\theta} - \theta^*\|_{L_2(\mathbb{P}_{X|T=1})}) = o_P(1).$$

As a result, we have

$$\frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} r(X_i)(\widehat{H}(X_i, Y_i) - H(X_i, Y_i)) \leq o_P(1/\sqrt{n}) + M\|\widehat{\theta} - \theta\|_{L_2(\mathbb{P}_{X|T=0})}^2 = o_P(1/\sqrt{n}).$$

Putting these pieces together, we conclude the proof of  $\widehat{\mu}_{1,0}^{(j)} = \widehat{\mu}_{1,0}^{*,(j)} + o_P(1/\sqrt{n})$  for each  $j$ . Therefore, averaging over the three folds, we have

$$\sqrt{n}(\widehat{\mu}_{1,0}^- - \mu_{1,0}^-) = \frac{\sqrt{n}}{n_1} \sum_{T_i=1} \left( r_{1,0}(X_i)(H(X_i, Y_i) - h(X_i)) - \mu_{1,0}^- \right) + \frac{\sqrt{n}}{n_0} \sum_{T_i=0} h(X_i) + o_P(1/\sqrt{n}),$$

which, by CLT and Slutsky's theorem, converges in distribution to  $N(0, \sigma^2)$ . Here  $n_1$  is the total number of treated samples, and  $n_0$  is the number of control samples. The asymptotic variance is

$$\sigma^2 = \frac{1}{p_1} \text{Var}\left(r_{1,0}(X)(H(X, Y(1)) - h(X)) \middle| T = 1\right) + \frac{1}{p_0} \text{Var}(h(X) | T = 0).$$

where  $p_1 = \mathbb{P}(T = 1)$ ,  $p_0 = \mathbb{P}(T = 0)$  and all the expectations (variances) are induced by the observed distribution.

It now remains to show that  $\hat{\sigma}^2 \rightarrow \sigma^2$  in the definition of Theorem 3. As  $\hat{p}_1 \xrightarrow{P} p_1$ ,  $\hat{p}_0 \xrightarrow{P} p_0$ , by the law of large numbers, it suffices to show that  $\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} (d_{1,i}^2 - (d_{1,i}^*)^2) = o_P(1)$  and  $\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} (d_{1,i} - d_{1,i}^*) = o_P(1)$  and similar for  $(d_{0,i}, d_{0,i}^*)$ , where we define the oracle counterparts

$$d_{1,i}^* = r_{1,0}(X_i)(H(X_i, Y_i) - h(X_i)), \quad d_{0,i}^* = h(X_i).$$

By Cauchy-Schwarz inequality,

$$\begin{aligned} \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} (d_{1,i}^2 - (d_{1,i}^*)^2) &= \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} (d_{1,i} - d_{1,i}^*)^2 + \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} 2(d_{1,i} - d_{1,i}^*)d_{1,i}^* \\ &\leq \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} (d_{1,i} - d_{1,i}^*)^2 + \sqrt{\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} (d_{1,i} - d_{1,i}^*)^2} \cdot \sqrt{\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} (d_{1,i}^*)^2}. \end{aligned}$$

Focusing on the summation within  $\mathcal{I}_1^{(j)}$ , we have

$$\frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} (d_{1,i} - d_{1,i}^*)^2 = O_P\left(\|\hat{r}^{(j)}(\hat{H}^{(j)} - \hat{h}^{(j)}) - r_{1,0}(H - h)\|_{L_2(\mathbb{P}_{\cdot|T=1})}^2\right),$$

where the right-handed side is  $o_P(1)$  under the conditions of Theorem 3. Other folds and other summation terms follow similar arguments hence  $\hat{\sigma}^2 \xrightarrow{P} \sigma^2$ . By Slutsky's lemma, we conclude the proof of Theorem 3.  $\square$

## C.8 Proof of Theorem 4

*Proof.* Proof of Theorem 4 The proof follows exactly the same arguments as the proof of Theorem 3 with  $\theta^\diamond$  in place of  $\theta^*$ , where all the errors are controlled in the same way; the only difference is to show that

$$\mathbb{E}[r(X)\nabla_{\theta}\ell(\theta^\diamond(X), X, Y(1))[\hat{\theta}(X) - \theta^\diamond(X)] \mid T = 1]$$

in parallel with (8) in the proof of Theorem 3. We note that this is directly implied by our local condition. Therefore, under the conditions of Theorem 4,

$$\hat{\mu}_{1,0}^- = -\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} r_{1,0}(X_i)[H^\diamond(X_i, Y_i(1)) - h^\diamond(X_i)] - \frac{1}{n_0} \sum_{i \in \mathcal{I}_0} h^\diamond(X_i).$$

The terms in the summation has expectation

$$\mu_{1,0}^\diamond := -\mathbb{E}[r_{1,0}(X_i)H^\diamond(X_i, Y_i(1)) \mid T = 1].$$

Since  $\alpha^*(x), \eta^*(x)$  is the per- $x$  minimizer of  $\mathbb{E}[\ell(\theta, X, Y(1)) \mid X = x, T = 1]$ , we have

$$\mathbb{E}[H^\diamond(X_i, Y_i(1)) \mid X = x, T = 1] \geq \mathbb{E}[H(X_i, Y_i(1)) \mid X = x, T = 1]$$

for  $\mathbb{P}_{X|T=1}$ -almost all  $x$ , hence by tower property, we have  $\mu_{1,0}^\diamond \leq \mu_{1,0}^-$ . On the other hand, the consistency of  $\hat{\sigma}^2$  to  $\text{Var}(\phi_{1,-}^\diamond(X, Y, T))$  also follows the same arguments as the proof of Theorem 3 with  $\theta^\diamond$ , which concludes our proof of Theorem 4.  $\square$

## D Technical proofs for marginal $f$ -sensitivity models

### D.1 Proof of Proposition B.1

*Proof.* Proof of Proposition B.1 We prove the claim by showing that  $\mathcal{Q}_{1,0}^m$  and the set on the right-handed side (RHS) of Proposition 1 contain each other.

**Step 1:**  $\mathcal{Q}_{1,0}^m \subseteq \text{RHS}$ .

By definition, any member of  $\mathcal{Q}_{1,0}^m$  is the induced distribution over  $(X, Y(1))$  given  $T = 0$  of some super-population  $\mathbb{Q} \in \mathcal{M}_{1,\text{mgn}}^{f,\rho}$  with  $\mathbb{Q}_{X,Y,T} = \mathbb{P}_{X,Y,T}$ . Thus, we write it as  $\mathbb{Q}_{X,Y(1)|T=0} \in \mathcal{Q}_{1,0}^m$ . Since  $\mathbb{Q}_{X,Y,T} = \mathbb{P}_{X,Y,T}$ , we know that  $\frac{d\mathbb{Q}_{X|T=0}}{d\mathbb{P}_{X|T=1}} = \frac{d\mathbb{P}_{X|T=0}}{d\mathbb{P}_{X|T=1}} = r_{1,0}(x)$ . Following the proof of Lemma 1,

$$\frac{d\mathbb{P}_{Y(1),U|X=x,T=0}}{d\mathbb{P}_{Y(1),U|X=x,T=1}} = \frac{\mathbb{P}(T=0|X=x,U)}{\mathbb{P}(T=1|X=x,U)} \cdot \frac{\mathbb{P}(T=1|X=x)}{\mathbb{P}(T=0|X=x)} = \text{OR}(x,U),$$

as well as (by condition (1) and the data-processing inequality)

$$\begin{aligned} D_f(\mathbb{P}_{Y(1)|X=x,T=0} \parallel \mathbb{P}_{Y(1)|X=x,T=1}) &\leq D_f(\mathbb{P}_{Y(1),U|X=x,T=0} \parallel \mathbb{P}_{Y(1),U|X=x,T=1}) \\ &= \mathbb{E}_{Y(1),U|X=x,T=1} \left[ f \left( \frac{d\mathbb{P}_{Y(1),U|X=x,T=0}}{d\mathbb{P}_{Y(1),U|X=x,T=1}} \right) \right]. \end{aligned}$$

