

Supplementary Material for “Integrative R -learner of heterogeneous treatment effects combining experimental and observational studies”

Outline of the appendices

- Appendix A: proof of Theorem 1.
- Appendix B: proof of Proposition 2.
- Appendix C: pseudocode of iterative learning algorithm used in the real data experiment.

Appendix A. Proof of Theorem 1

To prove Theorem 1, we show a useful lemma first.

Lemma 3 Under Assumptions 3-4, $\hat{L}_N(\beta) = L_N(\beta) + O_P(a_N^2)$.

Proof For the simplicity of the notations, we denote

$$\begin{aligned} A_{\mu,i} &= \mu(X_i, S_i) - \hat{\mu}^{-k(i)}(X_i, S_i), \quad B_{\mu,i} = Y_i - \mu(X_i, S_i), \\ A_{e,i} &= e(X_i, S_i) - \hat{e}^{-k(i)}(X_i, S_i), \quad B_{e,i} = A_i - e(X_i, S_i). \end{aligned}$$

By algebra, we have

$$\begin{aligned} \hat{L}_N(\beta) &= \mathbb{P}_N \{ B_{\mu,i} + A_{\mu,i} - g(X_i, S_i; \beta) B_{e,i} - g(X_i, S_i; \beta) A_{e,i} \}^2 \\ &= \mathbb{P}_N \{ B_{\mu,i} - g(X_i, S_i; \beta) B_{e,i} \}^2 + \mathbb{P}_N \{ A_{\mu,i} - g(X_i, S_i; \beta) A_{e,i} \}^2 + \\ &\quad 2\mathbb{P}_N \{ B_{\mu,i} - g(X_i, S_i; \beta) B_{e,i} \} \mathbb{P}_N \{ A_{\mu,i} - g(X_i, S_i; \beta) A_{e,i} \} \\ &= L_N(\beta) + \mathbb{P}_N A_{\mu,i}^2 + \mathbb{P}_N A_{e,i}^2 g^2(X_i, S_i) - 2\mathbb{P}_N A_{\mu,i} A_{e,i} g(X_i, S_i; \beta) + \\ &\quad 2\mathbb{P}_N B_{\mu,i} A_{\mu,i} - 2\mathbb{P}_N B_{\mu,i} A_{e,i} g(X_i, S_i; \beta) - 2\mathbb{P}_N B_{e,i} A_{\mu,i} g(X_i, S_i; \beta) + 2\mathbb{P}_N B_{e,i} A_{e,i} g^2(X_i, S_i). \end{aligned}$$

By Markov's inequality and Assumption 4, we have the second term above $\mathbb{P}_N A_{\mu,i}^2$ is $O_P(a_N^2)$. Plus Assumption 3, we have the third term $\mathbb{P}_N A_{e,i}^2$ is $O_P(a_N^2)$. As for the fourth term,

$$\mathbb{P}_N A_{\mu,i} A_{e,i} g(X_i, S_i; \beta) \leq C_1 \mathbb{P}_N A_{\mu,i} A_{e,i} \leq C_1 \sqrt{\mathbb{P}_N A_{\mu,i}^2 \mathbb{P}_N A_{e,i}^2} = O_P(a_N^2),$$

for some positive constant C_1 , and the last inequality is by Cauchy-Schwarz inequality.

Next, to deal with the last four terms, we tackle $\mathbb{P}_N B_{\mu,i} A_{\mu,i}$ first. Let

$$B_{\mu\mu}^k = \frac{\sum_{i:k(i)=k} B_{\mu,i} A_{\mu,i}}{|\{i : k(i) = k\}|},$$

and note that $|\mathbb{P}_N B_{\mu,i} A_{e,i}| \leq \sum_{k=1}^K |B_{\mu\mu}^k|$, where K is finite, thus it is suffice to show $B_{\mu\mu}^k = O_P(a_N^2)$. Let $\mathcal{I}^{-k} = \{X_i, A_i, Y_i, S_i : k(i) \neq k\}$. Then we have

$$\begin{aligned} \mathbb{E}(B_{\mu\mu}^k) &= \mathbb{E}(B_{\mu,i} A_{e,i}) = \mathbb{E}\{\mathbb{E}(B_{\mu,i} A_{e,i} \mid \mathcal{I}^{-k}, X_i, S_i)\} \\ &= \mathbb{E}\{A_{e,i} \mathbb{E}(B_{\mu,i} \mid \mathcal{I}^{-k}, X_i, S_i)\} = 0. \end{aligned}$$

