

1

## 2 **Supplementary Information for**

### 3 **Sensitivity analysis of individual treatment effects: a robust conformal inference** 4 **approach**

5 Ying Jin, Zhimei Ren, Emmanuel J. Candès

6 Emmanuel J. Candès.

7 E-mail: [candes@stanford.edu](mailto:candes@stanford.edu)

#### 8 **This PDF file includes:**

9 Figs. S1 to S12

10 Table S1

11 SI References

Counterfactual	Bound	ATE-type	ATT-type	ATC-type	General
Y(1)	$\ell(x)$	$p_1 \cdot \left(1 + \frac{1}{\Gamma \cdot r(x)}\right)$	1	$\frac{p_1}{p_0} \cdot \left(\frac{1}{\Gamma \cdot r(x)}\right)$	$p_1 \cdot \frac{dQ_X}{dP_X}(x) \cdot \left(1 + \frac{1}{\Gamma \cdot r(x)}\right)$
	$u(x)$	$p_1 \cdot \left(1 + \frac{\Gamma}{r(x)}\right)$	1	$\frac{p_1}{p_0} \cdot \frac{\Gamma}{r(x)}$	$p_1 \cdot \frac{dQ_X}{dP_X}(x) \cdot \left(1 + \frac{\Gamma}{r(x)}\right)$
Y(0)	$\ell(x)$	$p_0 \cdot \left(1 + \frac{r(x)}{\Gamma}\right)$	$\frac{p_0}{p_1} \cdot \frac{r(x)}{\Gamma}$	1	$p_0 \cdot \frac{dQ_X}{dP_X}(x) \cdot \left(1 + \frac{r(x)}{\Gamma}\right)$
	$u(x)$	$p_0 \cdot \left(1 + \Gamma \cdot r(x)\right)$	$\frac{p_0}{p_1} \cdot \Gamma \cdot r(x)$	1	$p_0 \cdot \frac{dQ_X}{dP_X}(x) \cdot \left(1 + \Gamma \cdot r(x)\right)$

**Table S1. Summary of the upper and lower bounds of the likelihood ratio for different inferential targets.** For  $t \in \{0, 1\}$ ,  $p_t = \mathbb{P}(T = t)$  and  $r(x) = e(x)/(1 - e(x))$  is the odds ratio of the propensity score. The training distribution for  $Y(t)$  is always  $\mathbb{P}_{X, Y(t) | T=t}$ . For target distributions, **ATE-type** refers to  $\mathbb{P}_{X, Y(t)}$ ; **ATT-type** refers to  $\mathbb{P}_{X, Y(t) | T=1}$ ; **ATC-type** refers to  $\mathbb{P}_{X, Y(t) | T=0}$ ; **General** refers to  $Q_X \times \mathbb{P}_{Y(t) | X}$ .

## 13 2. Additional results

14 **A. Constructing  $\hat{G}_n(\cdot)$  in the PAC-type procedure.** In this part, we provide the construction of  $\hat{G}_n(t)$  in Section A  
15 with Waudby-Smith–Ramdas bound (1). The proof is a slight modification of (1) and included in SI Appendix D for  
16 completeness.

**Proposition 2.1** (Waudby-Smith–Ramdas lower confidence bound for c.d.f.s). *Suppose  $\sup_x \hat{u}(x) \leq M$  for some constant  $M > 0$ . For  $t \in \mathbb{R}$  being any constant or any random variable in  $\sigma(\mathcal{D}_{\text{train}})$ , and any  $\delta \in (0, 1)$ , we define  $\hat{G}_n(t) = \max\{\hat{G}_n^L(t), \hat{G}_n^U(t)\}$ , where*

$$\hat{G}_n^L(t) = M \cdot \inf\{g \geq 0: \max_{1 \leq i \leq n} \mathcal{K}_i^L(g) \leq 2/\delta\},$$

$$\hat{G}_n^U(t) = 1 - M + M \cdot \inf\{g \geq 0: \max_{1 \leq i \leq n} \mathcal{K}_i^U(g) \leq 2/\delta\}.$$

For any  $g \geq 0$ , the thresholding functions for  $i = 1, \dots, n$  are defined as

$$\mathcal{K}_i^L(g) = \prod_{j=1}^i \left(1 + \nu_j^L \cdot [\mathbb{1}_{\{V_j \leq t\}} \hat{\ell}(X_j)/M - g]\right), \quad \mathcal{K}_i^U(g) = \prod_{j=1}^i \left(1 + \nu_j^U \cdot [1 - \mathbb{1}_{\{V_j > t\}} \hat{u}(X_j)/M - g]\right),$$

where  $\nu_j^L = \min\{1, \sqrt{2 \log(2/\delta)/[n(\hat{\sigma}_{j-1}^L)^2]}\}$ ,  $\nu_j^U = \min\{1, \sqrt{2 \log(2/\delta)/[n(\hat{\sigma}_{j-1}^U)^2]}\}$ , and

$$(\hat{\sigma}_i^L)^2 = \frac{\frac{1}{4} + \sum_{j=1}^i \left(\mathbb{1}_{\{V_j \leq t\}} \frac{\hat{\ell}(X_j)}{M} - \hat{\mu}_j^L\right)^2}{1+i}, \quad \hat{\mu}_i^L = \frac{\frac{1}{2} + \sum_{j=1}^i \mathbb{1}_{\{V_j \leq t\}} \frac{\hat{\ell}(X_j)}{M}}{1+i},$$

$$(\hat{\sigma}_i^U)^2 = \frac{\frac{1}{4} + \sum_{j=1}^i \left(1 - \mathbb{1}_{\{V_j > t\}} \frac{\hat{u}(X_j)}{M} - \hat{\mu}_j^U\right)^2}{1+i}, \quad \hat{\mu}_i^U = \frac{\frac{1}{2} + \sum_{j=1}^i (1 - \mathbb{1}_{\{V_j > t\}} \frac{\hat{u}(X_j)}{M})}{1+i}.$$

17 Then it holds that  $\mathbb{P}_{\mathcal{D}_{\text{calib}}}(\hat{G}_n(t) \leq G(t)) \geq 1 - \delta$  for  $G(t)$  defined in [13].

18 **B. Validity of prediction intervals for ITE.** In this part, we provide the coverage guarantee for prediction of ITEs  
19 omitted in Section A.

**Proposition 2.2.** *Consider a new test sample where we observe  $(X_{n+1}, T_{n+1}, Y_{n+1}(T_{n+1}))$ , with  $T_{n+1} = w$  for  $w \in \{0, 1\}$ . Let  $\mathcal{D}$  be the calibration data generated from  $\mathbb{P}^{\text{sup}}$  under confounding level  $\Gamma$ . If  $\hat{C}_{1-w}(X_{n+1}, \Gamma, 1 - \alpha)$  is constructed by Algorithm 1, then*

$$\mathbb{P}(Y_{n+1}(1) - Y_{n+1}(0) \in \hat{C}(X_{n+1}, \Gamma) \mid T_{n+1} = w) \geq 1 - \alpha - \hat{\Delta},$$

for the prediction set  $\hat{C}(X_{n+1}, \Gamma, 1 - \alpha)$  in [A], where the probability is over  $\mathcal{D}_{\text{calib}}$  as well as the test point, and  $\hat{\Delta}$  is the gap of coverage for  $Y_{n+1}(1 - w)$  in Theorem 1.2. If  $\hat{C}_{1-w}(X_{n+1}, \Gamma, 1 - \alpha)$  is constructed by Algorithm 2 with confidence level  $\delta \in (0, 1)$ , then with probability at least  $1 - \delta$  with respect to  $\mathcal{D}_{\text{calib}}$ , we have

$$\mathbb{P}(Y_{n+1}(1) - Y_{n+1}(0) \in \hat{C}(X_{n+1}, \Gamma) \mid T_{n+1} = w, \mathcal{D}_{\text{calib}}) \geq 1 - \alpha - \hat{\Delta},$$

20 where  $\hat{\Delta}$  is the gap of coverage for  $Y_{n+1}(1 - w)$  in Theorem 3.1.

**Proposition 2.3.** Consider a new test sample  $X_{n+1}$ . Let  $\mathcal{D}$  be the observations for which  $\mathbb{P}^{\text{sup}}$  is under confounding level  $\Gamma$ . If for  $w \in \{0, 1\}$ ,  $\hat{C}_w(X_{n+1}, \Gamma, 1 - \alpha/2, \delta/2)$  is constructed by Algorithm 1, then

$$\mathbb{P}(Y_{n+1}(1) - Y_{n+1}(0) \in \hat{C}(X_{n+1}, \Gamma)) \geq 1 - \alpha - \hat{\Delta}_1 - \hat{\Delta}_0,$$

where the probability is over  $\mathcal{D}$  as well as the test point;  $\hat{\Delta}_0, \hat{\Delta}_1$  is the coverage gap in Theorem 1.2 for counterfactual prediction of  $Y_{n+1}(1), Y_{n+1}(0)$  when the bound functions are estimated. If  $\hat{C}_{1-w}(X_{n+1}, \Gamma, 1 - \alpha/2, \delta/2)$  is constructed by Algorithm 2, then with probability at least  $1 - \delta$  with respect to  $\mathcal{D}_{\text{calib}}$ , we have

$$\mathbb{P}(Y_{n+1}(1) - Y_{n+1}(0) \in \hat{C}(X_{n+1}, \Gamma) \mid \mathcal{D}_{\text{calib}}) \geq 1 - \alpha - \hat{\Delta}_1 - \hat{\Delta}_0,$$

21 where  $\hat{\Delta}_0$  and  $\hat{\Delta}_1$  are the coverage gaps in Theorem 1.2 for counterfactual prediction of  $Y_{n+1}(1), Y_{n+1}(0)$  when the  
22 bound functions are estimated.

23 **C. Types of null hypotheses.** With the general recipe of assessing robustness of causal conclusions on ITE, we provide  
24 concrete examples of the target set  $C$  and the corresponding forms of  $\hat{C}(X_{n+1}, \Gamma)$ .

25 **Sharp null** One might be interested in the sharp null, i.e., whether the individual treatment effect is zero. In this case,  
26 one could let  $C = \{0\}$ , and the prediction set can take any form. Rejecting  $H_0(\Gamma)$  is saying that  $Y_{n+1}(1) \neq Y_{n+1}(0)$   
27 unless the true confounding level obeys  $\Gamma^* > \Gamma$ .

28 **Directional null** If we presume the stochastic nature of ITEs, the sharp null might be implausible. In this case, one  
29 may be interested in the directional null with  $Y_{n+1}(1) - Y_{n+1}(0) \leq 0$ , which is equivalent to choosing  $C = (-\infty, 0]$ .  
30 Rejecting  $H_0(\Gamma)$  here is saying that the individual treatment effect is positive unless  $\Gamma^* > \Gamma$ . More generally, one might  
31 consider  $C = (-\infty, a]$  for some  $a \in \mathbb{R}$  to test whether the ITE is above a certain value. To make sense of the multiple  
32 testing procedure, one-sided prediction intervals are constructed for ITE, i.e.,  $\hat{C}(X_{n+1}, \Gamma) = [\hat{I}(X_{n+1}, \Gamma), \infty)$  for some  
33  $\hat{I}(X_{n+1}, \Gamma) \in \mathbb{R}$ . It can be achieved by one-sided prediction intervals  $\hat{C}_0(X_{n+1}, \Gamma, 1 - \alpha, 1 - \delta) = (-\infty, \hat{Y}_0(X_{n+1}, \Gamma)]$   
34 if  $Y_{n+1}(0)$  is missing and  $\hat{C}_1(X_{n+1}, \Gamma, 1 - \alpha, 1 - \delta) = [\hat{Y}_1(X_{n+1}, \Gamma), \infty)$  if  $Y_{n+1}(1)$  is missing, with a Bonferroni  
35 correction as introduced in Section A if both outcomes are missing. The  $\Gamma$ -value is hence the smallest  $\Gamma$  such that  
36 the two prediction intervals overlap.

37 **D. Sharper results under parametric models.** From the sharpness results in Section 3.C.2 (main text), we see that  
38 in the worst case scenario, the odds ratio [3]  $\frac{\mathbb{P}(T=1 \mid X=x, U=u)\mathbb{P}(T=0 \mid X=x)}{\mathbb{P}(T=0 \mid X=x, U=u)\mathbb{P}(T=1 \mid X=x)}$  is equal to either  $\Gamma$  or  $1/\Gamma$ . Admittedly,  
39 we are unlikely to encounter this in practice. We now discuss how to get sharper prediction intervals by positing a  
40 parametric model.

Let  $s: \mathcal{X} \rightarrow \mathbb{R}^{d_1}$  and  $t: \mathcal{Y} \rightarrow \mathbb{R}^{d_2}$  be known functions. Imagine we are interested in predicting  $Y(0)$  for a treated unit. Now assume that the true (confounded) treatment rule  $\mathbb{P}(T = 1 \mid X, Y(0))$  follows the logistic model

$$\mathbb{P}(T = 1 \mid X, Y(0)) = \text{logit}^{-1}(\alpha^\top s(X) + \gamma^\top t(Y(0))), \quad [1]$$

41 where we define  $\text{logit}^{-1}(x) = e^x / (1 + e^x)$ , and  $\alpha \in \mathbb{R}^{d_1}, \gamma \in \mathbb{R}^{d_2}$  are unknown parameters. This parametric model may  
42 apply to healthcare data in which a doctor assigns a treatment to patients according to a logistic model depending  
43 on a full set of characteristics while only a few of those are included as covariates in the dataset due to privacy  
44 issues. The undocumented characteristics thus become the unmeasured confounder, and it may still be reasonable to  
45 assume that the treatment rule reduced to the outcome level, i.e.,  $\mathbb{P}(T = 1 \mid X, Y(0))$ , approximately follows a logistic  
46 model [1].

Under the parametric model, the  $\Gamma$ -selection condition translates to

$$\frac{1}{\Gamma} \cdot \frac{e(X)}{1 + e(X)} \leq \frac{d\mathbb{P}_{X, Y(0) \mid T=1}}{d\mathbb{P}_{X, Y(0) \mid T=0}} = e^{\alpha^\top s(X) + \gamma^\top t(Y(0))} \cdot \frac{\mathbb{P}(T = 0)}{\mathbb{P}(T = 1)} \leq \Gamma \cdot \frac{e(X)}{1 + e(X)}.$$

Now let  $\hat{e}(\cdot)$  be an estimator for  $e(\cdot)$ , and  $\hat{p} = \frac{1}{n} \sum_{i=1}^n T_i$ . We can modify the optimization problem in Remark 1.1 of the main text by replacing the constraints on  $W_i$ 's with those on  $(\alpha, \gamma)$ . Specifically, for each  $k$ , we now solve

$$\begin{aligned} & \underset{\alpha, \gamma}{\text{minimize}} && \frac{\sum_{i=1}^k \exp\{\alpha^\top s(X_{[i]}) + \gamma^\top t(Y_{[i]}(0))\}}{\sum_{i=1}^n \exp\{\alpha^\top s(X_i) + \gamma^\top t(Y_i(0))\}} + W_{n+1} \\ & \text{subject to} && \frac{\hat{e}(X_{n+1})}{\Gamma(1 + \hat{e}(X_{n+1}))} \leq W_{n+1} \leq \frac{\Gamma \cdot \hat{e}(X_{n+1})}{1 + \hat{e}(X_{n+1})} \\ & && \frac{\hat{e}(X_i)}{\Gamma(1 + \hat{e}(X_i))} \leq e^{\alpha^\top s(X_i) + \gamma^\top t(Y_i(0))} \cdot \frac{1 - \hat{p}}{\hat{p}} \leq \frac{\Gamma \cdot \hat{e}(X_i)}{1 + \hat{e}(X_i)}, \quad \forall i \in \mathcal{D}_{\text{calib}}. \end{aligned}$$

47 Here, as before,  $[1], [2], \dots, [n]$  is the permutation of  $\{1, \dots, n\}$  such that  $V_{[1]} \leq V_{[2]} \leq \dots \leq V_{[n]}$ . If we believe in the  
 48 treatment assignment model, then this gives an extension of Algorithm 1 which provides sharper prediction intervals  
 49 with valid coverage (up to estimation error in  $\hat{e}$  and  $\hat{p}$ ).

50 **E. Nested method for inference on ITEs under unmeasured confounding.** We provide details of the nested method  
 51 for constructing less conservative ITE prediction intervals introduced in Section 4.

52 Assume the calibration data consist of both treated and control units. The nested method first uses the robust  
 53 conformal inference methods in the main text to construct prediction intervals for the ITEs of all calibration data;  
 54 it then learns and calibrates these intervals to produce a prediction interval for the ITE of a test point. Below,  
 55 Algorithm 1 (resp. Algorithm 2) constructs a two-sided (resp. one-sided) prediction interval for the ITE of a test  
 56 point. The output prediction interval achieves  $(1 - \alpha - \gamma)$  coverage (up to estimation error of the propensity score),  
 57 where  $\gamma$  corresponds to the confidence of the second calibration step. Depending on the input preference, such  
 58 coverage is either marginalized over the calibration and test data, or calibration-conditionally valid with probability  
 59 at least  $1 - \delta$  over the calibration data.

