

**Improve Efficiency of Doubly Robust Estimator
when Propensity Score is Misspecified**

Liangbo Lyu

University of Michigan.

Molei Liu

Columbia University Mailman School of Public Health

Supplementary Material

S1 Numerical implementation of Algorithm 1

For the RWLS problem in Step 3 of Algorithm 1, inspired by Lemma 1, its solution can be calculated based on:

$$\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\Sigma}}^{-1} \hat{\mathbf{a}} (\hat{\mathbf{a}}^\top \hat{\boldsymbol{\Sigma}}^{-1} \hat{\mathbf{a}})^{-1} \hat{\mathbf{a}}^\top \hat{\boldsymbol{\Sigma}}^{-1} \hat{\mathbf{b}} - \hat{\boldsymbol{\Sigma}}^{-1} \hat{\mathbf{b}},$$

where $\hat{\mathbf{a}} = \hat{\mathbb{E}}_{\mathcal{S}} \hat{\boldsymbol{\Psi}} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \hat{\boldsymbol{\alpha}})$, $\hat{\mathbf{b}} = \hat{\mathbb{E}}_{\mathcal{S}} \hat{\boldsymbol{\Psi}} \exp(\mathbf{X}^\top \hat{\boldsymbol{\gamma}}) v_{\hat{\boldsymbol{\theta}}}(\mathbf{X}) + \hat{\mathbb{E}}_{\mathcal{S}} \hat{\boldsymbol{\Psi}} \mathbf{X}^\top v_{\hat{\boldsymbol{\theta}}}(\mathbf{X}) \hat{\mathbf{L}}$, and $\hat{\boldsymbol{\Sigma}} = \hat{\mathbb{E}}_{\mathcal{S}} \hat{\boldsymbol{\Psi}} \hat{\boldsymbol{\Psi}}^\top v_{\hat{\boldsymbol{\theta}}}(\mathbf{X})$. The other steps in Algorithm 1 can be implemented directly following their descriptions.

The two authors have equal contribution.

S2 Methodological extension

S2.1 Construction with alternative nuisance estimators

As emphasized in Remark 2, there are many potential choices of the estimating equations for the nuisance parameters $\hat{\gamma}$ and $\hat{\alpha}$, e.g., the maximum likelihood estimation. Under the correct specified OR model, the variance of $\hat{\gamma}$ (determined by its estimating equations) has no influence on the asymptotic properties of the PAD estimator. For $\hat{\alpha}$ obtained using alternative estimating equations to (2.1), we need to modify the form of $\hat{V}_\mu(\boldsymbol{\beta})$. Meanwhile, we can show that this modification will not change the double robustness and variance-reduction properties of our estimator.

Let $\hat{\mathbb{E}}_{\mathcal{S}}\mathbf{S}'(\boldsymbol{\alpha}) = \mathbf{0}$ be the estimating equation of the outcome model $g(\mathbf{X}^\top\boldsymbol{\alpha})$ with the score function $\mathbf{S}'(\boldsymbol{\alpha})$ different from $\mathbf{S}(\boldsymbol{\alpha}) = \mathbf{X}\{Y - g(\mathbf{X}^\top\boldsymbol{\alpha})\}$ used in the main text. Suppose $\hat{\alpha}$ converges to $\bar{\alpha}$, the objective variance function then needs to be changed to

$$\hat{V}_\mu(\boldsymbol{\beta}) = \hat{\mathbb{E}}_{\mathcal{S}}\{\exp(\mathbf{X}^\top\hat{\gamma}) + \hat{\Psi}^\top\boldsymbol{\beta}\}^2 v_{\hat{\theta}}(\mathbf{X}) + 2\hat{\mathbf{L}}^\top\hat{\mathbb{E}}_{\mathcal{S}}\{\exp(\mathbf{X}^\top\hat{\gamma}) + \hat{\Psi}^\top\boldsymbol{\beta}\}\{Y - g(\mathbf{X}^\top\hat{\alpha})\}\mathbf{S}'(\hat{\alpha})$$

correspondingly with

$$\hat{\mathbf{L}} = \mathbb{E}_{\mathcal{S}}\left\{\left.\frac{\partial\mathbf{S}'(\boldsymbol{\alpha})}{\partial\boldsymbol{\alpha}}\right|_{\hat{\alpha}}\right\}^{-1}\left\{\hat{\mathbb{E}}_{\mathcal{S}}\mathbf{X}\dot{g}(\mathbf{X}^\top\hat{\alpha})\exp(\mathbf{X}^\top\hat{\gamma}) - \hat{\mathbb{E}}_{\mathcal{T}}\mathbf{X}\dot{g}(\mathbf{X}^\top\hat{\alpha})\right\},$$

and we correspondingly have

$$V_\mu(\boldsymbol{\beta}) = \mathbb{E}_S\{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^\top \boldsymbol{\beta}\}^2 v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) + 2\mathbf{L}^\top \mathbb{E}_S\{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^\top \boldsymbol{\beta}\} \{Y - g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}})\} \mathbf{S}'(\bar{\boldsymbol{\alpha}})$$

the limiting function of $\widehat{V}_\mu(\boldsymbol{\beta})$ with

$$\mathbf{L} = \mathbb{E}_S \left\{ \frac{\partial \mathbf{S}(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}} \bigg|_{\bar{\boldsymbol{\alpha}}} \right\}^{-1} \{ \mathbb{E}_S \mathbf{X} \dot{g}(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}) \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) - \mathbb{E}_T \mathbf{X} \dot{g}(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}) \}$$

When the PS model is correct, we again have $\mathbf{L} = \mathbf{0}$ so the same construction of $\boldsymbol{\Psi}$ ensures double robustness. For correct OR model but wrong PS, the modified $V_\mu(\boldsymbol{\beta})$ can lead to variance reduction following a very similar discussion to Section 2.3.