Then we have its variance

$$\begin{aligned}\text{var}(B_{\mu\mu}^k) &= \mathbb{E}\{(B_{\mu\mu}^k)^2\} = \frac{\mathbb{E}\{\sum_{i:k(i)=k} B_{\mu,i}^2 A_{\mu,i}^2 + \sum_{i \neq j:k(i)=k, k(j)=k} B_{\mu,i} B_{\mu,j} A_{\mu,i} A_{\mu,j}\}}{|\{i : k(i) = k\}|^2} \\ &= \frac{\mathbb{E}(B_{\mu,i}^2 A_{\mu,i}^2)}{|\{i : k(i) = k\}|} + \frac{\sum_{i \neq j:k(i)=k, k(j)=k} \mathbb{E}(B_{\mu,i} B_{\mu,j} A_{\mu,i} A_{\mu,j})}{|\{i : k(i) = k\}|^2}.\end{aligned}$$

By Assumption 3-4 we have $\mathbb{E}(B_{\mu,i}^2 A_{e,i}^2) = \mathbb{E}\{\mathbb{E}(B_{\mu,i}^2 A_{e,i}^2 \mid \mathcal{I}^{-k}, X_i, S_i)\} = \mathbb{E}\{A_{e,i}^2 \mathbb{E}(B_{\mu,i}^2 \mid \mathcal{I}^{-k}, X_i, S_i)\} \leq C_2 \mathbb{E}A_{e,i}^2 = O(a_N^2)$, for some positive constant C_2 . As for the interaction terms, we have $\mathbb{E}(B_{\mu,i} B_{\mu,j} A_{\mu,i} A_{\mu,j}) = \mathbb{E}\{A_{\mu,i} A_{\mu,j} \mathbb{E}(B_{\mu,i} B_{\mu,j} \mid \mathcal{I}^{-k}, X_i, S_i)\} = \mathbb{E}\{A_{\mu,i} A_{\mu,j} \mathbb{E}(B_{\mu,j}) \mathbb{E}(B_{\mu,i} \mid \mathcal{I}^{-k}, X_i, S_i)\} = 0$. The second last equality is implied by $B_{\mu,i}$ and $B_{\mu,j}$ are independent for $i \neq j$, and the last equality comes from the definition of $B_{\mu,i}$. Therefore, we have $\text{var}(B_{\mu\mu}^k) = (K/N)O(a_N^2) = O(a_N^2/N)$, which is negligible with a faster diminishing rate than $O(a_N^2)$. Then, by Chebyshev' inequality, we have $B_{\mu\mu}^k = O_P(a_N^2/N)$, i.e., $\mathbb{P}_N B_{\mu,i} A_{\mu,i} = O_P(a_N^2/N)$. Similarly, under the Assumption 3 that $g(X_i, S_i; \beta)$ is uniformly bounded, we can get the same results for the left three terms equal to $O_P(a_N^2/N)$. Therefore, $\widehat{L}_N(\beta) - L_N(\beta)$ is dominated by the $O_P(a_N^2)$ -term $\mathbb{P}_N A_{\mu,i}^2 + \mathbb{P}_N A_{e,i}^2 g^2(X_i, S_i) - 2\mathbb{P}_N A_{\mu,i} A_{e,i} g(X_i, S_i; \beta)$, which finally leads to the conclusion in the lemma $\widehat{L}_N(\beta) - L_N(\beta) = O_P(a_N^2)$. \blacksquare

Next, we are showing the proof of Theorem 1 with the help of Lemma 3.

Proof of Theorem 1:

First, we let $\mathbb{G}_N = \sqrt{N}(\mathbb{P}_N - \mathbb{P})$, $p^\tau = \{A - e(X, A)\}p_\tau(X)$, and $p^c = \{A - e(X, A)\}(1 - S)p_c(X)$. Recall $p_i^\tau = \{(p^\tau)^\top, (p^c)^\top\}$, and then

$$\Gamma = \mathbb{P}(p_i p_i^\top) = \begin{Bmatrix} p^\tau (p^\tau)^\top & p^\tau (p^c)^\top \\ p^c (p^\tau)^\top & p^c (p^c)^\top \end{Bmatrix} \stackrel{\text{denoted as}}{=} \begin{pmatrix} \Gamma_{\tau\tau} & \Gamma_{\tau c} \\ \Gamma_{c\tau} & \Gamma_{cc} \end{pmatrix}.$$