Combining the above two facts and using the tower property,

$$\mathbb{E} \left[ D_f(\mathbb{P}_{Y(1)|X,T=0} \parallel \mathbb{P}_{Y(1)|X,T=1}) \mid T=1 \right] \leq \mathbb{E}_{Y(1),U|T=1} \left[ f(\text{OR}(X,U)) \right] \leq \rho.$$

Also, note that  $\mathbb{Q}_{Y(1)|X,T=1} = \mathbb{Q}_{Y|X,T=1} = \mathbb{P}_{Y|X,T=1}$  since  $Y = TY(1) + (1-T)Y(0)$  and  $\mathbb{Q}_{X,Y,T} = \mathbb{P}_{X,Y,T}$ . These altogether establish the ‘‘C’’ direction. It remains to prove the reverse.

**Step 2: RHS  $\subseteq \mathcal{Q}_{1,0}^m$ .**

Given any distribution  $Q$  over  $(X, Y(1))$  in RHS, we aim to find a distribution (super-population)  $\mathbb{Q}$  over  $(X, Y(0), Y(1), U, T)$  such that

- $(Y(1), Y(0)) \perp\!\!\!\perp T \mid X, U$  under  $\mathbb{Q}$ ;
- $\mathbb{Q}_{X,T,Y} = \mathbb{P}_{X,T,Y}$ ;
- $\mathbb{Q} \in \mathcal{M}_{1,\text{mgn}}^{f,\rho}$ ;
- $\mathbb{Q}_{X,Y(1)|T=0}(x, y) = Q(x, y)$ .

These conditions imply that  $Q \in \mathcal{Q}_{1,0}^m$ , and thus  $\text{RHS} \subseteq \mathcal{Q}_{1,0}^m$ . To construct  $\mathbb{Q}$ , we first set  $\mathbb{Q}_{X,T} = \mathbb{P}_{X,T}$ . Then we specify the distribution of  $Y(1) \mid X, T$  via

$$\mathbb{Q}_{Y(1)|X,T=1} = \mathbb{P}_{Y|X,T=1}, \quad \mathbb{Q}_{Y(1)|X,T=0} = \mathbb{Q}_{Y|X}.$$

So far,  $\mathbb{Q}_{X,T,Y(1)}$  has been determined. We let  $U = Y(1)$  be the unobserved confounder. Finally, the distribution of  $Y(0) \mid X, T, Y(1), U$  is specified via

$$\mathbb{Q}_{Y(0)|X,T,Y(1),U} = \mathbb{Q}_{Y(0)|X} = \mathbb{P}_{Y|X,T=0}.$$

Having constructed  $\mathbb{Q}$ , we proceed to check that it satisfies the conditions listed at the beginning of Step 2. Since  $Y(1) = U$  and  $\mathbb{Q}_{Y(0)|X,T,Y(1),U} = \mathbb{Q}_{Y(0)|X}$  by construction, it is easy to check that  $(Y(1), Y(0)) \perp\!\!\!\perp T \mid X, U$ . By construction, it also clearly holds that  $\mathbb{Q}_{X,T,Y} = \mathbb{P}_{X,T,Y}$ .

We now check that  $\mathbb{Q} \in \mathcal{M}_{1,\text{mgn}}^{f,\rho}$ . For any  $x$ , again by the construction of  $\mathbb{Q}$ ,

$$\begin{aligned} \mathbb{E}_{\mathbb{Q}_{X,Y(1)|T=1}} \left[ f \left( \frac{d\mathbb{Q}_{Y(1)|X=x,T=0}}{d\mathbb{Q}_{Y(1)|X=x,T=1}} \right) \right] &= \mathbb{E}_{\mathbb{P}_{X,Y|T=1}} \left[ f \left( \frac{dQ_{Y|X=x}}{d\mathbb{P}_{Y|X=x,T=1}} \right) \right] \\ &= \mathbb{E} \left[ D_f(Q_{Y|X=x} \parallel \mathbb{P}_{Y|X=x,T=1}) \mid T=1 \right] \leq \rho, \end{aligned}$$

where the last inequality is due to the fact that  $Q$  is in the set on the RHS. Therefore, we verify that  $\mathbb{Q} \in \mathcal{M}_{1,\text{mgn}}^{f,\rho}$ . By construction,  $\mathbb{Q}_{Y(1)|X,T=0} = \mathbb{Q}_{Y|X}$ . It remains to show that  $\mathbb{Q}_{X|T=0} = Q_X$ . For any measurable set  $A$ ,

$$\begin{aligned} \mathbb{Q}(X \in A \mid T=0) &= \mathbb{E}_{\mathbb{Q}} \left[ \frac{d\mathbb{Q}_{X|T=0}}{d\mathbb{Q}_{X|T=1}} \cdot \mathbf{1}\{X \in A\} \mid T=1 \right] = \mathbb{E}_{\mathbb{Q}} \left[ \frac{dQ_X|T=0}{d\mathbb{Q}_{X|T=1}} \cdot \mathbf{1}\{X \in A\} \mid T=1 \right] \\ &= \mathbb{E}_{\mathbb{Q}} \left[ r_{1,0}(X) \cdot \mathbf{1}\{X \in A\} \mid T=1 \right] = \mathbb{E}_{\mathbb{Q}} \left[ \frac{dQ_X}{d\mathbb{Q}_{X|T=1}} \cdot \mathbf{1}\{X \in A\} \mid T=1 \right] = Q(X \in A). \end{aligned}$$

Since the above holds for any measurable set  $A$ , we know that  $\mathbb{Q}_{X|T=0} = Q_X$ . Finally, switching the role of the treated and control groups yields similar results for the counterfactual distribution of  $(X, Y(0))$  given  $T = 1$  under  $\mathcal{M}_{0,\text{mgn}}^{f,\rho}$ . This completes the proof of Proposition 1.  $\square$

## D.2 Proof of Proposition B.3

*Proof.* Proof of Proposition B.3 The problem (5) is also convex. For notational simplicity, throughout the proof we write  $\mathbb{E}[\cdot]$  instead of  $\mathbb{E}[\cdot | T = 1]$ , and  $\mathbb{E}[\cdot | X = x]$  instead of  $\mathbb{E}[\cdot | X = x, T = 1]$ . According to Luenberger (1997, Theorem 8.6.1), the optimal value of the objective function in (5) can be written as

$$\min_{\substack{\mathbb{E}[L | X=x]=r_{1,0}(x), \forall x \\ \mathbb{E}[f(L/r_{1,0}(X))] - \rho \leq 0}} \mathbb{E}[Y(1)L] = \max_{\alpha \geq 0} \varphi(\alpha),$$

where the dual functional is given by (where  $L$  is essentially a measurable function of  $(X, Y(1))$ )

$$\begin{aligned} \varphi(\alpha) &= \min_{\mathbb{E}[L | X=x]=r_{1,0}(x), \forall x} \mathcal{L}(\alpha, L), \\ \mathcal{L}(\alpha, L) &= \mathbb{E}[Y(1)L] + \alpha(\mathbb{E}[f(L/r_{1,0}(X))] - \rho). \end{aligned} \quad (9)$$