60 Algorithms 1 and 2 use a subroutine to construct two-sided and one-sided robust prediction intervals for ITEs with  
 61 either marginal or PAC-type coverage. The choice of the subroutine is flexible and depends on the target. Algorithm 3  
 62 (resp. Algorithm 4) provides a concrete procedure that constructs two-sided (resp. one-sided) prediction intervals for  
 63 ITEs, where we use Conformalized Quantile Regression (CQR) (2) to compute the nonconformity score, and assume  
 64 we have estimates  $\hat{e}(x)$  for the propensity score  $e(x) = \mathbb{P}(T = 1 | X = x)$ ,  $\hat{p}$  for the treatment probability  $\mathbb{P}(T = 1)$ .  
 65 The  $\alpha/2$ -th and  $(1 - \alpha/2)$ -th quantiles for  $\mathbb{P}_{Y_i | X_i, T_i=1}$  and  $\mathbb{P}_{Y_i | X_i, T_i=0}$  on  $\mathcal{D}_{\text{train}}$  are denoted by  $\hat{q}_{\alpha/2}^{(1)}(x)$ ,  $\hat{q}_{1-\alpha/2}^{(1)}(x)$ ,  
 66  $\hat{q}_{\alpha/2}^{(0)}(x)$ , and  $\hat{q}_{1-\alpha/2}^{(0)}(x)$ , respectively.

---

**Algorithm 1** Nested (two-sided) robust conformal prediction of ITEs

---

1: Input: Calibration data  $\{(X_i, Y_i, T_i)\}_{i \in \mathcal{D}_{\text{calib}}}$ , test covariate  $x$ , target level  $\alpha, \gamma \in (0, 1)$ , confidence level  $\delta \in (0, 1)$ .

**Step I: sample splitting**

- 2: Randomly split  $\mathcal{D}_{\text{calib}}$  into  $\mathcal{D}_1$  and  $\mathcal{D}_2$ .
- 3: Randomly split  $\mathcal{D}_2$  into  $\mathcal{D}_{2,1}$  and  $\mathcal{D}_{2,2}$ .

**Step II: robust prediction of ITEs for  $\mathcal{D}_2$**

- 4: **if** target = marginal coverage **then**
- 5:   For all  $i \in \mathcal{D}_2$ , use  $\mathcal{D}_1$  as calibration data to obtain two-sided prediction interval  $[\hat{C}_i^L, \hat{C}_i^R] := [\hat{C}^L(X_i), \hat{C}^R(X_i)]$   
 for ITE with marginal coverage  $1 - \alpha$  (e.g., from Algorithm 3).
- 6: **else if** target =  $1 - \delta$  PAC-type guarantee **then**
- 7:   For all  $i \in \mathcal{D}_2$ , use  $\mathcal{D}_1$  as calibration data to obtain two-sided prediction interval  $[\hat{C}_i^L, \hat{C}_i^R] := [\hat{C}^L(X_i), \hat{C}^R(X_i)]$   
 for ITE with PAC-type coverage  $1 - \alpha$  at confidence level  $\delta/3$  (e.g., from Algorithm 3).

**Step III: nested prediction interval for ITE of the test point**

- 8: Fit  $\hat{m}^L(x)$  and  $\hat{m}^R(x)$ , the conditional mean/median of  $C^L, C^R$  given  $X_i = x$  using  $\mathcal{D}_{2,1}$ .
  - 9: For  $i \in \mathcal{D}_{2,2}$ , compute  $V_i = \max\{\hat{m}^L(X_i) - C_i^L, C_i^R - \hat{m}^R(X_i)\}$ .
  - 10: **if** target = marginal coverage **then**
  - 11:   Compute  $\eta$  as the  $(1 - \gamma)(1 + 1/|\mathcal{D}_{2,2}|)$ -th quantile of the empirical distribution of  $\{V_i : i \in \mathcal{D}_{2,2}\}$ .
  - 12: **else if** target =  $1 - \delta$  PAC-type guarantee **then**
  - 13:   Compute  $\eta$  as a  $(1 - \delta/3)$  lower confidence bound for the  $(1 - \gamma)$ -th population quantile viewing  $\{V_i : i \in \mathcal{D}_{2,2}\}$   
 as i.i.d. copies (using, e.g., Hoeffding's inequality or W-S-R inequality).
  - 14: Output: Prediction set  $\hat{C}(x) = [\hat{m}^L(x) - \eta, \hat{m}^R(x) + \eta]$ .
- 

67 **3. Technical proofs**

68 **A. Proof of Lemma 2.1.**

*Proof of Lemma 2.1.* For any measurable subset  $A \subset \mathcal{U}$ , any  $u \in A$  and any  $x \in \mathcal{X}$ , by the marginal  $\Gamma$ -selection condition [3],

$$\frac{1}{\Gamma} \cdot \mathbb{P}(T = 0 | X = x, U = u) \leq \mathbb{P}(T = 1 | X = x, U = u) \cdot \frac{\mathbb{P}(T = 0 | X = x)}{\mathbb{P}(T = 1 | X = x)} \leq \Gamma \cdot \mathbb{P}(T = 0 | X = x, U = u). \quad [2]$$

---

**Algorithm 2** Nested (one-sided) robust conformal prediction of ITEs
 

---

1: Input: Calibration data  $\{(X_i, Y_i, T_i)\}_{i \in \mathcal{D}_{\text{calib}}}$ , test covariate  $x$ , target level  $\alpha, \gamma \in (0, 1)$ , confidence level  $\delta \in (0, 1)$ .

**Step I: sample splitting**

2: Randomly split  $\mathcal{D}_{\text{calib}}$  into  $\mathcal{D}_1$  and  $\mathcal{D}_2$ .

3: Randomly split  $\mathcal{D}_2$  into  $\mathcal{D}_{2,1}$  and  $\mathcal{D}_{2,2}$ .

**Step II: robust prediction of ITEs for  $\mathcal{D}_2$**

4: **if** desire marginal coverage **then**

5: For all  $i \in \mathcal{D}_2$ , use  $\mathcal{D}_1$  as calibration data, obtain one-sided prediction interval  $[\hat{C}_i^L, +\infty) := [\hat{C}^L(X_i), +\infty)$  for ITEs with marginal coverage  $1 - \alpha$  (e.g., from Algorithm 4).

6: **else if** desire  $1 - \delta$  PAC-type guarantee **then**

7: For all  $i \in \mathcal{D}_2$ , use  $\mathcal{D}_1$  as calibration data, obtain one-sided prediction interval  $[\hat{C}_i^L, +\infty) := [\hat{C}^L(X_i), +\infty)$  for ITE with PAC-type coverage  $1 - \alpha$  at confidence level  $\delta/3$  (e.g., from Algorithm 4).

**Step III: nested prediction interval for ITE of test point**

8: Fit  $\hat{m}^L(x)$ , the conditional mean/median of  $C^L$  given  $X_i = x$  using  $\mathcal{D}_{2,1}$ .

9: For  $i \in \mathcal{D}_{2,2}$ , compute  $V_i = C_i^L - \hat{m}^L(X_i)$ .

10: **if** desire marginal coverage **then**

11: Compute  $\eta$  as the  $1 - (1 - \gamma)(1 + 1/|\mathcal{D}_{2,2}|)$ -th quantile of the empirical distribution of  $\{V_i : i \in \mathcal{D}_{2,2}\}$ .

12: **else if** desire  $1 - \delta$  PAC-type guarantee **then**

13: Compute  $\eta$  as a  $(1 - \delta/3)$  lower confidence bound for the  $\gamma$ -th population quantile viewing  $\{V_i : i \in \mathcal{D}_{2,2}\}$  as i.i.d. copies (using, e.g., Hoeffding's inequality or W-S-R inequality).

14: Output: Prediction set  $\hat{C}(x) = [\hat{m}^L(x) + \eta, \infty)$ .

---

Marginalizing over  $u \in A$  yields

$$\frac{1}{\Gamma} \leq \frac{\mathbb{P}(U \in A \mid X = x, T = 1)}{\mathbb{P}(U \in A \mid X = x, T = 0)} \leq \Gamma \quad [3]$$

for  $\mathbb{P}$ -almost  $x \in \mathcal{X}$ . Since [3] holds for any measurable set in  $\mathcal{U}$ , we have

$$\frac{1}{\Gamma} \leq \frac{d\mathbb{P}_{U \mid X, T=1}}{d\mathbb{P}_{U \mid X, T=0}}(u, x) \leq \Gamma,$$

for  $\mathbb{P}$ -almost all  $u \in \mathcal{U}$  and  $x \in \mathcal{X}$ . Meanwhile, for any measurable set  $B \subset \mathcal{Y}$ , by the tower property, we have for any  $t \in \{0, 1\}$  that

$$\begin{aligned} \mathbb{P}(Y(1) \in B, T = t \mid X) &= \mathbb{E} \left[ \mathbb{E}[\mathbf{1}_{\{Y(1) \in B\}} \mathbf{1}_{\{T=t\}} \mid X, U] \mid X \right] \\ &= \mathbb{E} \left[ \mathbb{E}[\mathbf{1}_{\{Y(1) \in B\}} \mid X, U] \cdot \mathbb{E}[\mathbf{1}_{\{T=t\}} \mid X, U] \mid X \right]. \end{aligned} \quad [4]$$

Rewriting [2], we have  $\mathbb{P}$ -almost surely that

$$\frac{1}{\Gamma} \cdot \mathbb{E}[\mathbf{1}_{\{T=0\}} \mid X, U] \cdot \frac{\mathbb{E}[\mathbf{1}_{\{T=1\}} \mid X]}{\mathbb{E}[\mathbf{1}_{\{T=0\}} \mid X]} \leq \mathbb{E}[\mathbf{1}_{\{T=1\}} \mid X, U] \leq \Gamma \cdot \mathbb{E}[\mathbf{1}_{\{T=0\}} \mid X, U] \cdot \frac{\mathbb{E}[\mathbf{1}_{\{T=1\}} \mid X]}{\mathbb{E}[\mathbf{1}_{\{T=0\}} \mid X]}.$$

Multiplying all sides by  $\mathbb{E}[\mathbf{1}_{\{Y(1) \in B\}} \mid X, U]$  and using [4], we know

$$\frac{1}{\Gamma} \cdot \mathbb{P}(Y(1) \in B, T = 0 \mid X) \cdot \frac{\mathbb{E}[\mathbf{1}_{\{T=1\}} \mid X]}{\mathbb{E}[\mathbf{1}_{\{T=0\}} \mid X]} \leq \mathbb{P}(Y(1) \in B, T = 1 \mid X) \leq \Gamma \cdot \mathbb{P}(Y(1) \in B, T = 0 \mid X) \cdot \frac{\mathbb{E}[\mathbf{1}_{\{T=1\}} \mid X]}{\mathbb{E}[\mathbf{1}_{\{T=0\}} \mid X]}$$

holds  $\mathbb{P}$ -almost surely, and for  $\mathbb{P}$ -almost all  $x \in \mathcal{X}$ ,

$$\frac{1}{\Gamma} \cdot \frac{1 - e(x)}{e(x)} \leq \frac{\mathbb{P}(Y(1) \in B, T = 0 \mid X = x)}{\mathbb{P}(Y(1) \in B, T = 1 \mid X = x)} \leq \Gamma \cdot \frac{1 - e(x)}{e(x)}.$$

Note that

$$\frac{\mathbb{P}(Y(1) \in B \mid X = x, T = 1)}{\mathbb{P}(Y(1) \in B \mid X = x, T = 0)} = \frac{\mathbb{P}(Y(1) \in B, T = 1 \mid X = x)}{\mathbb{P}(Y(1) \in B, T = 0 \mid X = x)} \cdot \frac{1 - e(x)}{e(x)}.$$

---

**Algorithm 3** Subroutine: Two-sided robust prediction for ITE
 

---

- 1: Input: Calibration data  $\{(X_i, Y_i, T_i)\}_{i \in \mathcal{D}}$ , estimated propensity score  $\hat{e}(x)$ , estimated  $\hat{p}$  for  $\mathbb{P}(T = 1)$ , estimated quantile functions  $\hat{q}_{\alpha/2}^{(1)}(x)$ ,  $\hat{q}_{1-\alpha/2}^{(1)}(x)$ , test covariate  $x$ , test treatment  $t \in \{0, 1\}$ , test outcome  $y(t)$ , target level  $\alpha \in (0, 1)$ , confidence level  $\delta \in (0, 1)$  if necessary.
  - 2: Define the nonconformity score  $V_t(x, y) = \max\{\hat{q}_{\alpha/2}^{(t)}(x) - y, y - \hat{q}_{1-\alpha/2}^{(t)}(x)\}$  for  $t = 0, 1$ .
  - 3: Set  $\hat{\ell}(\cdot) = \frac{1}{\Gamma} \left[ \frac{1-\hat{p}}{\hat{p}} \frac{\hat{e}(x)}{1-\hat{e}(x)} \right]^{2t-1}$ ,  $\hat{u}(\cdot) = \Gamma \left[ \frac{1-\hat{p}}{\hat{p}} \frac{\hat{e}(x)}{1-\hat{e}(x)} \right]^{2t-1}$ .
  - 4: **if** desire marginal coverage **then**
  - 5: Input  $\{j \in \mathcal{D} : T_j = 1 - t\}$ ,  $\hat{\ell}(\cdot)$ ,  $\hat{u}(\cdot)$ , score  $V_{1-t}$ , and target  $\alpha$  to Algorithm 1 (main text).
  - 6: Obtain  $[\hat{C}_{1-t}^L(x), \hat{C}_{1-t}^R(x)] := \hat{C}(x)$ , the output from Algorithm 1 (main text).
  - 7: **else if** desire  $1 - \delta$  PAC-type guarantee **then**
  - 8: Input  $\{j \in \mathcal{D} : T_j = 1 - t\}$ ,  $\hat{\ell}(\cdot)$ ,  $\hat{u}(\cdot)$ , score  $V_{1-t}$ , target  $\alpha$ ,  $\delta$  to Algorithm 2 (main text).
  - 9: Obtain  $[\hat{C}_{1-t}^L(x), \hat{C}_{1-t}^R(x)] := \hat{C}(x)$ , the output from Algorithm 2 (main text).
  - 10: Output:  $[\hat{C}^L(x), \hat{C}^R(x)] = [y(1) - \hat{C}_0^R(x), y(1) - \hat{C}_0^L(x)]$  if  $t = 1$ ,
  - 11:  $[\hat{C}^L(x), \hat{C}^R(x)] = [\hat{C}_1^L(x) - y(0), \hat{C}_1^R(x) - y(0)]$  if  $t = 0$ .
- 

**Algorithm 4** Subroutine: One-sided robust prediction for ITE
 

---

- 1: Input: Calibration data  $\{(X_i, Y_i, T_i)\}_{i \in \mathcal{D}}$ , estimated propensity score  $\hat{e}(x)$ , estimated  $\hat{p}$  for  $\mathbb{P}(T = 1)$ , estimated quantile functions  $\hat{q}_\alpha^{(0)}(x)$ ,  $\hat{q}_{1-\alpha}^{(1)}(x)$ , test covariate  $x$ , test treatment  $t \in \{0, 1\}$ , test outcome  $y(t)$ , target level  $\alpha \in (0, 1)$ , confidence level  $\delta \in (0, 1)$  if necessary.
  - 2: Define the nonconformity scores  $V_0(x, y) = y - \hat{q}_\alpha^{(0)}(x)$  and  $V_1(x, y) = \hat{q}_{1-\alpha}^{(1)}(x) - y$  for  $t = 0, 1$ .
  - 3: Set  $\hat{\ell}(\cdot) = \frac{1}{\Gamma} \left[ \frac{1-\hat{p}}{\hat{p}} \frac{\hat{e}(x)}{1-\hat{e}(x)} \right]^{2t-1}$ ,  $\hat{u}(\cdot) = \Gamma \left[ \frac{1-\hat{p}}{\hat{p}} \frac{\hat{e}(x)}{1-\hat{e}(x)} \right]^{2t-1}$ .
  - 4: **if** desire marginal coverage **then**
  - 5: Input  $\{j \in \mathcal{D} : T_j = 1 - t\}$ ,  $\hat{\ell}(\cdot)$ ,  $\hat{u}(\cdot)$ , score  $V_{1-t}$ , test covariate  $x$ , and target  $\alpha$  to Algorithm 1 (main text).
  - 6: Obtain  $(-\infty, \hat{C}_0(x)] := \hat{C}(x)$  if  $t = 1$  or  $[\hat{C}_1(x), +\infty) := \hat{C}(x)$  if  $t = 0$ , from Algorithm 1 (main text).
  - 7: **else if** desire  $1 - \delta$  PAC-type guarantee **then**
  - 8: Input  $\{j \in \mathcal{D} : T_j = 1 - t\}$ ,  $\hat{\ell}(\cdot)$ ,  $\hat{u}(\cdot)$ , score  $V_{1-t}$ , test covariate  $x$ , target  $\alpha$ ,  $\delta$  to Algorithm 2 (main text).
  - 9: Obtain  $(-\infty, \hat{C}_0(x)] := \hat{C}(x)$  if  $t = 1$  or  $[\hat{C}_1(x), +\infty) := \hat{C}(x)$  if  $t = 0$ , from Algorithm 2 (main text).
  - 10: Output:  $[\hat{C}^L(x), +\infty) = [y(1) - \hat{C}_0(x), +\infty)$  if  $t = 1$ , and  $[\hat{C}^L(x), +\infty) = [\hat{C}_1^L(x) - y(0), +\infty)$  if  $t = 0$ .
- 

Consequently,

$$\frac{1}{\Gamma} \leq \frac{\mathbb{P}(Y(1) \in B \mid X = x, T = 1)}{\mathbb{P}(Y(1) \in B \mid X = x, T = 0)} \leq \Gamma$$

holds for  $\mathbb{P}$ -almost all  $x \in \mathcal{X}$ . By the arbitrariness of  $B$ , we have

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

69 Repeating the above steps for  $Y(0)$  we conclude the proof of Lemma 2.1. □

70 **B. Proof of Theorem 1.2.**

71 *Proof of Theorem 1.2.* Fixing any  $\tilde{\mathbb{P}} \in \mathcal{P}(\mathbb{P}, \ell, u)$ , we denote the likelihood ratio  $w(x, y) = \frac{d\tilde{\mathbb{P}}}{d\mathbb{P}}(x, y)$ . Recall that the  
 72 calibration data is  $\{(X_i, Y_i)\}_{i \in \mathcal{D}_{\text{calib}}}$  with  $\mathcal{D}_{\text{calib}} = \{1, \dots, n\}$ , and the test data point is  $(X_{n+1}, Y_{n+1}) \sim \tilde{\mathbb{P}}$ . We denote  
 73 the random variables  $Z_i = (X_i, Y_i)$  and realized values  $z_i = (x_i, y_i)$  for  $i = 1, \dots, n$ .