By the definition of $\hat{\beta}$ and Lemma 3, we have

$$\hat{\beta} = \underset{b}{\text{argmin}} \mathbb{P}_N \{Y_i - \mu(X_i, S_i) - p_i^\top b\}^2 + O_P(a_N^2).$$

By the Pointwise Linearization of the series method (see Lemma 4.1 in Belloni et al., 2015), under Assumptions 5-7 and the fact $\mathbb{E}(\epsilon \mid X, A, S) = 0$ which is shown under Assumptions 1-2 in Section 3, we have for any $\alpha := (\alpha_\tau^\top, \alpha_c^\top)^\top \in \mathbb{R}^{2D}$,

$$\sqrt{N}\alpha^\top(\hat{\beta} - \beta) = \alpha^\top \begin{pmatrix} \Gamma_{\tau\tau} & \Gamma_{\tau c} \\ \Gamma_{c\tau} & \Gamma_{cc} \end{pmatrix}^{-1} \mathbb{G}_N \begin{pmatrix} p^\tau \epsilon \\ p^c \epsilon \end{pmatrix} + o_P(1) + O_P(\sqrt{N}a_N^2). \quad (12)$$

Let

$$\Sigma = \begin{pmatrix} \Sigma_{\tau\tau} & \Sigma_{\tau c} \\ \Sigma_{c\tau} & \Sigma_{cc} \end{pmatrix} := \begin{pmatrix} \Gamma_{\tau\tau} & \Gamma_{\tau c} \\ \Gamma_{c\tau} & \Gamma_{cc} \end{pmatrix}^{-1}.$$

Since under Assumption 4, $a_N = O(N^{-r})$, $r > 1/4$, thus $O_P(\sqrt{N}a_N^2)$ is negligible compared to $o_P(1)$. Therefore, we have

$$\sqrt{N}\alpha_\tau^\top(\hat{\beta}_\tau - \beta_\tau) = \alpha_\tau^\top \mathbb{G}_N (\Sigma_{\tau\tau} p^\tau \epsilon + \Sigma_{\tau c} p^c \epsilon) + o_P(1). \quad (13)$$

Under Assumptions 5-7 and the fact $\mathbb{E}(\epsilon \mid X, A, S) = 0$, by the Pointwise Normality of the series method (see Theorem 4.2 in Belloni et al., 2015), we have

$$\sqrt{N} \frac{\alpha_\tau^\top (\hat{\beta}_\tau - \beta_\tau)}{\|\alpha_\tau^\top \Omega^{1/2}\|} \xrightarrow{d} \mathcal{N}(0, 1) + o_P(1),$$

where $\Omega = \mathbb{E} \{ (\Sigma_{\tau\tau} p^\tau \epsilon + \Sigma_{\tau c} p^c \epsilon) (\Sigma_{\tau\tau} p^\tau \epsilon + \Sigma_{\tau c} p^c \epsilon)^\top \}$. Then, we take $\alpha_\tau = p^\tau$, for any $x \in \mathcal{X}$,

$$\sqrt{N} \frac{(p^\tau)^\top (\hat{\beta}_\tau - \beta_\tau)}{\|(p^\tau)^\top \Omega^{1/2}\|} \xrightarrow{d} \mathcal{N}(0, 1) + o_P(1).$$

Under Assumption 6(d), the approximation error is negligible relative to the estimation error, then

$$\sqrt{N} \frac{\hat{\tau}(x) - \tau(x)}{\|(p^\tau)^\top \Omega^{1/2}\|} \xrightarrow{d} \mathcal{N}(0, 1) + o_P(1),$$

which immediately arrives at the first part conclusion in Theorem 1, $\hat{\tau}(x) - \tau(x) = O(N^{-1/2})$, for any $x \in \mathcal{X}$. Besides, it also gives the asymptotic variance of $\hat{\tau}(x)$,