**Claim D.1.** *We claim that  $\varphi(\alpha) = \mathbb{E}[\varphi(\alpha, X)]$ , where  $\varphi(\alpha, x)$  is the optimal objective of*

$$\begin{aligned} \min \quad & \mathbb{E}[Y(1)L | X = x] + \alpha(\mathbb{E}[f(L/r_{1,0}(x)) | X = x] - \rho) \\ \text{s.t.} \quad & \mathbb{E}[L | X = x] = r_{1,0}(x). \end{aligned} \quad (10)$$

*Proof.* Proof of Claim D.1 The claim follows similar arguments as the Neyman-Pearson Lemma. To see why it is the case, let  $L^*(x, \cdot)$ , a function measurable with respect to  $Y(1)$ , be the optimizer of the per- $x$  optimization problem (10). Then  $L^*$ , as a function of both  $X$  and  $Y(1)$ , satisfies the constraint of (9), hence  $\varphi(\alpha) \leq \mathcal{L}(\alpha, L^*) = \mathbb{E}[\mathbb{E}[\mathcal{L}(\alpha, L^*) | X]] = \mathbb{E}[\varphi(\alpha, X)]$ . On the other hand, for every measurable function  $L(\cdot, \cdot)$  that satisfies the constraint of (9), if we define  $L_x = L(x, \cdot)$ , then  $L_x$  satisfies the constraint of (10), hence

$$\mathbb{E}[Y(1)L_x(Y(1)) | X = x] + \alpha(\mathbb{E}[f(L_x/r_{1,0}(x)) | X = x] - \rho) \geq \varphi(\alpha, x).$$

Integrating over  $X$ , we have  $\mathcal{L}(\alpha, L) \geq \mathbb{E}[\varphi(\alpha, X)]$ . Minimizing over such  $L$  we get  $\varphi(\alpha) \geq \mathbb{E}[\varphi(\alpha, X)]$ . By the two inequalities we conclude the proof of Claim D.1.  $\square$

Again invoking Luenberger (1997, Theorem 8.6.1), we can further solve for  $\varphi(\alpha, x)$  via its dual

$$\begin{aligned} \varphi(\alpha, x) &= \max_{\eta \in \mathbb{R}} \phi(\alpha, \eta, x) \\ &= \max_{\eta \in \mathbb{R}} \min_{L \geq 0} \left\{ \mathbb{E}[Y(1)L | X = x] + \alpha(\mathbb{E}[f(L/r_{1,0}(x)) | X = x] - \rho) + \eta \mathbb{E}[L - r_{1,0}(x) | X = x] \right\} \\ &= \max_{\eta \in \mathbb{R}} \min_{L \geq 0} \left\{ \mathbb{E}[Y(1)L + \alpha f(L/r_{1,0}(x)) + \alpha \eta L | X = x] - \alpha \rho - \eta r_{1,0}(x) \right\}. \end{aligned}$$

Similar to the idea of Claim D.1, the inner minimization problem can also be solved inside the expectation, which leads to

$$\begin{aligned} \varphi(\alpha, x) &= \max_{\eta \in \mathbb{R}} \left\{ \mathbb{E} \left[ \min_{L \geq 0} \{ Y(1)L + \alpha f(L/r_{1,0}(x)) + \alpha \eta L \} \mid X = x \right] - \alpha \rho - \eta r_{1,0}(x) \right\} \\ &= \max_{\eta \in \mathbb{R}} \left\{ -\alpha \mathbb{E} \left[ f^* \left( -\frac{r_{1,0}(x)}{\alpha} (Y(1) + \eta) \right) \mid X = x \right] - \alpha \rho - \eta r_{1,0}(x) \right\}, \end{aligned}$$

where  $f^*(s) = \sup_{t \geq 0} \{st - f(t)\}$  is the conjugate function. The solution to (2) is thus given by (3). This completes the proof of Proposition B.3.  $\square$

## D.3 Proof of sieve estimation

*Proof.* Proof of Proposition B.4 To simplify notations, sometimes we write  $\mathbb{E}[\cdot]$  instead of  $\mathbb{E}[\cdot | T = 1]$ ,  $\mathbb{E}[\cdot | X = x]$  instead of  $\mathbb{E}[\cdot | X = x, T = 1]$ ,  $\mathbb{P}(\cdot)$  instead of  $\mathbb{P}(\cdot | T = 1)$  and  $\mathbb{P}(\cdot | X = x)$  instead of  $\mathbb{P}(\cdot | X = x, T = 1)$ . We will invoke Lemma C.1 again and verify the conditions therein.

To be in line with the notations of Lemma C.1, recall that viewing  $\hat{r}$  as fixed, the risk function is  $\ell(\theta, x, y) = \bar{m}(x, y, \theta)$ , and the population risk minimizer is  $\theta^* = \bar{\theta} = \bar{\alpha}, \bar{\eta}(\cdot)$ . We define the  $L_2$ -norm  $d$  on

$\Theta = \mathbb{R}^+ \times \Lambda_c^p(\mathcal{X})$  as  $d(\theta, \theta') = \|\theta - \theta'\|_{L_2(\mathbb{P})} := ((\alpha - \alpha')^2 + \|\eta' - \eta\|_{L_2(\mathbb{P}, \cdot|_{T=1})}^2)^{1/2}$ . By the positive density condition, we have  $\|\cdot\|_{L_2(\lambda)} \asymp \|\cdot\|_{L_2(\mathbb{P})}$ . Hence for  $\theta, \bar{\theta} \in \Theta$ , we have  $\|\theta - \bar{\theta}\|_\infty = o(1)$  once  $\|\theta - \bar{\theta}\|_{L_2(\mathbb{P})} = o(1)$ . By Lemma 2 of [Chen and Shen \(1998\)](#), we have  $\|\theta\|_\infty \lesssim \|\theta\|_{L_2(\lambda)}^{2p/(2p+d)}$  for any  $\theta \in \Theta$ , where  $\lambda$  is the Lebesgue measure. Therefore, sufficiently small  $d(\theta - \bar{\theta})$  implies sufficiently small  $\|\theta - \bar{\theta}\|_\infty$ .