As a starting point, we elaborate on the weighted conformal inference introduced in (3), which paves the way for the analysis of marginal coverage later on. Following (3), the random variables  $\{Z_i\}_{i=1}^{n+1}$  are *weighted exchangeable*, meaning that the density of their joint distribution can be factorized as

$$f(z_1, \dots, z_{n+1}) = \prod_{i=1}^{n+1} w_i(z_i) \cdot g(z_1, \dots, z_{n+1}), \quad [5]$$

where  $g$  is some permutation-invariant function, i.e.,  $g(z_{\sigma(1)}, \dots, z_{\sigma(n+1)}) = g(z_1, \dots, z_{n+1})$  for any permutation  $\sigma$  of  $1, \dots, n+1$ . Specifically, here  $w_i(z) = 1$  for  $1 \leq i \leq n$  and  $w_{n+1}(z) = w(x, y)$ . For a set of values  $z_1, \dots, z_{n+1}$  where there may be repeated elements, we denote the unordered set  $z = [z_1, \dots, z_{n+1}]$  and the event

$$\mathcal{E}_z = \{[Z_1, \dots, Z_{n+1}] = [z_1, \dots, z_{n+1}]\}.$$

Let  $\Pi_{n+1}$  be the set of all permutations of  $\{1, \dots, n+1\}$ . Writing  $v_i = V(x_i, y_i) = V(z_i)$ , for each  $1 \leq i \leq n+1$ , it holds that

$$\mathbb{P}(V_{n+1} = v_i \mid \mathcal{E}_z) = \mathbb{P}(Z_{n+1} = z_i \mid \mathcal{E}_z) = \frac{\sum_{\sigma \in \Pi: \sigma(n+1)=i} f(z_{\sigma(1)}, \dots, z_{\sigma(n+1)})}{\sum_{\sigma \in \Pi} f(z_{\sigma(1)}, \dots, z_{\sigma(n+1)})},$$

where  $\mathbb{P}$  is induced by the joint distribution of  $\mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{calib}} \cup Z_{n+1}$ . By the factorization [5], we have

$$\frac{\sum_{\sigma \in \Pi: \sigma(n+1)=i} f(z_{\sigma(1)}, \dots, z_{\sigma(n+1)})}{\sum_{\sigma \in \Pi} f(z_{\sigma(1)}, \dots, z_{\sigma(n+1)})} = \frac{\sum_{\sigma \in \Pi: \sigma(n+1)=i} w_{n+1}(z_i) g(z_{\sigma(1)}, \dots, z_{\sigma(n+1)})}{\sum_{\sigma \in \Pi} w_{n+1}(z_{\sigma(n+1)}) g(z_{\sigma(1)}, \dots, z_{\sigma(n+1)})} = \frac{w_{n+1}(z_i)}{\sum_{j=1}^{n+1} w_{n+1}(z_j)}.$$

Therefore, the distribution of  $V_{n+1}$  conditional on the event  $\mathcal{E}_z$  is

$$V_{n+1} \mid \mathcal{E}_z \sim \sum_{i=1}^{n+1} \delta_{v_i} p_i^w, \quad \text{where } p_i^w = \frac{w_{n+1}(z_i)}{\sum_{j=1}^{n+1} w_{n+1}(z_j)} = \frac{w(x_i, y_i)}{\sum_{j=1}^{n+1} w(x_i, y_i)}.$$

Here  $\delta_v$  denotes the point mass at  $\{v\}$ . For any unordered set  $z = [z_1, \dots, z_n, z_{n+1}]$  and the corresponding  $k^*$  as defined in [B], it holds that

$$\mathbb{P}(Y_{n+1} \in \hat{C}(X_{n+1}) \mid \mathcal{E}_z) = \mathbb{P}(V_{n+1} \leq v_{[k^*]} \mid \mathcal{E}_z) = \sum_{i=1}^{n+1} p_i^w \mathbb{1}_{\{v_i \leq v_{[k^*]}\}} = \frac{\sum_{i=1}^{n+1} w(x_i, y_i) \mathbb{1}_{\{v_i \leq v_{[k^*]}\}}}{\sum_{j=1}^{n+1} w(x_i, y_i)}.$$

By the tower property of conditional expectations, we have

$$\mathbb{P}(Y_{n+1} \in \hat{C}(X_{n+1})) = \mathbb{E} \left[ \mathbb{P}(Y_{n+1} \in \hat{C}(X_{n+1}) \mid \mathcal{E}_z) \right] = \mathbb{E} \left[ \frac{\sum_{i=1}^{n+1} w(X_i, Y_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}}}{\sum_{j=1}^{n+1} w(X_i, Y_i)} \right]. \quad [6]$$

Equipped with the above preparations, we show the coverage guarantee in our setting. By definition [B], we know

$$V_{[k^*]} = \inf \left\{ v: \frac{\sum_{i=1}^n \hat{\ell}(X_i) \mathbb{1}_{\{V_i \leq v\}}}{\sum_{i=1}^n \hat{\ell}(X_i) \mathbb{1}_{\{V_i \leq v\}} + \sum_{i=1}^n \hat{u}(X_i) \mathbb{1}_{\{V_i > v\}} + \hat{u}(X_{n+1})} \geq 1 - \alpha \right\},$$

hence

$$\mathbb{E} \left[ \frac{\sum_{i=1}^n \hat{\ell}(X_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}}}{\sum_{i=1}^n \hat{\ell}(X_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} + \sum_{i=1}^n \hat{u}(X_i) \mathbb{1}_{\{V_i > V_{[k^*]}\}} + \hat{u}(X_{n+1})} \right] \geq 1 - \alpha$$

since the inner random variable is always no smaller than  $1 - \alpha$ . Combined with [6] and by the non-negativity of  $w(X_i, Y_i)$ , we have

$$\begin{aligned} & \mathbb{P}(Y_{n+1} \in \hat{C}(X_{n+1})) - (1 - \alpha) \\ & \geq \mathbb{E} \left[ \frac{\sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}}}{\sum_{i=1}^{n+1} w(X_i, Y_i)} \right] - \mathbb{E} \left[ \frac{\sum_{i=1}^n \hat{\ell}(X_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}}}{\sum_{i=1}^n \hat{\ell}(X_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} + \sum_{i=1}^n \hat{u}(X_i) \mathbb{1}_{\{V_i > V_{[k^*]}\}} + \hat{u}(X_{n+1})} \right] := \mathbb{E} \left[ \begin{array}{l} \text{(i)} \\ \text{(ii)} \end{array} \right], \end{aligned}$$

where we denote

$$\begin{aligned} \text{(i)} & = -w(X_{n+1}, Y_{n+1}) \cdot \sum_{i=1}^n \hat{\ell}(X_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} + \hat{u}(X_{n+1}) \cdot \sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} \\ & \quad + \left[ \sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} \right] \left[ \sum_{i=1}^n \hat{u}(X_i) \mathbb{1}_{\{V_i > V_{[k^*]}\}} \right] - \left[ \sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i > V_{[k^*]}\}} \right] \left[ \sum_{i=1}^n \hat{\ell}(X_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} \right] \end{aligned}$$

and

$$(ii) = \left[ \sum_{i=1}^{n+1} w(X_i, Y_i) \right] \left[ \sum_{i=1}^n \hat{\ell}(X_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} + \sum_{i=1}^n \hat{u}(X_i) \mathbb{1}_{\{V_i > V_{[k^*]}\}} + \hat{u}(X_{n+1}) \right].$$

We first establish a lower bound for the term (i). To this end, we define the random variables

$$\begin{aligned} \tilde{\ell}_i &= \max \{ w(X_i, Y_i), \hat{\ell}(X_i) \} \quad \text{and} \\ \tilde{u}_i &= \min \{ w(X_i, Y_i), \hat{u}(X_i) \}, \quad i = 1, \dots, n+1. \end{aligned}$$

By the above definition, it holds that  $0 \leq \tilde{u}_i \leq \hat{u}(X_i)$  and  $\tilde{\ell}_i \geq \hat{\ell}(X_i) \geq 0$ . We also define the differences

$$\begin{aligned} \Delta \tilde{\ell}_i &= \tilde{\ell}_i - w(X_i, Y_i) = [\hat{\ell}(X_i) - w(X_i, Y_i)]_+ \quad \text{and} \\ \Delta \tilde{u}_i &= \tilde{u}_i - w(X_i, Y_i) = -[\hat{u}(X_i) - w(X_i, Y_i)]_-, \quad i = 1, \dots, n+1. \end{aligned}$$

Using the above notation, we have

$$\begin{aligned} & \left[ \sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} \right] \left[ \sum_{i=1}^n \hat{u}(X_i) \mathbb{1}_{\{V_i > V_{[k^*]}\}} \right] - \left[ \sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i > V_{[k^*]}\}} \right] \left[ \sum_{i=1}^n \hat{\ell}(X_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} \right] \\ & \geq \left[ \sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} \right] \left[ \sum_{i=1}^n (w(X_i, Y_i) + \Delta \tilde{u}_i) \mathbb{1}_{\{V_i > V_{[k^*]}\}} \right] \\ & \quad - \left[ \sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i > V_{[k^*]}\}} \right] \left[ \sum_{i=1}^n (w(X_i, Y_i) + \Delta \tilde{\ell}_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} \right] \\ & = \left[ \sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} \right] \left[ \sum_{i=1}^n \Delta \tilde{u}_i \mathbb{1}_{\{V_i > V_{[k^*]}\}} \right] - \left[ \sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i > V_{[k^*]}\}} \right] \left[ \sum_{i=1}^n \Delta \tilde{\ell}_i \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} \right]. \end{aligned}$$

Following similar arguments, we have

$$\begin{aligned} & -w(X_{n+1}, Y_{n+1}) \cdot \sum_{i=1}^n \hat{\ell}(X_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} + \hat{u}(X_{n+1}) \cdot \sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} \\ & \geq -w(X_{n+1}, Y_{n+1}) \cdot \sum_{i=1}^n (w(X_i, Y_i) + \Delta \tilde{\ell}_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} + (w(X_{n+1}, Y_{n+1}) + \Delta \tilde{u}_{n+1}) \cdot \sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} \\ & = -w(X_{n+1}, Y_{n+1}) \cdot \sum_{i=1}^n \Delta \tilde{\ell}_i \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} + \Delta \tilde{u}_{n+1} \cdot \sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} \end{aligned}$$

By construction, we know that  $\Delta \tilde{\ell}_i \geq 0$  and  $\Delta \tilde{u}_i \leq 0$  for  $i = 1, \dots, n+1$ . Putting the lower bounds together, we obtain

$$\begin{aligned} (i) & \geq \left[ \sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} \right] \left[ \Delta \tilde{u}_{n+1} + \sum_{i=1}^n \Delta \tilde{u}_i \mathbb{1}_{\{V_i > V_{[k^*]}\}} \right] \\ & \quad - \left[ w(X_{n+1}, Y_{n+1}) + \sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i > V_{[k^*]}\}} \right] \left[ \sum_{i=1}^n \Delta \tilde{\ell}_i \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} \right] \\ & \geq \left[ \sum_{i=1}^n w(X_i, Y_i) \right] \left[ \sum_{i=1}^{n+1} \Delta \tilde{u}_i \right] - \left[ \sum_{i=1}^{n+1} w(X_i, Y_i) \right] \left[ \sum_{i=1}^n \Delta \tilde{\ell}_i \right]. \end{aligned}$$

Since the term (ii) is non-negative, we have the lower bound

$$\begin{aligned}
\mathbb{E} \left[ \begin{array}{l} \text{(i)} \\ \text{(ii)} \end{array} \right] &= \mathbb{E} \left[ \frac{\sum_{i=1}^n w(X_i, Y_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}}}{\sum_{i=1}^{n+1} w(X_i, Y_i)} - \frac{\sum_{i=1}^n \hat{\ell}(X_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}}}{\sum_{i=1}^n \hat{\ell}(X_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} + \sum_{i=1}^n \hat{u}(X_i) \mathbb{1}_{\{V_i > V_{[k^*]}\}} + \hat{u}(X_{n+1})} \right] \\
&\geq \mathbb{E} \left[ \frac{[\sum_{i=1}^n w(X_i, Y_i)] [\sum_{i=1}^{n+1} \Delta \tilde{u}_i] - [\sum_{i=1}^{n+1} w(X_i, Y_i)] [\sum_{i=1}^n \Delta \tilde{\ell}_i]}{[\sum_{i=1}^{n+1} w(X_i, Y_i)] [\sum_{i=1}^n \hat{\ell}(X_i) \mathbb{1}_{\{V_i \leq V_{[k^*]}\}} + \sum_{i=1}^n \hat{u}(X_i) \mathbb{1}_{\{V_i > V_{[k^*]}\}} + \hat{u}(X_{n+1})]} \right] \\
&\stackrel{\text{(a)}}{\geq} \mathbb{E} \left[ \frac{[\sum_{i=1}^n w(X_i, Y_i)] [\sum_{i=1}^{n+1} \Delta \tilde{u}_i] - [\sum_{i=1}^{n+1} w(X_i, Y_i)] [\sum_{i=1}^n \Delta \tilde{\ell}_i]}{[\sum_{i=1}^{n+1} w(X_i, Y_i)] [\sum_{i=1}^n \hat{\ell}(X_i)]} \right] \\
&\stackrel{\text{(b)}}{\geq} -\mathbb{E} \left[ \frac{\sum_{i=1}^{n+1} [\hat{u}(X_i) - w(X_i, Y_i)]_-}{\sum_{i=1}^n \hat{\ell}(X_i)} \right] - \mathbb{E} \left[ \frac{\sum_{i=1}^n [\hat{\ell}(X_i) - w(X_i, Y_i)]_+}{\sum_{i=1}^n \hat{\ell}(X_i)} \right].
\end{aligned}$$

Above, step (a) follows from the fact that  $\hat{u}(X_i) \geq \hat{\ell}(X_i) \geq 0$ , and step (b) is due to the non-negativity of  $w(X_i, Y_i)$ . By Hölder's inequality,

$$\begin{aligned}
\mathbb{E} \left[ \frac{\sum_{i=1}^n [\hat{\ell}(X_i) - w(X_i, Y_i)]_+}{\sum_{i=1}^n \hat{\ell}(X_i)} \right] &\leq \left\| \frac{1}{n} \sum_{i=1}^n (\hat{\ell}(X_i) - w(X_i, Y_i))_+ \right\|_p \cdot \left\| \frac{n}{\sum_{i=1}^n \hat{\ell}(X_i)} \right\|_q \\
&\stackrel{\text{(a)}}{\leq} \left\| (\hat{\ell}(X_i) - w(X_i, Y_i))_+ \right\|_p \cdot \left\| \frac{n}{\sum_{i=1}^n \hat{\ell}(X_i)} \right\|_q \\
&\stackrel{\text{(b)}}{\leq} \left\| (\hat{\ell}(X_i) - w(X_i, Y_i))_+ \right\|_p \cdot \left\| \frac{1}{\hat{\ell}(X_i)} \right\|_q,
\end{aligned}$$

where step (a) follows from Minkowski's inequality, and the step (b) follows from

$$\frac{n}{\sum_{i=1}^n \hat{\ell}(X_i)} \leq \frac{1}{n} \sum_{i=1}^n \frac{1}{\hat{\ell}(X_i)}$$

as implied by Cauchy-Schwarz inequality. Similarly,

$$\mathbb{E} \left[ \frac{\sum_{i=1}^{n+1} [\hat{u}(X_i) - w(X_i, Y_i)]_-}{\sum_{i=1}^n \hat{\ell}(X_i)} \right] \leq \left\| (\hat{u}(X_i) - u(X_i))_- \right\|_p \cdot \left\| \frac{1}{\hat{\ell}(X_i)} \right\|_q + \frac{1}{n} \left\| (\hat{u}(X_{n+1}) - u(X_{n+1}))_- \right\|_p \cdot \left\| \frac{1}{\hat{\ell}(X_i)} \right\|_q,$$

where the  $L_p$  norm for  $X_{n+1}$  is with respect to  $\tilde{\mathbb{P}}$ , hence

$$\frac{1}{n} \left\| (\hat{u}(X_{n+1}) - u(X_{n+1}))_- \right\|_p = \left\| \frac{w(X_i, Y_i)^{1/p}}{n} \cdot (\hat{u}(X_i) - u(X_i))_- \right\|_p$$

Combining the above results, we have

$$\mathbb{P}(Y_{n+1} \in \hat{C}(X_{n+1})) \geq 1 - \alpha - \hat{\Delta} \cdot \|1/\hat{\ell}(X_i)\|_q,$$

where

$$\hat{\Delta} = \left\| (\hat{\ell}(X_i) - \ell(X_i))_+ \right\|_p + \left\| (\hat{u}(X_i) - u(X_i))_- \right\|_p + \left\| \frac{w(X_i, Y_i)^{1/p}}{n} \cdot (\hat{u}(X_i) - u(X_i))_- \right\|_p,$$

74 which completes the proof of Theorem 1.2.  $\square$

75 **C. Proof of Theorem 3.1.**

76 *Proof of Theorem 3.1.* Throughout the proof, all statements are conditional on  $\mathcal{D}_{\text{train}}$ . By the independence of  
77  $\mathcal{D}_{\text{calib}} \cup \{(X_{n+1}, Y_{n+1})\}$ , the scores  $\{V(X_i, Y_i)\}_{i \in \mathcal{D}_{\text{calib}}}$  are i.i.d. and independent of  $V(X_{n+1}, Y_{n+1})$ .