$$\mathbb{V}\{\hat{\tau}(x)\} = N^{-1} (p^\tau)^\top \Omega p^\tau. \quad (14)$$

Expanding Ω , under Assumption 7, we have

$$\begin{aligned} \Omega &= \Sigma_{\tau\tau} \mathbb{E} \{ p^\tau (p^\tau)^\top \epsilon^2 \} \Sigma_{\tau\tau} + \Sigma_{\tau c} \mathbb{E} \{ p^\tau (p^c)^\top \epsilon^2 \} \Sigma_{\tau c} + \\ &\quad \Sigma_{\tau c} \mathbb{E} \{ p^c (p^c)^\top \epsilon^2 \} \Sigma_{\tau c} + \Sigma_{\tau\tau} \mathbb{E} \{ p^\tau (p^c)^\top \epsilon^2 \} \Sigma_{\tau c} \\ &= (\Sigma_{\tau\tau} \Gamma_{\tau\tau} + \Sigma_{\tau c} \Gamma_{\tau c}^\top) \Sigma_{\tau\tau} \sigma^2 + (\Sigma_{\tau c} \Gamma_{cc} + \Sigma_{\tau\tau} \Gamma_{\tau c}) \Sigma_{\tau c}^\top \sigma^2 \\ &= \Sigma_{\tau\tau} \sigma^2 + \mathbf{0}_{D \times D} \\ &= (\Gamma_{\tau\tau} - \Gamma_{\tau c} \Gamma_{cc}^{-1} \Gamma_{\tau c}^\top)^{-1} \sigma^2 \\ &= [\mathbb{E}\{S p^\tau (p^\tau)^\top\} + \mathbb{E}\{(1-S) p^\tau (p^\tau)^\top\} - \Gamma_{\tau c} \Gamma_{cc}^{-1} \Gamma_{\tau c}^\top]^{-1} \sigma^2. \end{aligned} \quad (15)$$

Next, we aim to obtain the asymptotic variance of the HTE estimator $\hat{\tau}_{\text{rct}} = (p^\tau)^\top \hat{\beta}_{\text{rct}}$ with only the RCT, where $\hat{\beta}_{\text{rct}} = \text{argmin}_{b \in \mathbb{R}^D} \mathbb{P}_N S_i [Y_i - \hat{\mu}^{-k(i)}(X_i, S_i) - \{A_i - \hat{e}^{-k(i)}(X_i, S_i)\} p_i^\top b]^2$. Similarly, we can replace the estimated nuisance functions with the true ones based on Lemma 3, and following the same strategy as the integrative R -learner above to obtain the asymptotic variance, $\mathbb{V}(\hat{\tau}_{\text{rct}}) = N^{-1} (p^\tau)^\top \Omega_{\text{rct}} p^\tau$, where

$$\Omega_{\text{rct}} = [\mathbb{E}\{S p^\tau (p^\tau)^\top\}]^{-1} \sigma^2. \quad (16)$$

By Hölder's inequality, $\mathbb{E}\{(1-S) p^\tau (p^\tau)^\top\} - \Gamma_{\tau c} \Gamma_{cc}^{-1} \Gamma_{\tau c}^\top$ is non-negative definitive; i.e., for any $v \in \mathbb{R}^D$,

$$v^\top (\Omega^{-1} - \Omega_{\text{rct}}^{-1}) v \geq 0, \quad (17)$$

where the inequality becomes an equality if and only if $p_\tau(X) = M p_c(X)$ for some constant matrix M . From (14), we have

$$p^\tau \mathbb{V}\{\hat{\tau}(x)\} (p^\tau)^\top = N^{-1} p^\tau (p^\tau)^\top \Omega p^\tau (p^\tau)^\top$$

$$\begin{aligned}
&\implies \mathbf{I}_D \mathbb{V}\{\hat{\tau}(x)\} = N^{-1} p^\tau (p^\tau)^\top \Omega \\
&\implies \mathbf{I}_D \mathbb{V}^{-1}\{\hat{\tau}(x)\} = N \Omega^{-1} \{p^\tau (p^\tau)^\top\}^{-1} \\
&\implies (p^\tau)^\top \{p^\tau (p^\tau)^\top\}^{-1} \mathbb{V}^{-1}\{\hat{\tau}(x)\} p^\tau = N (p^\tau)^\top \{p^\tau (p^\tau)^\top\}^{-1} \Omega^{-1} \{p^\tau (p^\tau)^\top\}^{-1} p^\tau \\
&\stackrel{\text{by (17)}}{\implies} (p^\tau)^\top \{p^\tau (p^\tau)^\top\}^{-1} [\mathbb{V}^{-1}\{\hat{\tau}(x)\} - \mathbb{V}^{-1}\{\hat{\tau}_{\text{rct}}(x)\}] p^\tau \geq 0 \\
&\quad (\text{Multiply a positive number } (p^\tau)^\top p^\tau \text{ on the both sides and get the below formula}) \\
&\implies (p^\tau)^\top p^\tau \{p^\tau (p^\tau)^\top\}^{-1} p^\tau [\mathbb{V}^{-1}\{\hat{\tau}(x)\} - \mathbb{V}^{-1}\{\hat{\tau}_{\text{rct}}(x)\}] \geq 0 \\
&\implies (p^\tau)^\top p^\tau [\mathbb{V}^{-1}\{\hat{\tau}(x)\} - \mathbb{V}^{-1}\{\hat{\tau}_{\text{rct}}(x)\}] \geq 0 \\
&\implies \mathbb{V}^{-1}\{\hat{\tau}(x)\} - \mathbb{V}^{-1}\{\hat{\tau}_{\text{rct}}(x)\} \geq 0 \\
&\implies \mathbb{V}\{\hat{\tau}(x)\} \leq \mathbb{V}\{\hat{\tau}_{\text{rct}}(x)\},
\end{aligned}$$