We first verify the equivalence of  $d(\theta, \bar{\theta})^2$  and  $\mathbb{E}[\ell(\theta, X, Y) - \ell(\theta^*, X, Y)]$ . By the optimality of  $\bar{\alpha}, \bar{\eta}(\cdot)$ , we know that  $\bar{\eta}(x)$  is also the per- $x$  optimizer of the risk function at  $\bar{\alpha}$ , that is, for  $\mathbb{P}_X|_{T=1}$ -almost all  $x$ ,

$$\bar{\eta}(x) = \arg \min_{\eta \in \mathbb{R}} \mathbb{E} \left[ \bar{m}(x, Y(1), \bar{\alpha}, \eta) \mid X = x, T = 1 \right].$$

Suppose  $d(\theta - \bar{\theta})$  is sufficiently small, which means  $\|\theta - \bar{\theta}\|_\infty < \epsilon$  for sufficiently small  $\epsilon > 0$ . Therefore, by the second order expansion of  $\bar{I}(x, a, b) = \mathbb{E}[\bar{m}(x, Y(1), a, b) \mid X = x, T = 1]$  in Assumption [B.1](#), we have

$$\left| \bar{I}(x, y, \alpha, \eta(x)) - \bar{I}(x, y, \bar{\alpha}, \bar{\eta}(x)) - \nabla_{a,b} \bar{I}(x, y, \bar{\alpha}, \bar{\eta}(x))[\alpha - \bar{\alpha}, \eta(x) - \bar{\eta}(x)] \right| \geq \lambda \|(\alpha - \bar{\alpha}, \eta(x) - \bar{\eta}(x))\|_2^2.$$

Since  $\mathbb{E}[\bar{I}(X, Y(1), \alpha, \eta(X)) \mid T = 1] = \ell(\theta, X, Y)$ , invoking Jensen's inequality we have

$$\left| \mathbb{E}[\ell(\theta, X, Y)] - \mathbb{E}[\ell(\bar{\theta}, X, Y)] - \mathbb{E}[\nabla_{a,b} \bar{I}(x, y, \bar{\alpha}, \bar{\eta}(X))[\alpha - \bar{\alpha}, \eta(X) - \bar{\eta}(X)]] \right| \geq \lambda \cdot d(\theta, \bar{\theta})^2. \quad (11)$$

Since  $\bar{m}$  is a convex function in  $(a, b)$  for all  $(x, y)$ , the conditional expectation is also convex, and by the differentiability and exchangeability of differentiation and expectation, we have

$$\mathbb{E} \left[ \nabla_b \bar{m}(x, Y(1), \bar{\alpha}, \bar{\eta}(x)) \mid X = x, T = 1 \right] = 0.$$

Similarly,  $\bar{\alpha}$  is also the minimizer of the convex risk given  $\bar{\eta}$ , and the regularity conditions imply

$$\mathbb{E} \left[ \nabla_a \bar{m}(X, Y(1), \bar{\alpha}, \bar{\eta}(X)) \mid T = 1 \right] = \nabla_a \mathbb{E} \left[ \bar{m}(X, Y(1), \bar{\alpha}, \bar{\eta}(X)) \mid T = 1 \right] = 0.$$

Returning to [\(11\)](#), we then have  $\mathbb{E}[\ell(\theta, X, Y) - \ell(\theta^*, X, Y)] \geq \lambda \cdot d(\theta, \bar{\theta})^2$  when  $\theta \in \Theta$  and  $d(\theta, \bar{\theta})$  is sufficiently small. On the other hand, when  $d(\theta - \bar{\theta})$  is sufficiently small, the second-order expansion of  $\bar{m}$  near  $(\bar{\alpha}, \bar{\eta}(x))$  in Assumption [B.1](#) implies

$$\begin{aligned} & \left| \mathbb{E}[\ell(\theta, X, Y)] - \mathbb{E}[\ell(\bar{\theta}, X, Y)] - \mathbb{E}[\nabla_{a,b} \bar{m}(X, Y, \bar{\alpha}, \bar{\eta}(X))[a - \bar{\alpha}, b - \bar{\eta}(X)]] \right| \\ & \leq \mathbb{E}[g(X, Y) \|(a - \bar{\alpha}, b - \bar{\eta}(X))\|_2^2]. \end{aligned}$$

Here similar to previous arguments,  $\mathbb{E}[\nabla_{a,b} \bar{m}(X, Y, \bar{\alpha}, \bar{\eta}(X))[a - \bar{\alpha}, b - \bar{\eta}(X)]] = 0$  by the exchangeability of differentiation and expectation; by the finite expectation of  $g(x, y)$  in Assumption [B.1](#), the right-hand side can be bounded as

$$\mathbb{E}[g(X, Y) \|(a - \bar{\alpha}, b - \bar{\eta}(X))\|_2^2] \leq \mathbb{E}[\mathbb{E}[g(X, Y) \mid X = x] \|(a - \bar{\alpha}, b - \bar{\eta}(X))\|_2^2] \leq M \cdot d(\theta, \bar{\theta})^2.$$

We have thus verified the equivalence of  $d(\theta, \bar{\theta})$  and  $|\mathbb{E}[\ell(\theta, X, Y)] - \mathbb{E}[\ell(\bar{\theta}, X, Y)]|$ .

We then proceed to verify condition (i) of Lemma [C.1](#). For  $\theta \in \Theta$  such that  $d(\theta, \bar{\theta})$  hence  $\|\theta - \bar{\theta}\|_\infty$  is sufficiently small, since  $\bar{m}(\theta, x, y) - \bar{m}(\bar{\theta}, x, y) \leq \bar{g}(x, y) \|\theta(x) - \bar{\theta}(x)\|_2$  where  $\mathbb{E}[\bar{g}(x, Y)^2 \mid X = x] \leq M$  for all  $x$ , we have

$$\begin{aligned} \text{Var}(\ell(\theta, X, Y) - \ell(\bar{\theta}, X, Y)) & \leq \mathbb{E}[|\bar{m}(\theta, X, Y) - \bar{m}(\bar{\theta}, X, Y)|^2] \leq M \mathbb{E}[(\theta(X) - \bar{\theta}(X))^2] \\ & \leq \mathbb{E}[\mathbb{E}[\bar{g}(X, Y)^2 \mid X] \cdot \|\theta(X) - \bar{\theta}(X)\|_2^2] \leq M \cdot d(\theta, \bar{\theta}), \end{aligned}$$

hence condition (i) is satisfied. For condition (ii), it suffices to let  $U_n(x, y) = \bar{g}(x, y)$ .

Therefore, applying Lemma [C.1](#), we know that  $d(\hat{\theta}_n - \bar{\theta}) = O(\max\{\delta_n, \inf_{\theta' \in \Theta_n} d(\theta', \bar{\theta})\})$ . Using a similar calculation as in the proof of Theorem [1](#), we have

$$\delta_n \asymp \sqrt{\frac{J_n^d \log n}{n}},$$

and the approximation error using  $\Theta_n$  is upper bounded as  $\inf_{\theta' \in \Theta_n} d(\theta', \bar{\theta}) \leq O(J_n^p)$ . We can thus set  $J_n = (n/\log n)^{1/(2p+d)}$ , so that  $\|\hat{\theta} - \bar{\theta}\|_{L_2(\mathbb{P})} = O_P((\log n/n)^{p/(2p+d)})$ , which completes our proof.  $\square$

## D.4 Proof of Theorem B.1

*Proof.* Proof of Theorem B.1 We are to show that, for each  $j \in \{1, 2, 3, 4\}$ ,  $\sqrt{n}(\widehat{\mu}_{1,0}^{m,(j)} - \widehat{\mu}_{1,0}^*) = o(1)$ , where we define

$$\widehat{\mu}_{1,0}^* = \frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} \left\{ m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i)) - \alpha^* r_{0,1}(X_i) h(X_i) \right\} + \frac{\alpha^*}{|\mathcal{I}_0^{(j)}|} \sum_{i \in \mathcal{I}_0^{(j)}} h(X_i).$$

Here we define  $h(x) = \mathbb{E}[H(x, Y(1)) | X = x, T = 1]$ , where

$$H(x, y) = \frac{\partial f^*}{\partial r}(x, y, \alpha^*, \eta^*(x), r_{0,1}(x)).$$

In the following, we suppress the dependency on  $j$  as the result applies to all  $j$ . We also implicitly condition on  $\mathcal{I}_1^{(j+1)} \cup \mathcal{I}_0^{(j+1)} \cup \mathcal{I}_1^{(j+2)} \cup \mathcal{I}_0^{(j+2)}$ , so that  $\widehat{\alpha}, \widehat{\eta}, \widehat{r}$  are viewed as fixed. Until further notice, we write  $\mathcal{I}_1 = \mathcal{I}_1^{(j)}$  and  $\mathcal{I}_0 = \mathcal{I}_0^{(j)}$  for notational simplicity, and write  $n_1 = |\mathcal{I}_1|$ ,  $n_0 = |\mathcal{I}_0|$ .