Recall that  $G(\cdot)$  is defined in [13]. To begin with, we define

$$\hat{q} = \inf \{t: G(t) \geq 1 - \alpha\}, \quad \hat{q}_n = \inf \{t: \hat{G}_n(t) \geq 1 - \alpha\}.$$

For any fixed  $\epsilon > 0$ , we have

$$\begin{aligned} \mathbb{P}(\hat{q}_n \leq \hat{q} - \epsilon) &= \mathbb{P}(\hat{G}_n(\hat{q} - \epsilon) \geq 1 - \alpha) \\ &\leq \mathbb{P}(G(\hat{q} - \epsilon) \geq \hat{G}_n(\hat{q} - \epsilon) \geq 1 - \alpha) + \mathbb{P}(G(\hat{q} - \epsilon) < \hat{G}_n(\hat{q} - \epsilon)) \leq \delta. \end{aligned}$$

Here the last inequality follows from the fact that  $G(\hat{q} - \epsilon) < 1 - \alpha$  for any fixed  $\epsilon > 0$  and  $\mathbb{P}(G(\hat{q} - \epsilon) < \hat{G}_n(\hat{q} - \epsilon)) \leq \delta$  by [14] with  $t = \hat{q} - \epsilon$ . Therefore, by the continuity of probability measures, we have

$$\mathbb{P}(\hat{q}_n \geq \hat{q}) = 1 - \lim_{\epsilon \rightarrow 0^+} \mathbb{P}(\hat{q}_n \leq \hat{q} - \epsilon) \geq 1 - \delta.$$

Moreover, on the event  $\{\hat{q}_n \geq \hat{q}\}$ , by the definition of  $\hat{C}(X_{n+1})$ , it holds for any  $\tilde{\mathbb{P}} \in \mathcal{P}(\mathbb{P}, \ell, u)$  that

$$\begin{aligned} \tilde{\mathbb{P}}(Y_{n+1} \in \hat{C}(X_{n+1}) \mid \mathcal{D}_{\text{calib}}) &= \tilde{\mathbb{P}}(V(X_{n+1}, Y_{n+1}) \leq \hat{q}_n \mid \mathcal{D}_{\text{calib}}) \\ &\geq \tilde{\mathbb{P}}(V(X_{n+1}, Y_{n+1}) \leq \hat{q} \mid \mathcal{D}_{\text{calib}}) \\ &= \mathbb{E}[\mathbf{1}_{\{V(X, Y) \leq \hat{q}\}} w(X, Y) \mid \mathcal{D}_{\text{calib}}], \end{aligned}$$

where the expectation is with respect to  $(X, Y) \sim \mathbb{P}$  independent of  $\mathcal{D}_{\text{calib}}$ . By the definition of  $G(t)$

$$\hat{q} = \inf \left\{ t: \max \left\{ \mathbb{E}[\mathbf{1}_{\{V(X, Y) \leq t\}} \hat{\ell}(X)], 1 - \mathbb{E}[\mathbf{1}_{\{V(X, Y) > t\}} \hat{u}(X)] \right\} \geq 1 - \alpha \right\} = \min\{\hat{q}_1, \hat{q}_2\},$$

where  $(X, Y) \sim \mathbb{P}$  is an independent copy and we define

$$\begin{aligned} \hat{q}_1 &= \inf \left\{ t: \mathbb{E}[\mathbf{1}_{\{V(X, Y) \leq t\}} \hat{\ell}(X)] \geq 1 - \alpha \right\}, \\ \hat{q}_2 &= \inf \left\{ t: 1 - \mathbb{E}[\mathbf{1}_{\{V(X, Y) > t\}} \hat{u}(X)] \geq 1 - \alpha \right\}. \end{aligned}$$

For constants  $\hat{q}_1, \hat{q}_2$ , we have

$$\begin{aligned} &\mathbb{E}[\mathbf{1}_{\{V(X, Y) \leq \hat{q}\}} w(X, Y) \mid \mathcal{D}_{\text{calib}}] \\ &= \min \left\{ \mathbb{E}[\mathbf{1}_{\{V(X, Y) \leq \hat{q}_1\}} w(X, Y) \mid \mathcal{D}_{\text{calib}}], \mathbb{E}[\mathbf{1}_{\{V(X, Y) \leq \hat{q}_2\}} w(X, Y) \mid \mathcal{D}_{\text{calib}}] \right\}. \end{aligned}$$

We analyze the two terms separately. Firstly,

$$\begin{aligned} &\mathbb{E}[\mathbf{1}_{\{V(X, Y) \leq \hat{q}_1\}} w(X, Y) \mid \mathcal{D}_{\text{calib}}] \\ &= \mathbb{E}[\mathbf{1}_{\{V(X, Y) \leq \hat{q}_1\}} \hat{\ell}(X) \mid \mathcal{D}_{\text{calib}}] - \mathbb{E}[\mathbf{1}_{\{V(X, Y) \leq \hat{q}_1\}} (\hat{\ell}(X) - w(X, Y)) \mid \mathcal{D}_{\text{calib}}] \\ &\stackrel{(a)}{\geq} 1 - \alpha - \mathbb{E}[\mathbf{1}_{\{V(X, Y) \leq \hat{q}_1\}} (\hat{\ell}(X) - w(X, Y)) \mid \mathcal{D}_{\text{calib}}] \\ &\geq 1 - \alpha - \mathbb{E}[(\hat{\ell}(X) - w(X, Y))_+], \end{aligned}$$

where the step (a) follows from the definition of  $\hat{q}_1$ . Similarly,

$$\begin{aligned} &\mathbb{E}[\mathbf{1}_{\{V(X, Y) \leq \hat{q}_2\}} w(X, Y) \mid \mathcal{D}_{\text{calib}}] \\ &= 1 - \mathbb{E}[\mathbf{1}_{\{V(X, Y) > \hat{q}_2\}} w(X, Y) \mid \mathcal{D}_{\text{calib}}] \\ &= 1 - \mathbb{E}[\mathbf{1}_{\{V(X, Y) \leq \hat{q}_2\}} \hat{u}(X) \mid \mathcal{D}_{\text{calib}}] + \mathbb{E}[\mathbf{1}_{\{V(X, Y) \leq \hat{q}_2\}} (\hat{u}(X) - w(X, Y)) \mid \mathcal{D}_{\text{calib}}] \\ &\geq 1 - \alpha + \mathbb{E}[\mathbf{1}_{\{V(X, Y) \leq \hat{q}_2\}} (\hat{u}(X) - w(X, Y)) \mid \mathcal{D}_{\text{calib}}] \\ &\geq 1 - \alpha - \mathbb{E}[(\hat{u}(X) - w(X, Y))_-], \end{aligned}$$

where the first equality follows from the fact that  $w(x, y)$  is a likelihood ratio, and the first inequality follows from the definition of  $\hat{q}_2$ . Putting them together, on the event  $\{\hat{q}_n \geq \hat{q}\}$  which happens with probability at least  $1 - \delta$  with respect to  $\mathcal{D}_{\text{calib}}$ , it holds that

$$\tilde{\mathbb{P}}(Y_{n+1} \in \hat{C}(X_{n+1}) \mid \mathcal{D}_{\text{calib}}) \geq 1 - \alpha - \hat{\Delta},$$

where the gap is

$$\hat{\Delta} = \max \left\{ \mathbb{E} \left[ \left( \hat{\ell}(X) - w(X, Y) \right)_+ \right], \mathbb{E} \left[ \left( \hat{u}(X) - w(X, Y) \right)_- \right] \right\},$$

78 and the expectations are with respect to an independent copy  $(X, Y) \sim \mathbb{P}$ . Therefore, we conclude the proof of  
79 Theorem 3.1.  $\square$

#### 80 D. Proof of Proposition 2.1.

81 *Proof of Proposition 2.1.* Let  $t \in \mathbb{R}$  be any fixed constant or any random variable in  $\sigma(\mathcal{D}_{\text{train}})$ , and fix any  $\delta \in (0, 1)$ .  
82 We condition on  $\mathcal{D}_{\text{train}}$  throughout the proof, so that  $V$ ,  $\hat{\ell}$  and  $\hat{u}$  can be viewed as fixed and  $t$  can be viewed as  
83 constant.

Since  $\hat{G}_n(t) = \max \{ \hat{G}_n^\ell(t), \hat{G}_n^U(t) \}$ , by the definition [13], it suffices to show that

$$\begin{aligned} \mathbb{P}_{\mathcal{D}_{\text{calib}}} \left( \hat{G}_n^\ell(t) \leq \mathbb{E} \left[ \mathbf{1}_{\{V_i \leq t\}} \hat{\ell}(X_i) \right] \right) &\geq 1 - \delta/2 \\ \text{and } \mathbb{P}_{\mathcal{D}_{\text{calib}}} \left( \hat{G}_n^U(t) \leq 1 - \mathbb{E} \left[ \mathbf{1}_{\{V_i > t\}} \hat{u}(X_i) \right] \right) &\geq 1 - \delta/2, \end{aligned}$$

Now let

$$f(X_i) = \mathbf{1}_{\{V_i \leq t\}} \hat{\ell}(X_i)/M, \quad h(X_i) = 1 - \mathbf{1}_{\{V_i > t\}} \hat{u}(X_i)/M,$$

so that  $f(X_i)$  and  $h(X_i)$ ,  $1 \leq i \leq n$ , are i.i.d. random variables in  $[0, 1]$ . We rescale the bounds into  $\hat{f}_n = \hat{G}_n^\ell(t)/M$  and  $\hat{h}_n = (\hat{G}_n^U(t) - 1 + M)/M$ , hence by the tower property of conditional expectations, it suffices to show that

$$\mathbb{P}_{\mathcal{D}_{\text{calib}}} \left( \hat{f}_n \leq \mathbb{E} [f(X_i)] \right) \geq 1 - \delta/2 \quad \text{and} \quad \mathbb{P}_{\mathcal{D}_{\text{calib}}} \left( \hat{h}_n \leq \mathbb{E} [h(X_i)] \right) \geq 1 - \delta/2.$$

84 We show the desired result for  $f(\cdot)$  and that for  $h(\cdot)$  naturally applies.

Consider the filtration  $\{\mathcal{F}_i\}_{i \geq 1}$ , where the  $\sigma$ -algebra  $\mathcal{F}_i = \sigma\{X_1, Y_1, \dots, X_i, Y_i\}$ . Then by definition,  $\hat{\sigma}_i^L$  and  $\hat{\mu}_i^L$  are measurable with respect to  $\mathcal{F}_i$ , and  $\{\nu_i^L\}_{i \geq 1}$  is a predictable sequence with respect to  $\{\mathcal{F}_i\}_{i \geq 1}$ . Thus letting

$$f_0 = \mathbb{E} [f(X_i)] = \mathbb{E} \left[ \mathbf{1}_{\{V_i \leq t\}} \ell(X_i) \right] / M,$$

we have

$$\mathbb{E} [\mathcal{K}_i^\ell(f_0) \mid \mathcal{F}_{i-1}] = \mathcal{K}_{i-1}^\ell(f_0) \cdot \left( 1 + \nu_j^L \cdot \mathbb{E} [f(X_i) - f_0 \mid \mathcal{F}_{i-1}] \right) = \mathcal{K}_{i-1}^\ell(f_0).$$

Thus  $\{\mathcal{K}_i^\ell(f_0)\}_{i=1}^n$  is a martingale. Meanwhile, since  $f(X_i) \in [0, 1]$ , we know  $f(X_i) - f_0 \geq -1$ . Hence  $\mathcal{K}_i^\ell(f_0) \geq 0$  for all  $i \in [n]$ . Therefore by Ville's inequality,

$$\mathbb{P}_{\mathcal{D}_{\text{calib}}} \left( \max_{1 \leq i \leq n} \mathcal{K}_i^L(f_0) > \frac{2}{\delta} \right) \leq \frac{\delta}{2}.$$

Also note that  $\mathcal{K}_i(g)$  is a decreasing function in  $g \in \mathbb{R}$ , hence

$$\mathbb{P}_{\mathcal{D}_{\text{calib}}} \left( \hat{f}_n > \mathbb{E} [f(X_i)] \right) \leq \mathbb{P} \left( \max_{1 \leq i \leq n} \mathcal{K}_i(f_0) > \frac{2}{\delta} \right) \leq \delta/2,$$

85 Therefore, we conclude the proof of the desired results.  $\square$

86 **E. Proof of Proposition 3.2.**

87 *Proof of Proposition 3.2.* Throughout the proof, we denote the generic random variables  $(X, Y) \sim \mathbb{P}$  and  $V =$   
 88  $V(X, Y)$ . Denoting the right-hand side of [18] as  $F^*(t)$ , We are to show that

- 89 • (i)  $F^*(\cdot): \mathbb{R} \rightarrow [0, 1]$  is a distribution function.
- 90 • (ii)  $F^*(t)$  is a lower bound for  $F(t; \mathcal{P}(\mathbb{P}, \ell, u))$  for all  $t \in \mathbb{R}$ .
- 91 • (iii)  $F^*(t)$  can be achieved by some element in  $\mathcal{P}(\mathbb{P}, \ell, u)$ .

First of all, (i) is straightforward by noting that  $F^*(\cdot)$  is right-continuous due to the continuity of probability measures and  $\lim_{t \rightarrow -\infty} F^*(t) = 0$ ,  $\lim_{t \rightarrow +\infty} F^*(t) = 1$ . Also, similar to the arguments in the proof of Theorem 3.1, for any  $\tilde{\mathbb{P}} \in \mathcal{P}(\mathbb{P}, \ell, u)$ , we have

$$F(t; V, \tilde{\mathbb{P}}) = \mathbb{E} \left[ \mathbf{1}_{\{V(X, Y) \leq t\}} \frac{d\tilde{\mathbb{P}}}{d\mathbb{P}}(X, Y) \right] \geq \mathbb{E} \left[ \mathbf{1}_{\{V(X, Y) \leq t\}} \ell(X) \right],$$

$$F(t; V, \tilde{\mathbb{P}}) = 1 - \mathbb{E} \left[ \mathbf{1}_{\{V(X, Y) > t\}} \frac{d\tilde{\mathbb{P}}}{d\mathbb{P}}(X, Y) \right] \geq 1 - \mathbb{E} \left[ \mathbf{1}_{\{V(X, Y) > t\}} u(X) \right],$$

92 hence (ii) follows. For (iii), we are to construct one distribution  $\mathbb{P}^* \in \mathcal{P}(\mathbb{P}, \ell, u)$  so that  $F(\cdot; V, \mathbb{P}^*) = F^*(\cdot)$ .