with the equality holding when $p_\tau(X) = M p_c(X)$ for some constant matrix M . ■

Appendix B. Proof of Proposition 2

Proof The proof of (10) is mainly based on the inverse probability weights (IPW) component, thus we tackle it first.

a) IPW-adjusted outcomes: We have

$$\begin{aligned}
&\mathbb{E} \left\{ \frac{AY}{e(X, S)} \mid X, S = 0 \right\} \\
&= \frac{1}{e(X, 0)} [\mathbb{E}\{AY \mid X, A = 1, S = 0\} e(X, 0) + \mathbb{E}\{AY \mid X, A = 0, S = 0\} \{1 - e(X, 0)\}] \\
&= \mathbb{E}\{Y \mid X, A = 1, S = 0\}.
\end{aligned}$$

Similarly, we have

$$\mathbb{E} \left\{ \frac{(1 - A)Y}{1 - e(X, S)} \mid X, S = 0 \right\} = \mathbb{E}\{Y \mid X, A = 0, S = 0\}.$$

b) Augmented IPW-adjusted outcomes: Then we have

$$\begin{aligned}
&\mathbb{E} \left[\frac{A\{Y - Q_S(X, 1)\}}{e(X, S)} + Q_S(X, 1) \mid X, S = 0 \right] \\
&= \mathbb{E} \left\{ \frac{AY}{e(X, S)} \mid X, S = 0 \right\} + Q_0(X, 1) \mathbb{E} \left[\left\{ 1 - \frac{A}{e(X, S)} \right\} \mid X, S = 0 \right] \\
&= \mathbb{E}(Y \mid X, A = 1, S = 0).
\end{aligned}$$

Similarly, we have

$$\mathbb{E} \left[\frac{(1 - A)\{Y - Q_S(X, 0)\}}{1 - e(X, S)} + Q_S(X, 0) \mid X, S = 0 \right] = \mathbb{E}(Y \mid X, A = 0, S = 0).$$

Finally, by taking the difference of the above two formulas and based on the definition of $c(X)$ in (1), we arrive at the conclusion $\mathbb{E}(\tilde{Y} \mid X, S = 0) = \tau(X) + c(X)$ ■

Appendix C. Iterative learning for the integrative R -learner in the real data experiment

Algorithm 2 Iterative learning for the integrative R -learner

Data: The RCT and the OS data $\{(X_i, A_i, Y_i, S_i)\}_{i \in \mathcal{I}_n \cup \mathcal{I}_m}$; the number of iterations B

Result: $\tau(\cdot)$

Initialize confounding function $c(X_i)$.

Estimate the nuisance functions with cross-fitting based on Xgboost, denoted as $\hat{\mu}(X_i, S_i)$ and $\hat{e}(X_i, S_i)$, resulting in $\hat{e}_{Y_i} = Y_i - \hat{\mu}(X_i, S_i)$ and $\hat{e}_{A_i} = A_i - \hat{e}(X_i, S_i)$.

for $b = 1, \dots, B$ **do**

Calculate the pseudo outcomes $\tilde{Y}_i \leftarrow \hat{e}_{Y_i} / \hat{e}_{A_i} - (1 - S_i)c(X_i), i \in \mathcal{I}_n \cup \mathcal{I}_m$;

Obtain $\tau(\cdot)$ by fitting \tilde{Y}_i on X_i using weighted random forest with weights equal to $\hat{e}_{A_i}^2, i \in \mathcal{I}_n \cup \mathcal{I}_m$;

Update the pseudo outcomes $\tilde{Y}_i \leftarrow \hat{e}_{Y_i} / \hat{e}_{A_i} - \tau(X_i), i \in \mathcal{I}_m$;

Update $c(X_i)$ by fitting \tilde{Y}_i on X_i using weighted random forest with weights equal to $\hat{e}_{A_i}^2, i \in \mathcal{I}_m$.

end