With similar arguments as in the proof of Proposition B.4, when  $\widehat{\eta}$  and  $\bar{\eta}$  are sufficiently smooth, sufficiently small  $\|\widehat{\theta} - \bar{\theta}\|_{L_2(\mathbb{P}_{X|T=1})}$  implies sufficiently small  $\|\widehat{\theta} - \bar{\theta}\|_\infty$ . By the second-order expansion at  $\bar{\theta}(x) = (\bar{\alpha}, \bar{\eta}(x))$  in Assumption B.3, when  $\|\widehat{\theta}(x) - \bar{\eta}(x)\|_2$  is sufficiently small, we have

$$\begin{aligned} & \left| \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} m(X_i, Y_i, \widehat{\alpha}, \widehat{\eta}(X_i), \widehat{r}(X_i)) - \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \widehat{r}(X_i)) \right. \\ & \quad \left. - \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} \nabla_{a,b} m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \widehat{r}(X_i)) [\widehat{\alpha} - \bar{\alpha}, \widehat{\eta}(X_i) - \bar{\eta}(X_i)] \right| \leq \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} g(X_i, Y_i) \|\widehat{\theta}(X_i) - \bar{\theta}(X_i)\|_2^2. \end{aligned}$$

Here by the optimality of  $\bar{\alpha}$  and  $\bar{\eta}(x)$ , we know that

$$\mathbb{E}[\nabla_a m(X_i, Y_i(1), \bar{\alpha}, \bar{\eta}(X_i), \widehat{r}(X_i)) | T_i = 1] = 0$$

and for  $\mathbb{P}_{X|T=1}$ -almost all  $x$ ,

$$\mathbb{E}[\nabla_b m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \widehat{r}(X_i)) | X_i = x, T_i = 1] = 0.$$

Consequently, we have

$$\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} \nabla_a m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \widehat{r}(X_i)) (\widehat{\alpha} - \bar{\alpha}) = O_p(|\widehat{\alpha} - \bar{\alpha}|/\sqrt{n}) = o_P(1/\sqrt{n})$$

and

$$\begin{aligned} & \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} \nabla_a m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \widehat{r}(X_i)) (\widehat{\eta}(X_i) - \bar{\eta}(X_i)) \\ & = O_P\left(\left\{ \text{Var}(\nabla_a m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \widehat{r}(X_i)) (\widehat{\eta}(X_i) - \bar{\eta}(X_i)) | T_i = 1) / n \right\}^{1/2}\right) = o(1/\sqrt{n}) \end{aligned}$$

since the latter, due to Assumption B.3, is further bounded as

$$\begin{aligned} & \text{Var}(\nabla_b m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \widehat{r}(X_i)) (\widehat{\eta}(X_i) - \bar{\eta}(X_i))) \\ & \leq \mathbb{E}\left[\mathbb{E}[\nabla_b m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \widehat{r}(X_i))^2 | X_i, T_i = 1] (\widehat{\eta}(X_i) - \bar{\eta}(X_i))^2 | T_i = 1\right] \leq M \|\widehat{\eta} - \bar{\eta}\|_{L_2(\mathbb{P}_{X|T=1})}^2. \end{aligned}$$

Furthermore, by Markov's inequality we have

$$\begin{aligned} & \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} g(X_i, Y_i) \|\widehat{\theta}(X_i) - \bar{\theta}(X_i)\|_2^2 \\ & = O_P(\mathbb{E}[g(X_i, Y_i) \|\widehat{\theta}(X_i) - \bar{\theta}(X_i)\|_2^2 | T_i = 1]) = O_P(\|\widehat{\theta} - \bar{\theta}\|_2^2) = o_P(1/\sqrt{n}), \end{aligned}$$

with the last rate implied by Assumption B.3. As a result, we have

$$\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} m(X_i, Y_i, \hat{\alpha}, \hat{\eta}(X_i), \hat{r}(X_i)) = \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \hat{r}(X_i)) + o_P(1/\sqrt{n}). \quad (12)$$

On the other hand, we write

$$\bar{H}(x, y) = \frac{\partial f^*}{\partial r}(x, y, \bar{\alpha}, \bar{\eta}(x), \hat{r}^{(j)}(x)).$$

Then the expansion with respect to  $r$  in Assumption B.3 implies

$$\begin{aligned} & \left| m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \hat{r}(X_i)) - m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), r_{1,0}(X_i)) - \bar{\alpha} \bar{H}(X_i, Y_i) (\hat{r}(X_i) - r_{1,0}(X_i)) \right| \\ &= \left| \bar{\alpha} \left[ f^* \left( \frac{\hat{r}(X_i)}{-\bar{\alpha}} (Y_i + \bar{\eta}(X_i)) \right) - f^* \left( \frac{r}{-\bar{\alpha}} (Y_i + \bar{\eta}(X_i)) \right) \right] \right. \\ & \quad \left. - \bar{\alpha} \frac{\partial f^*}{\partial r}(x, y, \bar{\alpha}, \bar{\eta}(x), \hat{r}^{(j)}(x)) (\hat{r}(X_i) - r_{1,0}(X_i)) \right| \\ & \leq \bar{\alpha} \cdot \bar{g}(X_i, Y_i) (\hat{r}(X_i) - r_{1,0}(X_i))^2, \end{aligned} \quad (13)$$

where by the finite conditional expectation of  $\bar{g}(X, Y)$  and the convergence rate of  $\hat{r}$ , we have

$$\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} \bar{g}(X_i, Y_i) (\hat{r}(X_i) - r_{1,0}(X_i))^2 = O_P(\mathbb{E}[(\hat{r}(X_i) - r_{1,0}(X_i))^2]) = o(1/\sqrt{n}).$$

Therefore, combining the above residual bound with (12), the difference between  $\hat{\mu}_{1,0}^{m,(j)}$  and  $\hat{\mu}_{1,0}^*$  can be written as

$$\begin{aligned} \hat{\mu}_{1,0}^{m,(j)} - \hat{\mu}_{1,0}^* &= \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), r_{1,0}(X_i)) - \frac{1}{n_1} \sum_{i \in \mathcal{I}_1^{(j)}} m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i)) \\ & \quad + \frac{\bar{\alpha}}{n_1} \sum_{i \in \mathcal{I}_1} \bar{H}(X_i, Y_i) (\hat{r}(X_i) - r_{1,0}(X_i)) - \frac{\hat{\alpha}}{n_1} \sum_{i \in \mathcal{I}_1} \hat{r}(X_i) \hat{h}(X_i) + \frac{\hat{\alpha}}{n_0} \sum_{i \in \mathcal{I}_0} \hat{h}(X_i) \\ & \quad + \frac{\alpha^*}{n_1} \sum_{i \in \mathcal{I}_1} r_{0,1}(X_i) h(X_i) - \frac{\alpha^*}{n_0} \sum_{i \in \mathcal{I}_0} h(X_i) + o_P(1/\sqrt{n}). \end{aligned} \quad (14)$$