If  $\mathbb{E}[u(X)] = 1$ , then since  $\mathbb{E}[w(X, Y)] = 1$  for the likelihood ratio function, the collection  $\mathcal{P}(\mathbb{P}, \ell, u) = \{\mathbb{P}^*\}$  is a singleton with  $d\mathbb{P}^*/d\mathbb{P}(x, y) = u(x)$ , and

$$F(t; \mathcal{P}(\mathbb{P}, \ell, u)) = F(t; \mathbb{P}^*) = \mathbb{E} \left[ \mathbf{1}_{\{V(X, Y) \leq t\}} u(X) \right] = F^*(t),$$

where the last equality follows from  $\ell(x) \leq u(x)$ . Then  $\mathbb{P}^*$  satisfies (iii). Similarly, for the case where  $\mathbb{E}[\ell(X)] = 1$ , the collection  $\mathcal{P}(\mathbb{P}, \ell, u) = \{\mathbb{P}^*\}$  is a singleton with  $d\mathbb{P}^*/d\mathbb{P}(x, y) = \ell(x)$ , and

$$F(t; \mathcal{P}(\mathbb{P}, \ell, u)) = F(t; \mathbb{P}^*) = \mathbb{E} \left[ \mathbf{1}_{\{V(X, Y) \leq t\}} \ell(X) \right] = F^*(t),$$

hence  $\tilde{\mathbb{P}}^*$  satisfies (iii). In the sequel, we consider the case where  $\mathbb{E}[\ell(X)] < 1 < \mathbb{E}[u(X)]$ . We define

$$H(t) := \mathbb{E} \left[ \ell(X) \cdot \mathbf{1}_{\{V \leq t\}} + u(X) \cdot \mathbf{1}_{\{V > t\}} \right].$$

By the construction and the continuity of probability measures, we have

$$\lim_{t \rightarrow \infty} H(t) = \mathbb{E}[\ell(X)] < 1 < \mathbb{E}[u(X)] = \lim_{t \rightarrow -\infty} H(t).$$

We additionally define  $t^* = \inf\{t \in \mathbb{R} : H(t) \leq 1\}$ . Note that  $t^* < \infty$  and  $H(t^*) \leq 1$  by the right continuity of  $H(t)$ . We also define the left limit of  $H(t)$  at  $t^*$  as

$$H_-(t^*) = \lim_{t \uparrow t^*} H(t),$$

and define the weight as

$$\gamma = \frac{1 - H(t^*)}{H_-(t^*) - H(t^*)} \mathbf{1}_{\{H_-(t^*) > 1\}}$$

Here by the definition of  $t^*$ , we have  $H_-(t^*) \geq 1 \geq H(t^*)$ , and the left limit takes the form

$$H_-(t^*) = \mathbb{E} \left[ \ell(X) \cdot \mathbf{1}_{\{V < t^*\}} + u(X) \cdot \mathbf{1}_{\{V \geq t^*\}} \right].$$

We now construct the worst-case distribution  $\mathbb{P}^*$  by

$$\frac{d\mathbb{P}^*}{d\mathbb{P}}(x, y) = \gamma \cdot \left[ \ell(x) \mathbf{1}_{\{V(x, y) < t^*\}} + u(x) \mathbf{1}_{\{V(x, y) \geq t^*\}} \right] + (1 - \gamma) \cdot \left[ \ell(x) \mathbf{1}_{\{V(x, y) \leq t^*\}} + u(x) \mathbf{1}_{\{V(x, y) > t^*\}} \right].$$

We denote the likelihood ratio  $w^*(x, y) = d\mathbb{P}^*/d\mathbb{P}(x, y)$  as constructed above. Note that

$$\mathbb{P}^*(\mathcal{X} \times \mathcal{Y}) = \mathbb{E} \left[ w^*(X, Y) \right] = \gamma \cdot H_-(t^*) + (1 - \gamma) \cdot H(t^*) = 1,$$

and also  $\ell(x) \leq w^*(x, y) \leq u(x)$  for all  $x \in \mathcal{X}$ . Therefore  $\mathbb{P}^*$  is a probability measure and is an element of  $\mathcal{P}(\mathbb{P}, \ell, u)$ . In the following, we check that  $F(t; V, \mathbb{P}^*) = F^*(t)$  for all  $t \in \mathbb{R}$ . Recall that we work with the case  $\mathbb{E}[\ell(X)] < 1 < \mathbb{E}[u(X)]$ . For any constant  $t < t^*$ , by the construction of  $w^*$ ,

$$\begin{aligned} F(t; V, \mathbb{P}^*) &= \mathbb{E}\left[w^*(X, Y) \mathbf{1}_{\{V(X, Y) \leq t\}}\right] \\ &= \gamma \cdot \mathbb{E}\left[\ell(X) \mathbf{1}_{\{V(X, Y) \leq t\}}\right] + (1 - \gamma) \cdot \mathbb{E}\left[\ell(X) \mathbf{1}_{\{V(X, Y) \leq t\}}\right] \\ &= \mathbb{E}\left[\ell(X) \cdot \mathbf{1}_{\{V \leq t\}}\right] = F^*(t), \end{aligned}$$

where the last equality follows from the fact that

$$\mathbb{E}\left[\ell(X) \cdot \mathbf{1}_{\{V \leq t\}}\right] > 1 - \mathbb{E}\left[u(X) \cdot \mathbf{1}_{\{V > t\}}\right]$$

since  $H(t) > 1$  for  $t < t^*$ . When  $t = t^*$ , note that

$$\mathbb{E}\left[\ell(X) \cdot \mathbf{1}_{\{V \leq t^*\}}\right] - 1 + \mathbb{E}\left[u(X) \cdot \mathbf{1}_{\{V > t^*\}}\right] = H(t^*) - 1 \leq 0,$$

hence the right-hand side of [18] admits the form

$$F^*(t^*) = 1 - \mathbb{E}\left[u(X) \cdot \mathbf{1}_{\{V > t^*\}}\right].$$

Meanwhile, by the construction of  $w^*(x, y)$ , we have

$$\begin{aligned} F(t^*; V, \mathbb{P}^*) &= \mathbb{E}\left[w^*(X, Y) \mathbf{1}_{\{V(X, Y) \leq t^*\}}\right] \\ &= \gamma \cdot \mathbb{E}\left[\ell(X) \mathbf{1}_{\{V(X, Y) < t^*\}} + u(X) \mathbf{1}_{\{V(X, Y) = t^*\}}\right] + (1 - \gamma) \cdot \mathbb{E}\left[\ell(X) \mathbf{1}_{\{V(X, Y) \leq t^*\}}\right] \\ &= \mathbb{E}\left[\ell(X) \mathbf{1}_{\{V(X, Y) \leq t^*\}}\right] - \gamma \cdot \left(\mathbb{E}\left[\ell(X) \mathbf{1}_{\{V(X, Y) = t^*\}}\right] - u(X) \cdot \mathbf{1}_{\{V(X, Y) = t^*\}}\right) \\ &= \mathbb{E}\left[\ell(X) \mathbf{1}_{\{V(X, Y) \leq t^*\}}\right] - \gamma \cdot (H_-(t^*) - H(t^*)) \\ &= \mathbb{E}\left[\ell(X) \mathbf{1}_{\{V(X, Y) \leq t^*\}}\right] + 1 - H(t^*) = 1 - \mathbb{E}\left[u(X) \cdot \mathbf{1}_{\{V > t^*\}}\right] = F^*(t^*). \end{aligned}$$

Similarly, when  $t > t^*$ , by the construction of  $w^*(x, y)$  we have

$$\begin{aligned} F(t^*; V, \mathbb{P}^*) &= 1 - \mathbb{E}\left[w^*(X, Y) \mathbf{1}_{\{V(X, Y) > t\}}\right] \\ &= 1 - \gamma \cdot \mathbb{E}\left[u(X) \mathbf{1}_{\{V(X, Y) > t\}}\right] + (1 - \gamma) \cdot \mathbb{E}\left[u(X) \mathbf{1}_{\{V(X, Y) > t\}}\right] \\ &= 1 - \mathbb{E}\left[u(X) \mathbf{1}_{\{V(X, Y) > t\}}\right] = F^*(t), \end{aligned}$$

93 where the last equality follows from the fact that  $H(t) \leq 1$  thus  $1 - \mathbb{E}\left[u(X) \mathbf{1}_{\{V(X, Y) > t\}}\right] \geq \mathbb{E}\left[\ell(X) \mathbf{1}_{\{V(X, Y) \leq t\}}\right]$ .  
 94 Combining the three cases, we arrive at  $F^*(\cdot) = F(\cdot; V, \mathbb{P}^*)$ , hence [18] follows and we conclude the proof of Proposi-  
 95 tion 3.2.  $\square$

### 96 F. Proof of Proposition 3.3.

97 *Proof of Proposition 3.3.* The proof proceeds by showing that  $\mathcal{P} \subset \mathcal{P}(\mathbb{P}, f, \ell_0, u_0)$  and  $\mathcal{P}(\mathbb{P}, f, \ell_0, u_0) \subset \mathcal{P}$ , which  
 98 together lead to the desired result.

**Step 1:**  $\mathcal{P} \subset \mathcal{P}(\mathbb{P}, f, \ell_0, u_0)$ . Let  $\mathbb{P}^{\text{sup}}$  be any super-population that satisfies [3] and [19]. Due to the partial observation of potential outcomes, the observed distribution admits the decomposition

$$\mathbb{P}_{X, Y, T}^{\text{obs}} = \mathbb{P}_{T=1}^{\text{obs}} \times \mathbb{P}_{X, Y(1) | T=1}^{\text{obs}} + \mathbb{P}_{T=0}^{\text{obs}} \times \mathbb{P}_{X, Y(0) | T=0}^{\text{obs}}.$$

Therefore, the data-compatibility condition [19] is equivalent to

$$\mathbb{P}_T^{\text{sup}} = \mathbb{P}_T^{\text{obs}}, \quad \mathbb{P}_{X, Y(1) | T=1}^{\text{sup}} = \mathbb{P}_{X, Y(1) | T=1}^{\text{obs}}, \quad \mathbb{P}_{X, Y(0) | T=0}^{\text{sup}} = \mathbb{P}_{X, Y(0) | T=0}^{\text{obs}}$$

where the latter two are further equivalent to

$$\mathbb{P}_{X | T=w}^{\text{sup}} = \mathbb{P}_{X | T=w}^{\text{obs}}, \quad \mathbb{P}_{Y(1) | X, T=w}^{\text{sup}} = \mathbb{P}_{Y(1) | X, T=w}^{\text{obs}}, \quad w \in \{0, 1\}. \quad [7]$$

Recall that  $\tilde{\mathbb{P}} = \mathbb{P}_{X, Y(1) | T=0}^{\text{sup}}$  and  $\mathbb{P} = \mathbb{P}_{X, Y(1) | T=1}^{\text{obs}}$ . Then we have

$$\frac{d\tilde{\mathbb{P}}_X}{d\mathbb{P}_X} = \frac{d\mathbb{P}_X^{\text{sup}} | T=0}{d\mathbb{P}_X^{\text{obs}} | T=1} = \frac{d\mathbb{P}_X^{\text{obs}} | T=0}{d\mathbb{P}_X^{\text{obs}} | T=1} = \frac{\mathbb{P}^{\text{obs}}(T=1) \cdot \mathbb{P}^{\text{obs}}(T=0)}{\mathbb{P}^{\text{obs}}(T=0) \cdot \mathbb{P}^{\text{obs}}(T=1)} = f(X),$$

where the second equality follows from [7] and the third equality follows from the Bayes rule. On the other hand, the shift of conditional distribution is

$$\frac{d\tilde{\mathbb{P}}_{Y(1)|X}}{d\mathbb{P}_{Y(1)|X}} = \frac{d\mathbb{P}_{Y(1)|X,T=0}^{\text{sup}}}{d\mathbb{P}_{Y(1)|X,T=1}^{\text{obs}}} \in [1/\Gamma, \Gamma]$$

99 according to Lemma 2.1. Therefore,  $\mathbb{P}^{\text{sup}} \in \mathcal{P}(\mathbb{P}, f, \ell_0, u_0)$  by the definition. Hence we have  $\mathcal{P} \subset \mathcal{P}(\mathbb{P}, f, \ell_0, u_0)$ .

**Step 2:**  $\mathcal{P}(\mathbb{P}, f, \ell_0, u_0) \subset \mathcal{P}$ . For this part, we are to show that for any  $\tilde{\mathbb{P}} \in \mathcal{P}(\mathbb{P}, f, \ell_0, u_0)$ , there exists some  $\mathbb{P}^{\text{sup}}$  satisfying [3] and [19] such that  $\tilde{\mathbb{P}} = \mathbb{P}_{X,Y(1)|T=0}^{\text{sup}}$ . Fixing an arbitrary probability distribution  $\tilde{\mathbb{P}} \in \mathcal{P}(\mathbb{P}, f, \ell_0, u_0)$ , we define the function

$$w(y|x) = \frac{d\tilde{\mathbb{P}}_{Y(1)|X}}{d\mathbb{P}_{Y(1)|X}}(y|x),$$

so that  $w(y|x) \in [1/\Gamma, \Gamma]$  for  $\mathbb{P}^{\text{obs}}$ -almost all  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ . Also, since  $\tilde{\mathbb{P}}$  is a distribution, we have

$$\mathbb{E}[w(Y(1)|x) | X = x] = \tilde{\mathbb{P}}(Y(1) \in \mathcal{Y} | X = x) = 1$$

for  $\mathbb{P}$ -almost all  $x \in \mathcal{X}$ , and the conditional expectation is induced by  $(X, Y(1)) \sim \mathbb{P} = \mathbb{P}_{X,Y(1)|T=1}^{\text{obs}}$ . Applying Lemma 3.1 with  $r(x, y) = w(y|x)$  and  $t = 1$ , we know there exists a distribution  $\mathbb{P}^{\text{sup}}$  over  $(X, Y(0), Y(1), U, T)$  for some confounder  $U$  that satisfies [19] and [3], and

$$\frac{d\mathbb{P}_{Y(1),X|T=0}^{\text{sup}}}{d\mathbb{P}_{Y(1),X|T=1}^{\text{sup}}}(y|x) = w(y|x).$$

Since  $\mathbb{P}^{\text{sup}}$  satisfies [19], we have  $\mathbb{P}_{Y(1),X|T=1}^{\text{sup}} = \mathbb{P}_{Y(1),X|T=1}^{\text{obs}} = \mathbb{P}_{Y(1)|X}$  where we recall the definition of  $\mathbb{P}$ , the distribution at hand. Hence

$$w(y|x) = \frac{d\mathbb{P}_{Y(1),X|T=0}^{\text{sup}}}{d\mathbb{P}_{Y(1)|X}}(y|x) = \frac{d\tilde{\mathbb{P}}_{Y(1)|X}}{d\mathbb{P}_{Y(1)|X}}(y|x).$$

Therefore, we have  $\tilde{\mathbb{P}}_{Y(1)|X} = \mathbb{P}_{Y(1),X|T=0}^{\text{sup}}$ . Furthermore, since  $\mathbb{P}^{\text{sup}}$  satisfies [19], we know

$$\frac{d\mathbb{P}_X|T=0}^{\text{sup}}}{d\mathbb{P}_X} = \frac{d\mathbb{P}_X|T=0}^{\text{obs}}}{d\mathbb{P}_X|T=1}^{\text{obs}} = f(X) = \frac{d\tilde{\mathbb{P}}_X}{d\mathbb{P}_X},$$

100 where the second equality follows from the Bayes rule and the last equality follows from the fact that  $\tilde{\mathbb{P}} \in \mathcal{P}(\mathbb{P}, f, \ell_0, u_0)$ .  
 101 Thus, we have  $\tilde{\mathbb{P}}_X = \mathbb{P}_X^{\text{sup}}|T=0$ . Putting the two parts together, we have  $\tilde{\mathbb{P}} = \mathbb{P}_{X,Y(1)|T=0}^{\text{sup}}$ . By the arbitrariness of  $\tilde{\mathbb{P}}$ ,  
 102 we arrive at  $\mathcal{P}(\mathbb{P}, f, \ell_0, u_0) \subset \mathcal{P}$ .

103 Combining the two steps, we conclude the proof of Proposition 3.3.  $\square$

## 104 G. Sharpness of the identification set.

**Lemma 3.1** (Sharpness of Lemma 2.1). *Given  $t \in \{0, 1\}$ , a marginal distribution  $\mathbb{P}^{\text{obs}}$  over  $(X, Y, T)$  and a function  $r(x, y) \in [1/\Gamma, \Gamma]$  such that*

$$\mathbb{E}^{\text{obs}}[r(X, Y(t)) | X, T = t] = 1, \quad \mathbb{P}^{\text{obs}}\text{-almost surely,}$$

105 *there exists a distribution  $\mathbb{P}^{\text{sup}}$  over  $(X, Y(0), Y(1), U, T)$  for some confounder  $U$  such that*

- 106 • (i)  $\mathbb{P}_{X,Y,T}^{\text{sup}} = \mathbb{P}_{X,Y,T}^{\text{obs}}$  for  $Y = Y(T)$ ;
- 107 • (ii)  $\mathbb{P}^{\text{sup}}$  satisfies the marginal  $\Gamma$ -selection condition;
- 108 • (iii) the likelihood ratio is exactly  $r(x, y)$ , so that  $r(x, y) = \frac{d\mathbb{P}_{Y(t)|X,T=1-t}^{\text{sup}}}{d\mathbb{P}_{Y(t)|X,T=t}^{\text{sup}}}(x, y)$  for  $\mathbb{P}^{\text{sup}}$ -almost all  $x, y$ .