We now proceed to treat the summations in (14). Firstly, we note that

$$\begin{aligned} & \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} r_{0,1}(X_i) \hat{h}(X_i) + \frac{1}{n_0} \sum_{i \in \mathcal{I}_0} \hat{h}(X_i) - \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} r_{0,1}(X_i) h(X_i) + \frac{1}{n_0} \sum_{i \in \mathcal{I}_0} h(X_i) \\ &= \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} \left( r_{0,1}(X_i) (\hat{h}(X_i) - h(X_i)) - \mathbb{E}[\hat{h}(X) - h(X) | T = 1] \right) \\ & \quad + \frac{1}{n_0} \sum_{i \in \mathcal{I}_0} \left( \hat{h}(X_i) - h(X_i) - \mathbb{E}[\hat{h}(X) - h(X) | T = 0] \right) \\ &= O_P(\|r_{1,0}(\hat{h} - h)\|_{L_2(\mathbb{P}_{X|T=1})}/\sqrt{n} + \|\hat{h} - h\|_{L_2(\mathbb{P}_{X|T=0})}/\sqrt{n}) = o_P(1/\sqrt{n}), \end{aligned}$$

where the first line uses the fact that  $\mathbb{P}_{X|T=1}$  to  $\mathbb{P}_{X|T=0}$  admits a covariate shift  $r_{1,0}(\cdot)$ , hence the two expectations are equivalent, and the two summations consist of i.i.d. mean zero samples. The last rate uses the fact that

$$\|r_{1,0}(\hat{h} - h)\|_{L_2(\mathbb{P}_{X|T=1})} \leq \|r_{1,0}(\hat{h} - \bar{h})\|_{L_2(\mathbb{P}_{X|T=1})} + \|r_{1,0}(\bar{h} - h)\|_{L_2(\mathbb{P}_{X|T=1})} = o_P(1)$$

and similar for the other term. Again by the covariate shift  $r_{1,0}(x)$  between two groups, we have

$$(\hat{\alpha} - \alpha^*) \cdot \left( \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} r_{0,1}(X_i) h(X_i) - \frac{1}{n_0} \sum_{i \in \mathcal{I}_0} h(X_i) \right) = O_P(|\hat{\alpha} - \alpha^*|/\sqrt{n}) = o_P(1/\sqrt{n})$$

since  $|\hat{\alpha} - \alpha^*| \leq |\hat{\alpha} - \bar{\alpha}| + |\bar{\alpha} - \alpha^*| = o_P(1)$ . Therefore, we have

$$\begin{aligned} & \frac{\hat{\alpha}}{n_1} \sum_{i \in \mathcal{I}_1} \hat{r}(X_i) \hat{h}(X_i) - \frac{\hat{\alpha}}{n_0} \sum_{i \in \mathcal{I}_0} \hat{h}(X_i) \\ &= \frac{\hat{\alpha}}{n_1} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r_{0,1}(X_i)) \hat{h}(X_i) + \frac{\alpha^*}{n_1} \sum_{i \in \mathcal{I}_1} r_{0,1}(X_i) h(X_i) - \frac{\alpha^*}{n_0} \sum_{i \in \mathcal{I}_0} h(X_i) + o_P(1/\sqrt{n}). \end{aligned}$$

Combined with (14), the above error bounds imply

$$\begin{aligned} \hat{\mu}_{1,0}^{m,(j)} - \hat{\mu}_{1,0}^* &= \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), r_{1,0}(X_i)) - \frac{1}{n_1} \sum_{i \in \mathcal{I}_1^{(j)}} m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i)) \\ &+ \frac{\bar{\alpha}}{n_1} \sum_{i \in \mathcal{I}_1} \bar{H}(X_i, Y_i) (\hat{r}(X_i) - r_{1,0}(X_i)) \\ &- \frac{\hat{\alpha}}{n_1} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r_{0,1}(X_i)) \hat{h}(X_i) + o_P(1/\sqrt{n}). \end{aligned} \quad (15)$$

We further note that since  $\mathbb{E}[\bar{H}(x, Y(1)) | X = x, T = 1] \leq M$  for all  $x \in \mathcal{X}$ , we have

$$\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} \bar{H}(X_i, Y_i) (\hat{r}(X_i) - r_{1,0}(X_i)) = O_P(\|\hat{r} - r_{1,0}\|_{L_2(\mathbb{P}_{X|T=1})}),$$

hence by the convergence rate in Assumption B.3, we have

$$\frac{\hat{\alpha} - \bar{\alpha}}{n_1} \sum_{i \in \mathcal{I}_1} \bar{H}(X_i, Y_i) (\hat{r}(X_i) - r_{1,0}(X_i)) = O_P(|\hat{\alpha} - \bar{\alpha}| \cdot \|\hat{r} - r_{1,0}\|_{L_2(\mathbb{P}_{X|T=1})}) = o_P(1/\sqrt{n}).$$

As a result, we obtain

$$\begin{aligned} & \frac{\bar{\alpha}}{n_1} \sum_{i \in \mathcal{I}_1} \bar{H}(X_i, Y_i) (\hat{r}(X_i) - r_{1,0}(X_i)) - \frac{\hat{\alpha}}{n_1} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r_{0,1}(X_i)) \hat{h}(X_i) \\ &= \frac{\hat{\alpha}}{n_1} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r_{1,0}(X_i)) (\bar{H}(X_i, Y_i) - \hat{h}(X_i)) + o_P(1/\sqrt{n}). \end{aligned} \quad (16)$$

The above term excluding the  $\hat{\alpha}$  factor can be decomposed as

$$\begin{aligned} & \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r_{1,0}(X_i)) (\bar{H}(X_i, Y_i) - \hat{h}(X_i)) \\ &= \underbrace{\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r_{1,0}(X_i)) (\bar{H}(X_i, Y_i) - \bar{h}(X_i))}_{(i)} + \underbrace{\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r_{1,0}(X_i)) (\bar{h}(X_i, Y_i) - \hat{h}(X_i))}_{(ii)}. \end{aligned}$$

The first term is unbiased since  $\mathbb{E}[\bar{H}(x, Y(1)) | X = x, T = 1] = \bar{h}(x)$  almost surely. We thus have

$$\begin{aligned} |(i)| &= O_P\left(\left\{\mathbb{E}[(\hat{r}(X_i) - r_{1,0}(X_i))^2 (\bar{H}(X_i, Y_i) - \bar{h}(X_i))^2 | T_i = 1]\right\}^{1/2}/\sqrt{n}\right) \\ &= O_P(\|\hat{r} - r_{1,0}\|_{L_2(\mathbb{P}_{\cdot|T=1})}/\sqrt{n}) = o_P(1/\sqrt{n}) \end{aligned}$$

since the conditional variance of  $\bar{H}(X, Y)$  is uniformly bounded. By Cauchy-Schwarz inequality, the other term is bounded as

$$\begin{aligned} |(ii)| &\leq \sqrt{\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} (\hat{r}(X_i) - r_{1,0}(X_i))^2} \sqrt{\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} (\bar{h}(X_i, Y_i) - \hat{h}(X_i))^2} \\ &= O_P(\|\hat{r} - r_{1,0}\|_{L_2(\mathbb{P}_{\cdot|T=1})} \cdot \|\hat{h} - \bar{h}\|_{L_2(\mathbb{P}_{\cdot|T=1})}) = o_P(1/\sqrt{n}) \end{aligned}$$

by the convergence rates in Assumption B.3. Hence (16) is also  $o_P(1/\sqrt{n})$ . This further leads (15) to

$$\begin{aligned}\widehat{\mu}_{1,0}^{m,(j)} - \widehat{\mu}_{1,0}^* &= \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), r_{1,0}(X_i)) \\ &\quad - \frac{1}{n_1} \sum_{i \in \mathcal{I}_1^{(j)}} m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i)) + o_P(1/\sqrt{n}).\end{aligned}\quad (17)$$