109 *hold simulatenously.*

110 *Proof of Lemma 3.1.* Fix any marginal distribution  $\mathbb{P}^{\text{obs}}$  over  $(X, Y, T)$  for  $Y = Y(T)$ , and a function  $r(x, y) \in [1/\Gamma, \Gamma]$   
 111 satisfying the given condition. We show the result for  $t = 1$ , while that for  $t = 0$  follows exactly the same arguments.

**The construction of  $\mathbb{P}^{\text{sup}}$**  To begin with, we let the confounder be the counterfactual itself, so that  $U = Y(1)$ . The joint distribution is thus

$$\mathbb{P}_{(X, Y(1), Y(0), T)}^{\text{sup}} = \mathbb{P}_{X, T}^{\text{sup}} \times \mathbb{P}_{(Y(1), Y(0)) | X, T}^{\text{sup}}.$$

We set the two parts separately. Firstly, we set  $\mathbb{P}_{X, T}^{\text{sup}} = \mathbb{P}_{X, T}^{\text{obs}}$  for the distribution on  $(X, T)$ . The joint distribution of potential outcomes given  $(X, T)$  admits

$$\mathbb{P}_{(Y(1), Y(0)) | X, T}^{\text{sup}} = \mathbb{P}_{Y(1) | X, T}^{\text{sup}} \times \mathbb{P}_{Y(0) | X, T, Y(1)}^{\text{sup}},$$

Since our target is for  $Y(1)$ , we take a simple coupling where  $Y(0)$  is independent of  $(T, Y(1))$  conditional on  $X$ , so that we set

$$\mathbb{P}_{Y(0) | X, T, Y(1)}^{\text{sup}} = \mathbb{P}_{Y(0) | X}^{\text{sup}} = \mathbb{P}_{Y(0) | X, T=0}^{\text{sup}} = \mathbb{P}_{Y(0) | X, T=0}^{\text{obs}}. \quad [8]$$

On the other hand, we set  $\mathbb{P}_{Y(1) | X, T}^{\text{sup}}$  for  $T = 0, 1$  by

$$\mathbb{P}_{Y(1) | X, T=1}^{\text{sup}} = \mathbb{P}_{Y(1) | X, T=1}^{\text{obs}} \quad \text{and} \quad \frac{d\mathbb{P}_{Y(1) | X, T=0}^{\text{sup}}}{d\mathbb{P}_{Y(1) | X, T=1}^{\text{obs}}}(y | x) = r(x, y). \quad [9]$$

So far we've completed the pieces of constructing  $\mathbb{P}^{\text{sup}}$ . It remains to check that it is indeed a probability measure. By construction,  $\mathbb{P}_{X, T}^{\text{sup}} = \mathbb{P}_{X, T}^{\text{obs}}$  is a probability measure on  $(X, T)$ . Also, by the construction,

$$\mathbb{P}^{\text{sup}}(Y(1) \in \mathcal{Y} | X, T = 1) = \mathbb{P}^{\text{obs}}(Y(1) \in \mathcal{Y} | X, T = 1) = 1,$$

and

$$\begin{aligned} \mathbb{P}^{\text{sup}}(Y(1) \in \mathcal{Y} | X, T = 0) &= \mathbb{E}^{\text{obs}} \left[ \mathbf{1}_{\{Y(1) \in \mathcal{Y}\}} \frac{d\mathbb{P}_{Y(1) | X, T=0}^{\text{sup}}}{d\mathbb{P}_{Y(1) | X, T=1}^{\text{obs}}} \middle| X, T = 1 \right] \\ &= \mathbb{E}^{\text{obs}} \left[ \frac{d\mathbb{P}_{Y(1) | X, T=0}^{\text{sup}}}{d\mathbb{P}_{Y(1) | X, T=1}^{\text{obs}}} \middle| X, T = 1 \right] = \mathbb{E}^{\text{obs}} [r(X, Y(1)) | X, T = 1] = 1, \end{aligned}$$

where the first equality is by the change-of-measure formula, the second equality follows from  $1 = \mathbf{1}_{\{Y(1) \in \mathcal{Y}\}}$ , the third equality follows from the construction, and the last equality is the given condition on  $r(x, y)$ . Thus,  $\mathbb{P}_{Y(1) | X, T}^{\text{sup}}$  is a probability measure. Also, by the construction of  $\mathbb{P}_{Y(0) | X, T, Y(1)}^{\text{sup}}$ , we have

$$\mathbb{P}^{\text{sup}}(Y(0) \in \mathcal{Y} | X, T, Y(1)) = \mathbb{P}^{\text{obs}}(Y(0) \in \mathcal{Y} | X, T = 0) = 1,$$

hence  $\mathbb{P}_{Y(0) | X, T, Y(1)}^{\text{sup}}$  is also a probability measure. Putting them together, we know that

$$\mathbb{P}_{X, T, Y(1), Y(0)}^{\text{sup}} = \mathbb{P}_{X, T}^{\text{sup}} \times \mathbb{P}_{Y(1) | X, T}^{\text{sup}} \times \mathbb{P}_{Y(0) | X, T, Y(1)}^{\text{sup}}$$

112 is indeed a probability measure over  $(X, T, Y(1), Y(0))$ .

**Verify the properties** We now proceed to verify the three stated properties. For (i), due to partial observability, the observed distribution admits the decomposition

$$\mathbb{P}_{X, Y, T}^{\text{obs}} = \mathbb{P}_T^{\text{obs}} \times \mathbb{P}_{X, Y | T}^{\text{obs}} = \mathbb{P}_{T=1}^{\text{obs}} \times \mathbb{P}_{X, Y(1) | T=1}^{\text{obs}} + \mathbb{P}_{T=0}^{\text{obs}} \times \mathbb{P}_{X, Y(0) | T=0}^{\text{obs}}.$$

Similarly, the projection on the observable of  $\mathbb{P}^{\text{sup}}$  is

$$\mathbb{P}_{X, Y, T}^{\text{sup}} = \mathbb{P}_T^{\text{sup}} \times \mathbb{P}_{X, Y | T}^{\text{sup}} = \mathbb{P}_{T=1}^{\text{sup}} \times \mathbb{P}_{X, Y(1) | T=1}^{\text{sup}} + \mathbb{P}_{T=0}^{\text{sup}} \times \mathbb{P}_{X, Y(0) | T=0}^{\text{sup}}.$$

Here by the construction of  $\mathbb{P}_{X, T}^{\text{sup}} = \mathbb{P}_{X, T}^{\text{obs}}$ , we have  $\mathbb{P}_{T=w}^{\text{obs}} = \mathbb{P}_{T=w}^{\text{sup}}$  and  $\mathbb{P}_{X | T=w}^{\text{sup}} = \mathbb{P}_{X | T=w}^{\text{obs}}$  for  $w \in \{0, 1\}$ . Also,  $\mathbb{P}_{Y(w) | X, T=w}^{\text{sup}} = \mathbb{P}_{Y(w) | X, T=w}^{\text{obs}}$  holds for  $w \in \{0, 1\}$  by [8] and [9]. The equivalences altogether leads to  $\mathbb{P}_{X, Y, T}^{\text{sup}} = \mathbb{P}_{X, Y, T}^{\text{obs}}$ . For (ii), by the Bayes rule, we have

$$r(x, y) = \frac{d\mathbb{P}_{Y(1) | X, T=0}^{\text{sup}}}{d\mathbb{P}_{Y(1) | X, T=1}^{\text{sup}}}(y | x) = \frac{\mathbb{P}^{\text{sup}}(T = 0 | X = x, Y(1) = y)}{\mathbb{P}^{\text{sup}}(T = 1 | X = x, Y(1) = y)} \cdot \frac{\mathbb{P}^{\text{sup}}(T = 1 | X = x)}{\mathbb{P}^{\text{sup}}(T = 0 | X = x)} \in [1/\Gamma, \Gamma],$$

113 so  $\mathbb{P}^{\text{sup}}$  satisfies the marginal  $\Gamma$ -selection condition [3] hence (ii) is verified. Property (iii) has also been verified as  
114 above. So far, we've constructed  $\mathbb{P}^{\text{sup}}$  that satisfies all stated conditions and we conclude the proof of Lemma 3.1.  $\square$

115 **H. Proof of Proposition 3.4.**

116 *Proof of Proposition 3.4.* For simplicity, we denote [21] as  $F^*(t)$ , and aim to show that

- 117 • (i)  $F^*(t)$  is a distribution function;
- 118 • (ii)  $F^*(t)$  is a lower bound for  $F(t; V, \tilde{\mathbb{P}})$  for all  $t \in \mathbb{R}$  and all  $\tilde{\mathbb{P}} \in \mathcal{P}(\mathbb{P}, f, \ell_0, u_0)$ ;
- 119 • (iii)  $F^*(t)$  can be achieved by  $\mathbb{P}^* \in \mathcal{P}(\mathbb{P}, f, \ell_0, u_0)$  where  $d\mathbb{P}^*/d\mathbb{P}(x, y) = w^*(x, y)$  as defined in [3.4].

To verify (i), we note that  $F^*(t)$  is right continuous by the continuity of probability measures, as well as  $\lim_{t \rightarrow -\infty} F^*(t) = 0$  and  $\lim_{t \rightarrow +\infty} F^*(t) = 1$ . To show (ii), we are to show that for any  $\tilde{\mathbb{P}} \in \mathcal{P}(\mathbb{P}, f, \ell_0, u_0)$ , it holds  $\mathbb{P}$ -almost surely that

$$\mathbb{E}\left[\mathbf{1}_{\{V(X, Y) \leq t\}} w^*(X, Y) \mid X\right] \leq \tilde{\mathbb{P}}(V(X, Y) \leq t \mid X) = \mathbb{E}\left[\mathbf{1}_{\{V(X, Y) \leq t\}} \frac{d\tilde{\mathbb{P}}}{d\mathbb{P}}(X, Y) \mid X\right].$$

Fixing any  $\tilde{\mathbb{P}} \in \mathcal{P}(\mathbb{P}, f, \ell_0, u_0)$ , we denote the conditional likelihood as  $w_0(y \mid x) = d\tilde{\mathbb{P}}_{Y \mid X}/d\mathbb{P}_{Y \mid X}(y \mid x)$ , so that  $\ell_0(x) \leq w_0(y \mid x) \leq u_0(x)$  for  $\mathbb{P}$ -almost all  $x \in \mathcal{X}$ . Then the marginal likelihood ratio is  $w(x, y) := d\tilde{\mathbb{P}}/d\mathbb{P}(x, y) = f(x) \cdot w_0(y \mid x)$ . Hence for any  $t \in \mathbb{R}$ ,

$$\begin{aligned} & \mathbb{E}\left[\mathbf{1}_{\{V(X, Y) \leq t\}} w^*(X, Y) \mid X\right] - \tilde{\mathbb{P}}(V(X, Y) \leq t \mid X) \\ &= \mathbb{E}\left[\mathbf{1}_{\{V(X, Y) \leq t\}} w^*(X, Y) \mid X\right] - \mathbb{E}\left[\mathbf{1}_{\{V(X, Y) \leq t\}} w(X, Y) \mid X\right] \\ &= \mathbb{E}\left[\mathbf{1}_{\{V(X, Y) \leq t\}} (w^*(X, Y) - w(X, Y)) \mid X\right] \cdot \mathbf{1}_{\{t < q(\tau(X); \mathbb{X}, \mathbb{P})\}} \\ & \quad + \mathbb{E}\left[\mathbf{1}_{\{V(X, Y) \leq t\}} (w^*(X, Y) - w(X, Y)) \mid X\right] \cdot \mathbf{1}_{\{t \geq q(\tau(X); \mathbb{X}, \mathbb{P})\}}. \end{aligned}$$

We treat the two terms in the last summation separately. By the definition of  $w^*(x, y)$ ,

$$\begin{aligned} & \mathbb{E}\left[\mathbf{1}_{\{V(X, Y) \leq t\}} (w^*(X, Y) - w(X, Y)) \mid X\right] \cdot \mathbf{1}_{\{t < q(\tau(X); \mathbb{X}, \mathbb{P})\}} \\ &= f(X) \cdot \mathbb{E}\left[\mathbf{1}_{\{V(X, Y) \leq t\}} (\ell_0(X) - w_0(Y \mid X)) \mid X\right] \cdot \mathbf{1}_{\{t < q(\tau(X); \mathbb{X}, \mathbb{P})\}} \leq 0. \end{aligned}$$

Meanwhile, since  $1 = \mathbb{E}[w_0(Y \mid X) \mid X] = \mathbb{E}[w^*(X, Y)/f(X) \mid X] = 1$  holds  $\mathbb{P}$ -almost surely, we have

$$\begin{aligned} & \mathbb{E}\left[\mathbf{1}_{\{V(X, Y) \leq t\}} (w^*(X, Y) - w(X, Y)) \mid X\right] \cdot \mathbf{1}_{\{t \geq q(\tau(X); \mathbb{X}, \mathbb{P})\}} \\ &= f(X) \cdot \mathbb{E}\left[\mathbf{1}_{\{V(X, Y) > t\}} (w_0(Y \mid X) - w^*(X, Y)/f(X)) \mid X\right] \cdot \mathbf{1}_{\{t \geq q(\tau(X); \mathbb{X}, \mathbb{P})\}} \\ &= f(X) \cdot \mathbb{E}\left[\mathbf{1}_{\{V(X, Y) > t\}} (w_0(Y \mid X) - u_0(X)) \mid X\right] \cdot \mathbf{1}_{\{t \geq q(\tau(X); \mathbb{X}, \mathbb{P})\}} \leq 0. \end{aligned}$$

Summing them up and by the tower property of conditional expectations, it holds for any  $t \in \mathbb{R}$  that

$$F^*(t) = \mathbb{E}\left[\mathbf{1}_{\{V(X, Y) \leq t\}} w^*(X, Y)\right] \leq \tilde{\mathbb{P}}(V(X, Y) \leq t),$$

which verifies property (ii). Finally, we define  $\mathbb{P}^*$  by  $d\mathbb{P}^*/d\mathbb{P}(x, y) = w^*(x, y)$ . Then since  $\mathbb{E}[w^*(X, Y)/f(X) \mid X = x] = 1$ , the marginal likelihood ratio satisfies  $d\mathbb{P}_X^*/d\mathbb{P}_X = f(x)$ , hence  $d\mathbb{P}_{Y \mid X}^*/d\mathbb{P}_{Y \mid X}(x, y) = w^*(x, y)/f(x)$ . To verify  $\mathbb{P}^* \in \mathcal{P}(\mathbb{P}, f, \ell_0, u_0)$ , it remains to show  $\ell_0(x) \leq \gamma_0(x) \leq u_0(x)$  when it is nonzero, i.e.,  $\mathbb{P}(V(x, Y) = q(\tau(x); x, \mathbb{P}) \mid X = x) > 0$ . In this case,

$$\gamma_0(x) = \frac{1 - \ell_0(x) \cdot \mathbb{P}(V(x, Y) < q(\tau(x); x, \mathbb{P}) \mid X = x) - u_0(x) \cdot \mathbb{P}(V(x, Y) > q(\tau(x); x, \mathbb{P}) \mid X = x)}{\mathbb{P}(V(x, Y) = q(\tau(x); x, \mathbb{P}) \mid X = x)}.$$

Note that  $\mathbb{P}(V(x, Y) < q(\tau(x); x, \mathbb{P}) \mid X = x) \leq \tau(x) = (u_0(x) - 1)/(u_0(x) - \ell_0(x))$ , hence

$$\begin{aligned} & 1 - \ell_0(x) \cdot \mathbb{P}(V(x, Y) < q(\tau(x); x, \mathbb{P}) \mid X = x) \\ & \leq u_0(x) - u_0(x) \cdot \mathbb{P}(V(x, Y) < q(\tau(x); x, \mathbb{P}) \mid X = x) \\ & = u_0(x) \cdot \mathbb{P}(V(x, Y) \geq q(\tau(x); x, \mathbb{P}) \mid X = x) \\ & = u_0(x) \cdot \mathbb{P}(V(x, Y) > q(\tau(x); x, \mathbb{P}) \mid X = x) + u_0(x) \cdot \mathbb{P}(V(x, Y) = q(\tau(x); x, \mathbb{P}) \mid X = x), \end{aligned}$$

which leads to  $\gamma_0(x) \leq u_0(x)$ . On the other hand, by the definition of quantiles, we have  $\mathbb{P}(V(x, Y) \leq q(\tau(x); x, \mathbb{P}) \mid X = x) \geq \tau(x) = (u_0(x) - 1)/(u_0(x) - \ell_0(x))$ . Hence

$$\begin{aligned} & \ell_0(x) \cdot \mathbb{P}(V(x, Y) < q(\tau(x); x, \mathbb{P}) \mid X = x) + \ell_0(x) \cdot \mathbb{P}(V(x, Y) = q(\tau(x); x, \mathbb{P}) \mid X = x) \\ &= \ell_0(x) \cdot \mathbb{P}(V(x, Y) \leq q(\tau(x); x, \mathbb{P}) \mid X = x) \\ &\leq 1 - u_0(x) + u_0(x) \cdot \mathbb{P}(V(x, Y) \leq q(\tau(x); x, \mathbb{P}) \mid X = x) \\ &= 1 - u_0(x) \cdot \mathbb{P}(V(x, Y) > q(\tau(x); x, \mathbb{P}) \mid X = x), \end{aligned}$$

120 which leads to  $\gamma_0(x) \geq \ell_0(x)$ . Therefore, we conclude the proof of Proposition 3.4.  $\square$

## 121 I. Proofs of Proposition 4.1.

*Proof of Proposition 4.1.* Recall that  $\Gamma^*$  is the smallest sensitivity parameter such that  $\mathbb{P}^{\text{sup}} \in \mathcal{P}(\Gamma^*)$ . By the definition of  $\mathcal{R}$  in [B] and the nested property of the prediction sets, we know  $\mathcal{R} = [1, \hat{\Gamma}]$  if  $C \cap \hat{C}(X_{n+1}, 1) = \emptyset$ , where

$$\hat{\Gamma} = \sup\{\Gamma \geq 1 : C \cap \hat{C}(X_{n+1}, \Gamma) = \emptyset\},$$

and  $\mathcal{R} = \emptyset$  otherwise. Also, by the nested nature of  $H_0(\Gamma)$ , we know  $\mathcal{H}_0 = [\Gamma^*, \infty)$  when  $Y_{n+1}(1) - Y_{n+1}(0) \in C$  and  $\mathcal{H}_0 = \emptyset$  otherwise. Hence

$$\begin{aligned} \text{mErr} &= \mathbb{P}(\mathcal{R} \cap \mathcal{H}_0 \neq \emptyset) \\ &= \mathbb{P}\left(\{Y_{n+1}(1) - Y_{n+1}(0) \in C\} \cap \{\exists \Gamma \geq \Gamma^*, C \cap \hat{C}(X_{n+1}, \Gamma) = \emptyset\}\right) \\ &\leq \mathbb{P}\left(\{Y_{n+1}(1) - Y_{n+1}(0) \in C\} \cap \{C \cap \hat{C}(X_{n+1}, \Gamma^*) = \emptyset\}\right) \\ &\leq \mathbb{P}(Y_{n+1}(1) - Y_{n+1}(0) \notin \hat{C}(X_{n+1}, \Gamma^*)). \end{aligned}$$

Following exactly the same arguments, we have

$$\begin{aligned} \text{dErr} &= \mathbb{P}(\mathcal{R} \cap \mathcal{H}_0 \neq \emptyset \mid \mathcal{D}_{\text{calib}}) \\ &= \mathbb{P}\left(\{Y_{n+1}(1) - Y_{n+1}(0) \in C\} \cap \{\exists \Gamma \geq \Gamma^*, C \cap \hat{C}(X_{n+1}, \Gamma) = \emptyset\} \mid \mathcal{D}_{\text{calib}}\right) \\ &\leq \mathbb{P}\left(\{Y_{n+1}(1) - Y_{n+1}(0) \in C\} \cap \{C \cap \hat{C}(X_{n+1}, \Gamma^*) = \emptyset\} \mid \mathcal{D}_{\text{calib}}\right) \\ &\leq \mathbb{P}(Y_{n+1}(1) - Y_{n+1}(0) \notin \hat{C}(X_{n+1}, \Gamma^*) \mid \mathcal{D}_{\text{calib}}), \end{aligned}$$

122 completing the proof of Proposition 4.1.  $\square$

## 123 4. Additional simulation results

124 **A. Additional results for Section C.** In this part, we provide additional simulation results on the counterfactual pre-  
125 diction task in Section C.

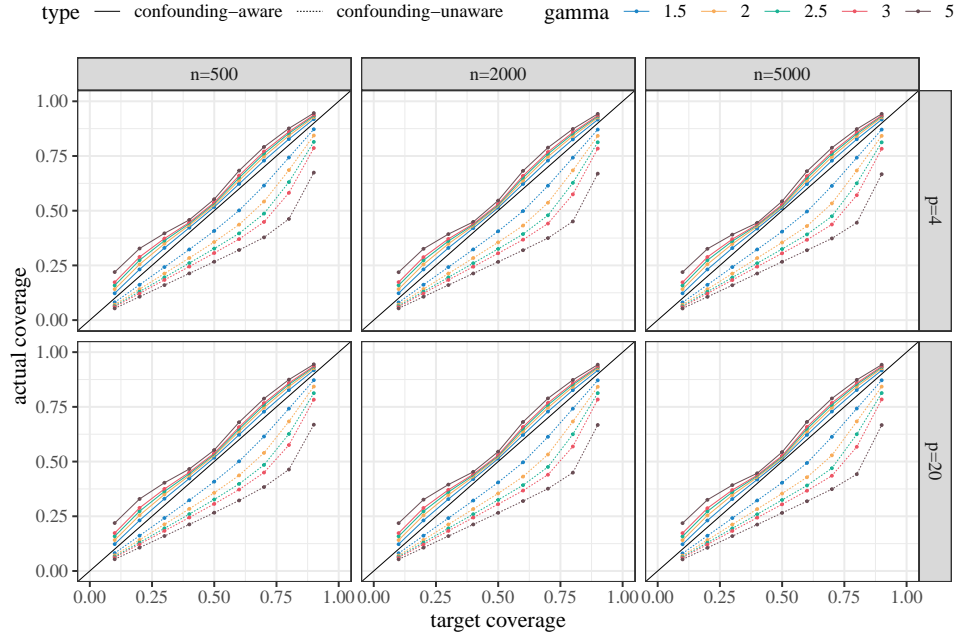


Fig. S1. Empirical (average) coverage of Algorithm 1 when  $\ell(\cdot)$  and  $u(\cdot)$  are known. The details are otherwise the same as in Figure 3.

126 **B. Additional simulation results for Section D.** In this part, we provide additional simulation results on the counter-  
 127 factual prediction task in Section D.

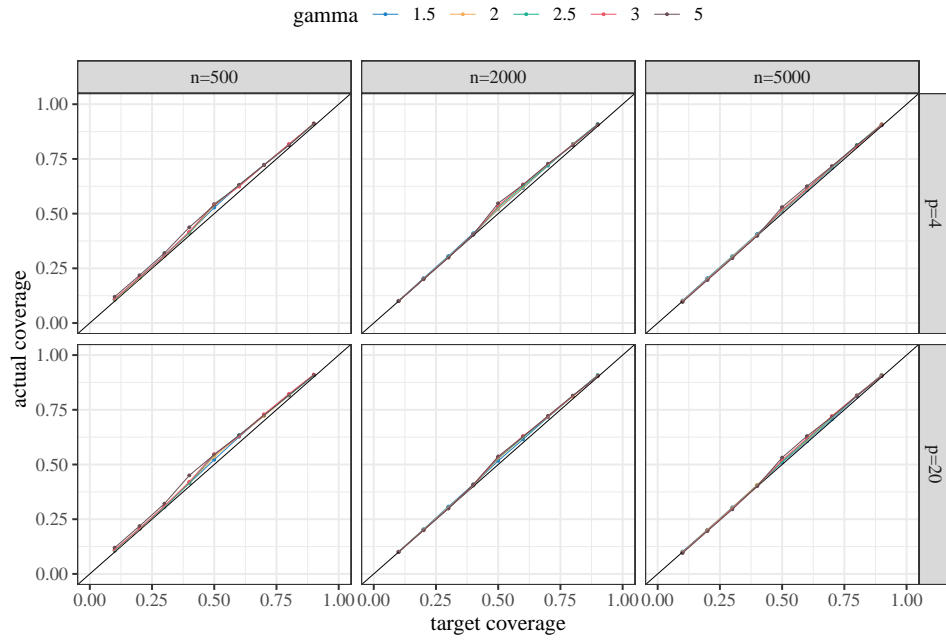


Fig. S2. 0.05-th quantile of empirical coverage on test samples in Algorithm 2 when  $\ell(\cdot)$ ,  $u(\cdot)$  are known. The details are otherwise the same as in Figure 5.

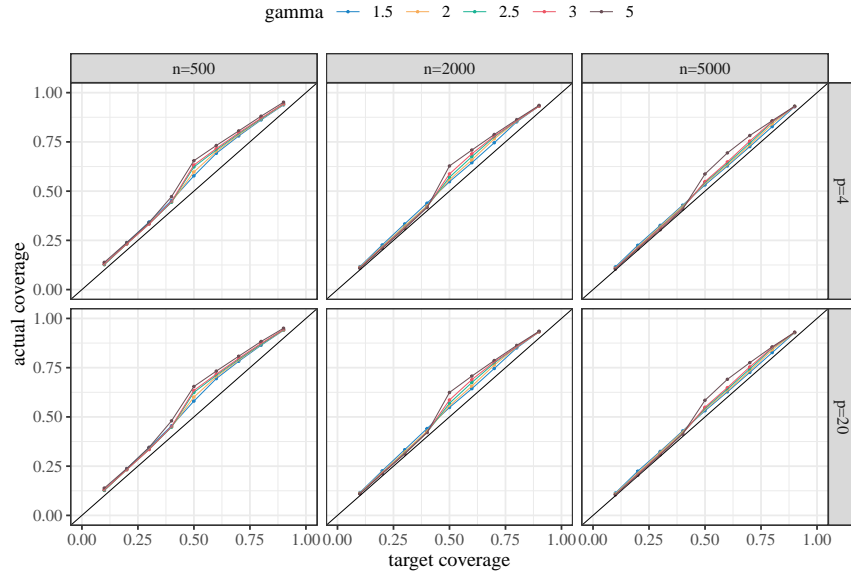


Fig. S3. Empirical (average) coverage of Algorithm 2 when  $\ell(\cdot)$  and  $u(\cdot)$  are known. The details are otherwise the same as in Figure 5.

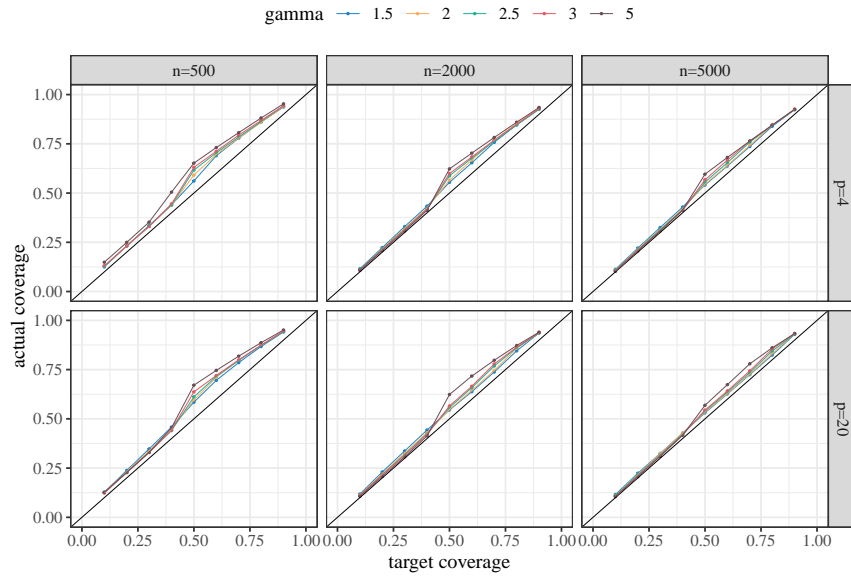


Fig. S4. Empirical (average) coverage of Algorithm 2 when  $\hat{\ell}(\cdot)$  and  $\hat{u}(\cdot)$  are estimated. The details are otherwise the same as in Figure 5.

128 **C. Additional results for Section E.** In this part, we collect the results in Section E of procedures with known bound  
 129 functions.

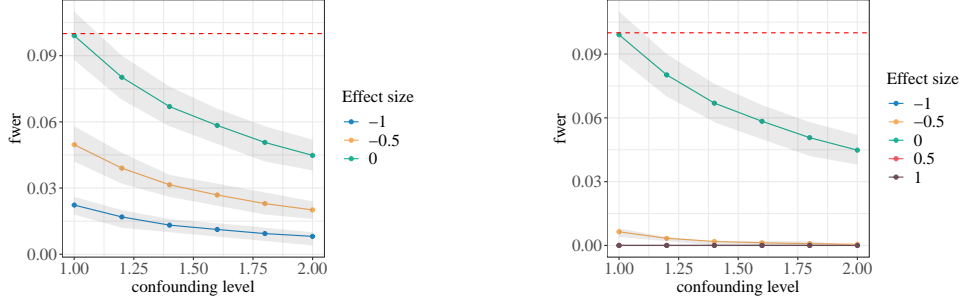


Fig. S5. Empirical FWER for fixed ITE (left) and random ITE (right) with Algorithm 1. The details are otherwise the same as in Figure S9.

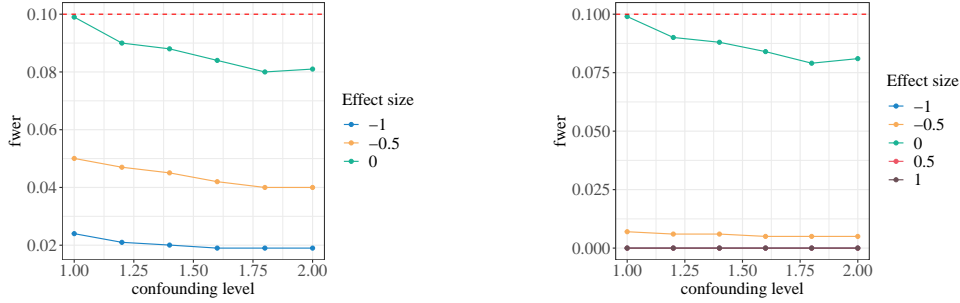


Fig. S6. Empirical FWER for fixed ITE (left) and random ITE (right) with Algorithm 2. The details are otherwise the same as in Figure S9.

130 **D. Counterfactual prediction on a semi-real dataset.** We consider an observational study dataset from (4), based on  
 131 which we generate confounded synthetic potential outcomes, and conduct the ATE-type prediction of  $Y(1)$ .

Randomly splitting the original dataset into two folds of sizes  $|\mathcal{Z}_1| = 2078$  and  $|\mathcal{Z}_2| = 8313$  respectively, we use the samples in  $\mathcal{Z}_1$  to fit a regression function  $\hat{\mu}_0(x)$  for  $\mathbb{E}[Y(0) | X = x, T = 0]$  and a propensity score function  $\hat{e}(x)$  for  $\mathbb{P}(T = 1 | X = x)$ . We sample  $n = 10000$  covariates  $\{X_i\}_{1 \leq i \leq n}$  from  $\mathcal{Z}_2$  with replacement. Then i.i.d. counterfactuals are generated via

$$Y_i(1) = \hat{\mu}_0(X_i) + \tau(X_i) + U_i, \quad Y_i(0) = \hat{\mu}_0(X_i) - U_i, \quad U_i \sim N(0, 0.2^2), \quad i = 1, \dots, n,$$

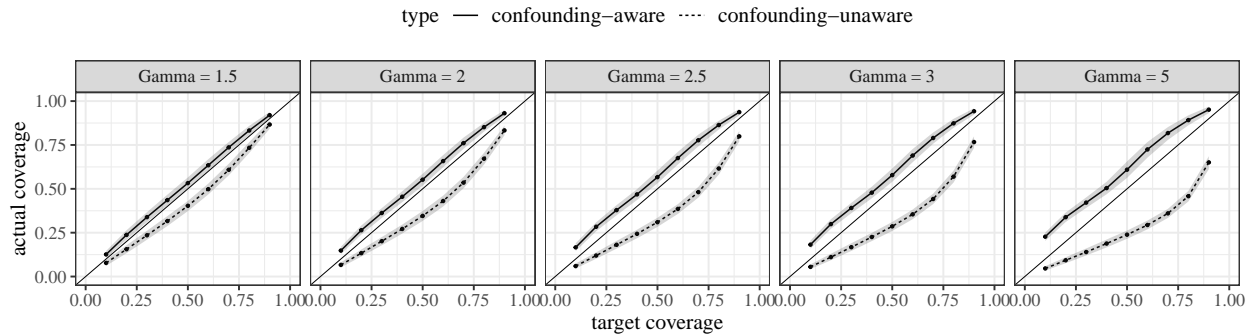
132 where the conditional treatment effect function  $\tau(x)$  is specified the same way as in equation (1) of (4). The propensity  
 133 scores are specified as  $e(X_i) := \hat{e}(X_i)$ . For each confounding level  $\Gamma \in \{1.5, 2, 2.5, 3, 5\}$ , both  $e(X_i, U_i)$  and  $T_i$  are  
 134 generated the same way as in Section C.