Again with similar arguments as in the proof of Proposition B.4, when  $\eta^*$  and  $\bar{\eta}$  are sufficiently smooth, sufficiently small  $\|\theta^* - \bar{\theta}\|_{L_2(\mathbb{P}_{X|T=1})}$  implies sufficiently small  $\|\theta^* - \bar{\theta}\|_\infty$ . Therefore, the second-order expansion at  $(\alpha^*, \eta^*(x))$  of Assumption B.3 implies

$$\begin{aligned}&\left| \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), r_{1,0}(X_i)) - \frac{1}{n_1} \sum_{i \in \mathcal{I}_1^{(j)}} m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i)) \right. \\ &\quad \left. - \frac{1}{n_1} \sum_{i \in \mathcal{I}_1^{(j)}} \nabla_{a,b} m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i)) [\bar{\alpha} - \alpha^*, \bar{\eta}(X_i) - \eta^*(X_i)] \right| \\ &\leq \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} \bar{g}(X_i, Y_i) \|\bar{\theta}(X_i) - \theta^*(X_i)\|_2^2.\end{aligned}\quad (18)$$

By the optimality of  $\alpha^*$  and the exchangeability of differentiation and expectation,  $\nabla_a m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i))$  has mean zero, thus we have

$$\frac{1}{n_1} \sum_{i \in \mathcal{I}_1^{(j)}} \nabla_a m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i)) [\bar{\alpha} - \alpha^*] = O_P(|\bar{\alpha} - \alpha^*|/\sqrt{n}) = o_P(1/\sqrt{n}).$$

Meanwhile, since  $\eta^*(x)$  is the per- $x$  optimizer given  $\alpha^*$ ,  $\nabla_b m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i))$  has conditional mean zero, hence

$$\begin{aligned}&\frac{1}{n_1} \sum_{i \in \mathcal{I}_1^{(j)}} \nabla_b m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i)) [\bar{\eta}(X_i) - \eta^*(X_i)] \\ &= O_P\left(\left\| \nabla_b m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i)) [\bar{\eta}(X_i) - \eta^*(X_i)] \right\|_{L_2(\mathbb{P}_{\cdot|T=1})} / \sqrt{n}\right) \\ &= O_P(\|\bar{\eta} - \eta^*\|_{L_2(\mathbb{P}_{\cdot|T=1})} / \sqrt{n}) = o_P(1/\sqrt{n}),\end{aligned}$$

where we use the fact that the second moment of  $\nabla_b m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i))$  conditional on  $X$  is uniformly upper bounded as stated in Assumption B.3. Also, the uniformly bounded conditional expectation of  $\bar{g}(X, Y)$  implies

$$\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} \bar{g}(X_i, Y_i) \|\bar{\theta}(X_i) - \theta^*(X_i)\|_2^2 = O_P(\|\bar{\theta} - \theta^*\|_{L_2(\mathbb{P}_{X|T=1})}^2) = o_P(1/\sqrt{n}).$$

Returning to (17), this leads to  $\sqrt{n}(\widehat{\mu}_{1,0}^{m,(j)} - \widehat{\mu}_{1,0}^*) = o_P(1)$  for each fold.

We now return to the original notations for all folds of data, where  $n_1 = |\mathcal{I}_1|$  and  $n_0 = |\mathcal{I}_0|$  represent all treated or control samples. Averaging over the four folds, we have  $\sqrt{n}(\widehat{\mu}_{1,0}^{m,-} - \widehat{\mu}_{1,0}^{m,-,*}) = o_P(1)$ , where

$$\widehat{\mu}_{1,0}^{m,-,*} = \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} \left\{ m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i)) - \alpha^* r_{0,1}(X_i) h(X_i) \right\} + \frac{\alpha^*}{|\mathcal{I}_0^{(j)}|} \sum_{i \in \mathcal{I}_0} h(X_i).$$

Since  $\mathbb{E}[\widehat{\mu}_{1,0}^{m,-,*}] = \mu_{1,0}^{m,-}$ , invoking Central Limit Theorem and Slutsky's lemma, we have  $\sqrt{n}(\widehat{\mu}_{1,0}^{m,-} - \mu_{1,0}^{m,-}) \rightsquigarrow N(0, \text{Var}(\phi_{1,-}^m(X, Y, T)))$ , where the influence function is given by

$$\phi_{1,-}^m(X, Y, T) = \frac{T}{p_1} \left\{ m(X, Y, \alpha^*, \eta^*(X), r_{0,1}(X)) - \alpha^* r_{0,1}(X) h(X) \right\} + \frac{1-T}{p_0} \alpha^* h(X),$$

with  $p_1 = \mathbb{P}(T = 1) = 1 - p_0$ . It is also straightforward to see that

$$\begin{aligned} \text{Var}(\phi_{1,-}^m(X, Y, T)) &= \frac{1}{p_1} \text{Var}\left(m(X, Y, \alpha^*, \eta^*(X), r_{0,1}(X)) - \alpha^* r_{0,1}(X) h(X) \mid T = 1\right) \\ &\quad + \frac{1}{p_0} \text{Var}(\alpha^* h(X) \mid T = 0). \end{aligned}$$

Therefore, a consistent variance estimator can be constructed by plugging in their empirical counterparts. To be specific, for  $i \in \mathcal{I}_1$  we let

$$\widehat{d}_{1,i} = m(X_i, Y_i, \widehat{\alpha}^{j[i]}, \widehat{\eta}^{j[i]}(X_i), \widehat{r}^{j[i]}(X_i)) - \widehat{\alpha}^{j[i]} \widehat{r}^{j[i]}(X_i) \widehat{h}^{j[i]}(X_i),$$

where  $j[i]$  is the fold that sample  $i$  belongs to. It aims at approximating

$$d_{1,i} := m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i)) - \alpha^* r_{0,1}(X_i) h(X_i).$$

We show  $\frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} (\widehat{d}_{1,i}^2 - d_{1,i}^2) = o_P(1)$  as an example. By Cauchy-Schwarz inequality,

$$\begin{aligned} \left| \frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} (\widehat{d}_{1,i}^2 - d_{1,i}^2) \right| &= \left| \frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} (\widehat{d}_{1,i} - d_{1,i})^2 + \frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} 2(\widehat{d}_{1,i} - d_{1,i}) d_{1,i} \right| \\ &\leq \frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} (\widehat{d}_{1,i} - d_{1,i})^2 + 2 \sqrt{\frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} (\widehat{d}_{1,i} - d_{1,i})^2} \cdot \sqrt{\frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} d_{1,i}^2}. \end{aligned}$$