135 We then conduct counterfactual inference on the synthetic dataset  $\mathcal{D}_{\text{obs}} = \{Y_i, X_i, T_i\}_{1 \leq i \leq n}$ . We randomly split  
 136  $\mathcal{D}_{\text{obs}}$  into three folds with 1 : 2 : 1 sizes. The treated samples in the first fold are used as  $\mathcal{D}_{\text{train}}$ . The treated samples  
 137 in the second fold are used as  $\mathcal{D}_{\text{calib}}$ , so that  $\mathcal{D} = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{calib}}$ . All samples in the third fold (where we have ground  
 138 truth even for those in the control group) are used as test samples. The process is repeated  $N = 1000$  times, where  
 139 there are approximately  $|\mathcal{D}_{\text{calib}}| = 1900$  calibration samples fed into the procedures and 2500 test samples.

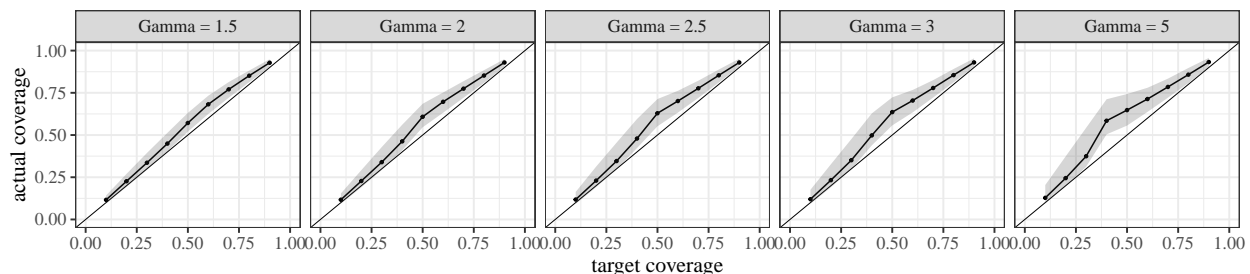
140 Figure S7 summarizes the empirical coverage using Algorithm 1 and estimated  $\hat{\ell}(\cdot)$ ,  $\hat{u}(\cdot)$ . The output of Algorithm 1  
 141 always achieves valid coverage. The solid lines are close to the 45°-line, showing the tightness of our procedure in  
 142 this setting.

143 The empirical coverage of Algorithm 2 with estimated  $\hat{\ell}(\cdot)$ ,  $\hat{u}(\cdot)$  is summarized in Figure S8, which validate the  
 144 PAC-type guarantee (referring to the lower boundary of the shaded area in each plot, which is the 0.05-th quantile of  
 145 empirical coverage). Due to the conservativeness of the confidence lower bound constructed by the WSR inequality,  
 146 the average coverage of Algorithm 2 is a bit higher than the target for targets around 0.5. However, the 0.05-th  
 147 quantile (lower boundary of the shaded area) is still very close to the target.

148 It is also worth pointing out that the shaded bands indicate the quantiles of the empirical coverage on test samples,  
 149 which could be understood as an estimate of  $\hat{c}(\mathcal{D}) := \mathbb{P}(Y_{n+1} \in \hat{C}(X_{n+1}) | \mathcal{D})$ . Although PAC-type guarantee is not  
 150 theoretically provided by Algorithm 1, most of the time,  $\hat{c}(\mathcal{D})$  is above the target in this example (referring to the  
 151 lower boundary of the shaded area, which is the 0.05-th quantile of  $\hat{c}(\mathcal{D})$ ). The widths of the 0.05-th and the 0.95-th



**Fig. S7.** Empirical coverage on the test sample. Each plot corresponds to a confounding level  $\Gamma$ . The points are average empirical coverage. The shaded bands corresponds to the 0.05-th and 0.95-th quantiles of coverage on test samples. The solid lines correspond to Algorithm 1. The dashed lines assume no confounding and are shown for comparison. In this case, counterfactual prediction intervals are invalid.



**Fig. S8.** Empirical coverage of Algorithm 2. Each plot corresponds to a confounding level  $\Gamma$ . The shaded bands corresponds to the 0.05-th and 0.95-th quantiles of coverage on test samples.

152 quantiles also provide empirical evidence that  $\hat{c}(\mathcal{D})$  from Algorithm 1 might be less variable than Algorithm 2. The  
 153 theoretical analysis of this phenomenon might deserve future investigation.

### E. Simulations for sensitivity analysis.

$$H_0(\Gamma): Y(1) - Y(0) \leq 0 \text{ and } \mathbb{P}^{\text{sup}} \in \mathcal{P}(\Gamma).$$

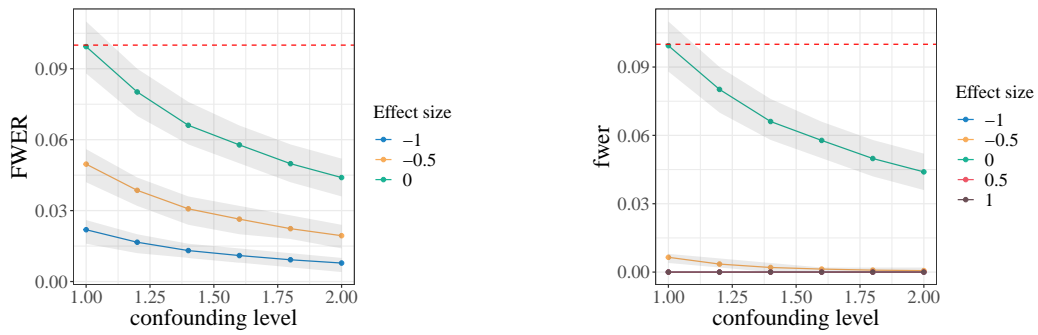
154 The test sample is from  $\mathbb{P}_{X, Y(0), Y(1) | T=1}$ , for which we observe  $(X_{n+1}, Y_{n+1}(1))$  and would like to predict  $Y_{n+1}(0)$ .  
 155 We fix  $n_{\text{train}} = n_{\text{calib}} = 2000$  and  $p = 4$ . The covariates  $X$ , unobserved confounders  $U$  and counterfactual  $Y(0)$   
 156 (instead of  $Y(1)$ ) are generated in the same way as in Section C. The treatment mechanism  $e(x, u)$  is also the same as  
 157 in [12] with confounding level  $\Gamma \in \{1.2, 1.4, \dots, 2\}$ . The training data are  $(X_i, Y_i(0))$  for those  $T_i = 0$ . We generate  
 158  $Y(1)$  in two ways: 1)  $Y(1) - Y(0) \equiv a$  (fixed ITE) and 2)  $Y(1) - Y(0) = a \cdot U$  (random ITE). Here  $a$  ranges in  
 159  $\{-1, -0.5, 0, 0.5, 1\}$ .

Fixing level  $\alpha = 0.1$  and  $\delta = 0.05$  (for Algorithm 2), for each fixed  $\Gamma \geq 1$ , we construct a one-sided ATT-type prediction interval for  $Y(0)$ , which takes the form  $\hat{C}(X_{n+1}, \Gamma) = (-\infty, \hat{Y}(X_{n+1}, \Gamma)]$ . The prediction interval is obtained via the non-conformity score

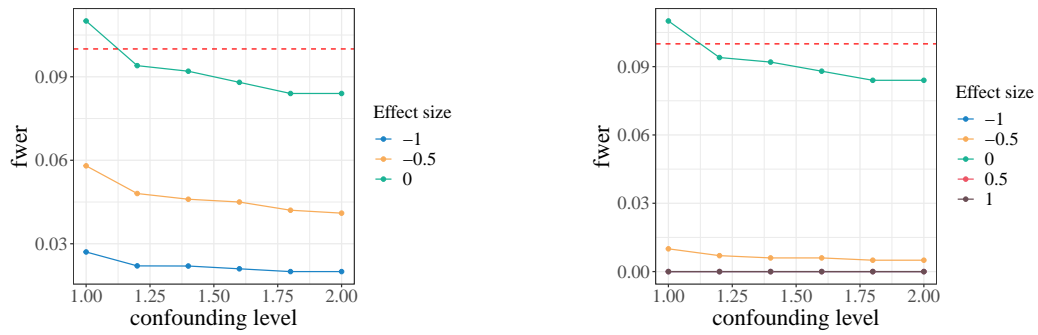
$$V(x, y) = y - \hat{q}(x, 1 - \alpha).$$

160 We reject no hypotheses if  $Y_{n+1}(1) \leq \hat{Y}(X_{n+1}, 1)$ ; otherwise, we reject all  $H_0(\Gamma)$  such that  $Y_{n+1}(1) > \hat{Y}(X_{n+1}, \Gamma)$ ,  
 161 hence the rejection set is  $\mathcal{R} = [1, \hat{\Gamma}]$  for some  $\hat{\Gamma} > 1$ , which we define as the  $\Gamma$ -value.

162 **E.1. FWER control.** We evaluate the empirical FWER, which is the proportion of making a false rejection among  
 163  $\{H_0(\Gamma)\}_{\Gamma \geq 1}$ , averaged over all test samples in all  $N = 1000$  independent runs. The results with estimated  $\hat{\ell}(\cdot)$  and  
 164  $\hat{u}(\cdot)$  are presented in Figures S9 and S10, showing control of the FWER even when the likelihood ratio bounds are  
 165 estimated; we omit the case of fixed ITE at  $a > 0$  since the hypotheses  $H_0(\Gamma)$  are always false. The results with the  
 166 ground truth  $\ell(\cdot)$  and  $u(\cdot)$  are in SI Appendix C.

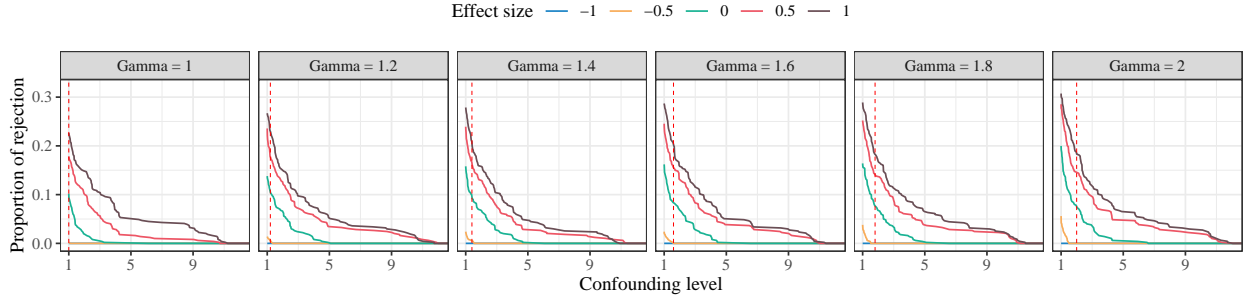


**Fig. S9.** Empirical FWER of Algorithm 1. The effect size  $a$  ranges in  $\{-1, -0.5, 0, 0.5, 1\}$ . The solid lines are averaged over  $N = 1000$  runs, while the 0.25-th and 0.75-th quantiles form the shaded area.



**Fig. S10.** 0.05-th quantile of empirical FWER using Algorithm 2 with estimated  $\hat{\ell}(\cdot)$  and  $\hat{u}(\cdot)$ , with the effect size  $a$  ranging in  $\{-1, -0.5, 0, 0.5, 1\}$  for fixed ITE (left) and random ITE (right).

167 **E.2.  $\Gamma$ -values.** We plot the estimated survival function  $\hat{\Gamma}$  defined as  $S(\Gamma) = \mathbb{P}(\hat{\Gamma} > \Gamma)$ , which characterizes the propor-  
 168 tion of test units that are identified as positive ITE with each confounding level  $\Gamma$ . Figure S11 presents the results  
 169 from one run of the procedure with Algorithm 1, where we focus on the random ITE:  $Y(1) - Y(0) = a \cdot U$ .

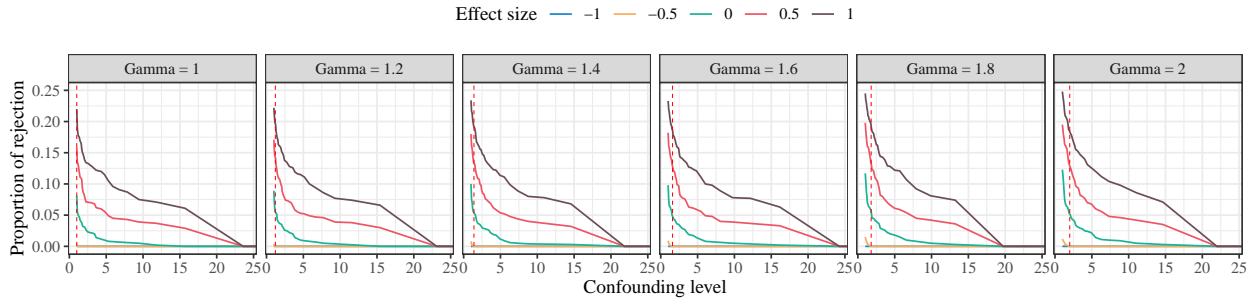


**Fig. S11.** Empirical evaluation of  $S(\Gamma)$  reported by (one run of) the sensitivity analysis procedure with Algorithm 1. The red dashed vertical lines are the true confounding levels.

170 To see how to interpret these plots, let us consider the example where  $a = 1$  and the true confounding level is 1.6.  
 171 We can see that around 10% of the samples have a  $\Gamma$ -value greater than 2.5, and 5% have a  $\Gamma$ -value greater than 5,  
 172 showing strong evidence for positive ITEs. Note also that the ITE is always zero when  $a = 0$ , and the  $\Gamma$ -value should  
 173 be a 90% lower confidence bound for the true confounding level. In Figure S11, we indeed observe that for around  
 174 90% of the samples, the  $\Gamma$ -value is below the true confounding level (see the green curves).

175 With random ITEs, both the magnitude of actual ITE and the gap between observed outcomes in treated and  
 176 control groups increase with the effects size  $a$ . Within each subplot, the reported  $\Gamma$ -values become larger as the effect  
 177 size increases. Thresholding at the true confounding level  $\Gamma$ , we see that larger magnitude of true effects also makes  
 178 it easier to detect positive ITEs at the correct confounding level  $\Gamma$ .

179 The results from one run of Algorithm 2 are in Figure S12. The patterns are similar to Figure S11 in general,  
 180 except that it is sometimes sharper than Algorithm 1 and provides slightly stronger evidence against unmeasured  
 181 confounding.



**Fig. S12.** Empirical evaluation of  $S(\Gamma)$  reported by (one run of) the sensitivity analysis procedure with Algorithm 2. The red dashed vertical lines are the true confounding levels.

182 In Figure S12, some test units have large  $\Gamma$ -values especially when the effect size is positive. Such strong evidence  
 183 would happen if there is a large gap between the observed  $Y(1)$  and the typical behavior of  $Y(0)$  predicted with  
 184 the training data—thus, the only way our procedures can output a prediction interval that overlaps with  $Y(1)$  is  
 185  $\hat{C}(X, \Gamma) = \{V(X, y) \leq \hat{v}\} = \mathbb{R}$ , where  $\hat{v} = \infty$ . In that case, there is a certain value  $\Gamma_\infty(X)$  such that  $\hat{v} = \infty$  once  
 186  $\Gamma \geq \Gamma_\infty(X)$ . Since Algorithm 2 outputs a same  $\hat{v}$  for all test samples, the value  $\Gamma_\infty$  is the same for all such extreme  
 187 individuals, leading to steep tails in Figure S12. In contrast, Algorithm 1 returns different  $\Gamma_\infty(X)$  for different  
 188 individuals, and produces a smoother curve (Figure S11).

189 **E.3. False discovery proportions.** We also track the empirical false discovery proportion  $\text{FDP}(\Gamma) = \frac{|\{j \in \mathcal{D}_{\text{test}} : \hat{\Gamma} > \Gamma, Y_j(1) \leq Y_j(0)\}|}{|\{j \in \mathcal{D}_{\text{test}} : \hat{\Gamma} > \Gamma\}|}$   
 190 for random ITE with  $a \neq 0$ , which is the proportion of false rejections among test units that are rejected at confound-  
 191 ing level  $\Gamma$ . With both the two procedures, we find that  $\text{FDP}(\Gamma)$  is always zero—the units that survive certain levels

192 of adjustment for confounding all have positive ITE. It could be explained by the conservativeness of the procedure:  
193 to survive the adjustment, the observed  $Y(1)$  needs to be larger than the whole  $1 - \alpha$  prediction interval for  $Y(0)$ .  
194 Therefore, a unit that survives the adjustment is much more likely to have a positive ITE, leading to vanishing FDPs.  
195

## 196 **References**

- 197 1. I Waudby-Smith, A Ramdas, Estimating means of bounded random variables by betting (2021).
- 198 2. Y Romano, E Patterson, E Candes, Conformalized quantile regression. *Adv. Neural Inf. Process. Syst.* **32**, 3543–  
199 3553 (2019).
- 200 3. RJ Tibshirani, RF Barber, EJ Candès, A Ramdas, Conformal prediction under covariate shift in *Advances in*  
201 *Neural Information Processing Systems 32*. (2019).
- 202 4. C Carvalho, A Feller, J Murray, S Woody, D Yeager, Assessing treatment effect variation in observational studies:  
203 Results from a data challenge. *Obs. Stud.* **5**, 21–35 (2019).