It thus boils down to show that  $\frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} (\widehat{d}_{1,i} - d_{1,i})^2 = o_P(1)$ . We suppress the dependence on  $j$  to simplify notations. By the second-order expansion at  $\bar{\theta}(x) = (\bar{\alpha}, \bar{\eta}(x))$ , we have

$$\begin{aligned} &\left| m(X_i, Y_i, \widehat{\alpha}, \widehat{\eta}(X_i), \widehat{r}(X_i)) - m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \widehat{r}(X_i)) \right| \\ &\leq \left| \nabla_{\alpha, b} m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \widehat{r}(X_i)) [\widehat{\alpha} - \bar{\alpha}, \widehat{\eta}(X_i) - \bar{\eta}(X_i)] \right| + g(X_i, Y_i) \|\widehat{\theta}(X_i) - \bar{\theta}(X_i)\|_2^2. \end{aligned}$$

when  $\|\widehat{\theta} - \bar{\theta}\|_{L_2(\mathbb{P}_{X|T=1})}$  is sufficiently small. By the finite second moment condition and the consistency of  $\widehat{\alpha}$  and  $\widehat{\eta}$ , we have

$$\frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} \left| m(X_i, Y_i, \widehat{\alpha}, \widehat{\eta}(X_i), \widehat{r}(X_i)) - m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \widehat{r}(X_i)) \right|^2 = o_P(1).$$

Furthermore, the expansion with respect to  $r$  in (13) implies

$$\begin{aligned} &\left| m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \widehat{r}(X_i)) - m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), r_{1,0}(X_i)) \right| \\ &\leq \bar{\alpha} \bar{H}(X_i, Y_i) (\widehat{r}(X_i) - r_{1,0}(X_i)) + \bar{\alpha} \cdot \bar{g}(X_i, Y_i) (\widehat{r}(X_i) - r_{1,0}(X_i))^2, \end{aligned}$$

the finite moment conditions of which imply

$$\frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} \left| m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), \widehat{r}(X_i)) - m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), r_{1,0}(X_i)) \right|^2 = o_P(1).$$

Finally, the expansion around  $(\alpha^*, \eta^*(x))$  similar to (18) leads to

$$\begin{aligned} &\left| m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), r_{1,0}(X_i)) - m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i)) \right| \\ &\leq \left| \nabla_{\alpha, b} m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i)) [\bar{\alpha} - \alpha^*, \bar{\eta}(X_i) - \eta^*(X_i)] \right| + \bar{g}(X_i, Y_i) \|\bar{\theta}(X_i) - \theta^*(X_i)\|_2^2. \end{aligned}$$

By the consistency of  $\bar{\theta}$  and finite moments, we thus have

$$\frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} \left| m(X_i, Y_i, \bar{\alpha}, \bar{\eta}(X_i), r_{1,0}(X_i)) - m(X_i, Y_i, \alpha^*, \eta^*(X_i), r_{0,1}(X_i)) \right|^2 = o_P(1).$$

Futhermore, we note that

$$\left| \hat{\alpha} \hat{r}(X_i) \hat{h}(X_i) - \alpha^* r_{1,0}(X_i) h(X_i) \right|^2 \leq 2\hat{\alpha}^2 \left| \hat{r}(X_i) \hat{h}(X_i) - r_{1,0}(X_i) h(X_i) \right|^2 + 2|\hat{\alpha} - \alpha^*|^2 \cdot r_{1,0}(X_i)^2 h(X_i)^2,$$

which is  $o_P(1)$  by the consistency of estimators. Invoking Cauchy-Schwarz inequality, we thus obtain  $\frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} (\hat{d}_{1,i} - d_{1,i})^2 = o_P(1)$ . This also implies  $\frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} (\hat{d}_{1,i} - d_{1,i}) = o_P(1)$ , hence  $\frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} \hat{d}_{1,i}^2 - \left( \frac{1}{|\mathcal{I}_1^{(j)}|} \sum_{i \in \mathcal{I}_1^{(j)}} \hat{d}_{1,i} \right)^2$  is consistent for  $\text{Var}(m(X, Y, \alpha^*, \eta^*(X), r_{0,1}(X)) - \alpha^* r_{0,1}(X) h(X) \mid T = 1)$ . Similar consistency can be obtained for the other variance term, which completes the proof of the consistency of the variance estimator. Therefore, we conclude the proof of Theorem B.1.  $\square$

## E Technical lemmas

**Lemma E.1.** *Let  $\mathcal{F}_n$  be a sequence of  $\sigma$ -algebra, and let  $A_n \geq 0$  be a sequence of nonnegative random variables. If  $\mathbb{E}[A_n \mid \mathcal{F}_n] = o_P(1)$ , then  $A_n = o_P(1)$ .*

*Proof.* Proof of Lemma E.1 By Markov's inequality, for any  $\epsilon > 0$ , we have

$$B_n := \mathbb{P}(A_n > \epsilon \mid \mathcal{F}_n) \leq \frac{\mathbb{E}[A_n \mid \mathcal{F}_n]}{\epsilon} = o_P(1),$$

and  $B_n \in [0, 1]$  are bounded random variables. For any subsequence  $\{n_k\}_{k \geq 1}$  of  $\mathbb{N}$ , since  $B_{n_k} \xrightarrow{P} 0$ , there exists a subsequence  $\{n_{k_i}\}_{i \geq 1} \subset \{n_k\}_{k \geq 1}$  such that  $B_{n_{k_i}} \xrightarrow{\text{a.s.}} 0$  as  $i \rightarrow \infty$ . By the dominated convergence theorem, we have  $\mathbb{E}[B_{n_{k_i}}] \rightarrow 0$ , or equivalently,  $\mathbb{P}(A_{n_{k_i}} > \epsilon) \rightarrow 0$ . Therefore, for any subsequence  $\{n_k\}_{k \geq 1}$  of  $\mathbb{N}$ , there exists a subsequence  $\{n_{k_i}\}_{i \geq 1} \subset \{n_k\}_{k \geq 1}$  such that  $A_{n_{k_i}} \xrightarrow{P} 0$  as  $i \rightarrow \infty$ . By the arbitrariness of  $\{n_k\}_{k \geq 1}$ , we know  $A_n \xrightarrow{P} 0$  as  $n \rightarrow \infty$ , which completes the proof.  $\square$

## References

- Chen, X. (2007). Large sample sieve estimation of semi-nonparametric models. *Handbook of econometrics*, 6:5549–5632.
- Chen, X. and Shen, X. (1998). Sieve extremum estimates for weakly dependent data. *Econometrica*, pages 289–314.
- Geer, S. A., van de Geer, S., and Williams, D. (2000). *Empirical Processes in M-estimation*, volume 6. Cambridge university press.
- Jin, Y., Ren, Z., and Candès, E. J. (2023). Sensitivity analysis of individual treatment effects: A robust conformal inference approach. *Proceedings of the National Academy of Sciences*, 120(6):e2214889120.
- Kullback, S. and Leibler, R. A. (1951). On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86.
- Luenberger, D. G. (1997). *Optimization by vector space methods*. John Wiley & Sons.
- Rudin, W. et al. (1976). *Principles of mathematical analysis*, volume 3. McGraw-hill New York.
- Timan, A. F. (2014). *Theory of approximation of functions of a real variable*. Elsevier.
- Van der Vaart, A. W. (2000). *Asymptotic statistics*, volume 3. Cambridge university press.
- Yadlowsky, S., Namkoong, H., Basu, S., Duchi, J., and Tian, L. (2018). Bounds on the conditional and average treatment effect with unobserved confounding factors. *arXiv preprint arXiv:1808.09521*.