S2.2 Dual construction to augment OR

In analogy to our PAD estimator, to improve the efficiency our the DR estimator under the correct PS and wrong OR models, we propose the Outcome regression Augmented Doubly robust (OAD) estimator in the following algorithm.

Algorithm 1 Outcome regression Augmented Doubly robust (OAD) estimation

[Step 1] Solve the estimating equations in (2.1) to obtain $\hat{\gamma}$ and $\hat{\alpha}$, and obtain the conditional variance estimator as $\hat{\theta}$.

[Step 2] Let $\Phi = \phi(\mathbf{X})$ with function $\phi(\cdot)$, $\tilde{g}(\mathbf{X}^\top \hat{\alpha}) = g(\mathbf{X}^\top \hat{\alpha}) - \hat{\mathbb{E}}_{\mathcal{T}} g(\mathbf{X}^\top \hat{\alpha})$ and

$$\hat{\Psi} = \Phi - \frac{\hat{\mathbb{E}}_{\mathcal{T}} \Phi \tilde{g}(\mathbf{X}^\top \hat{\alpha})}{\hat{\mathbb{E}}_{\mathcal{T}} \tilde{g}^2(\mathbf{X}^\top \hat{\alpha})} \tilde{g}(\mathbf{X}^\top \hat{\alpha}).$$

[Step 3] Solve the restricted weighted least square (RWLS) problem:

$$\hat{\beta}_{\text{OAD}} = \operatorname{argmin}_{\beta} \hat{V}_{\mu, \text{OAD}}(\beta), \quad \text{s.t.} \quad \hat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \hat{\Psi}^\top \beta \exp(\mathbf{X}^\top \hat{\gamma}) = \mathbf{0}, \quad (\text{S2.1})$$

where

$$\begin{aligned} \hat{V}_{\mu, \text{OAD}}(\beta) &= n^{-1} \widehat{\text{Var}}_{\mathcal{S}}[\{Y - g(\mathbf{X}^\top \hat{\alpha}) - \hat{\Psi}^\top \hat{\beta}\} \exp(\mathbf{X}^\top \hat{\gamma})] + N^{-1} \widehat{\text{Var}}_{\mathcal{T}}\{g(\mathbf{X}^\top \hat{\alpha}) + \hat{\Psi}^\top \hat{\beta}\} \\ &\quad + 2 \widehat{\mathbf{L}}^* [N^{-1} \widehat{\text{Cov}}_{\mathcal{T}}(\mathbf{X}, \hat{\Psi}^\top \hat{\beta}) + n^{-1} \widehat{\text{Cov}}_{\mathcal{S}}\{\mathbf{X} \exp(\mathbf{X}^\top \hat{\gamma}), \hat{\Psi}^\top \hat{\beta} \exp(\mathbf{X}^\top \hat{\gamma})\}], \end{aligned} \quad (\text{S2.2})$$

and $\widehat{\mathbf{L}}^* = \{\hat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \exp(\mathbf{X}^\top \hat{\gamma}) \mathbf{X}^\top\}^{-1} \hat{\mathbb{E}}_{\mathcal{S}} \{Y - g(\mathbf{X}^\top \hat{\alpha})\} \exp(\mathbf{X}^\top \hat{\gamma}) \mathbf{X}$.

[Step 4] Obtain the OAD estimator:

$$\hat{\mu}_{\text{OAD}} = \hat{\mathbb{E}}_{\mathcal{S}} \{Y - g(\mathbf{X}^\top \hat{\alpha}) - \hat{\Psi}^\top \hat{\beta}_{\text{OAD}}\} \exp(\mathbf{X}^\top \hat{\gamma}) + \hat{\mathbb{E}}_{\mathcal{T}} \{g(\mathbf{X}^\top \hat{\alpha}) + \hat{\Psi}^\top \hat{\beta}_{\text{OAD}}\}.$$

To demonstrate how Algorithm 1 works, we define that

$$\begin{aligned} \tilde{\mu}_{\text{OAD}} &= \hat{\mathbb{E}}_{\mathcal{S}} \{Y - g(\mathbf{X}^\top \bar{\alpha}) - \Psi^\top \bar{\beta}_{\text{OAD}}\} \exp(\mathbf{X}^\top \bar{\gamma}) + \hat{\mathbb{E}}_{\mathcal{T}} \{g(\mathbf{X}^\top \bar{\alpha}) + \Psi^\top \bar{\beta}_{\text{OAD}}\} \\ &\quad + \mathbb{E}_{\mathcal{S}} \{Y - g(\mathbf{X}^\top \bar{\alpha})\} \exp(\mathbf{X}^\top \bar{\gamma}) \mathbf{X}^\top \{\mathbb{E}_{\mathcal{S}} \mathbf{X} \exp(\mathbf{X}^\top \bar{\gamma}) \mathbf{X}^\top\}^{-1} \{\hat{\mathbb{E}}_{\mathcal{T}} \mathbf{X} - \hat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \exp(\mathbf{X}^\top \bar{\gamma})\}. \end{aligned}$$

Then similar to our previous analysis, when the PS model is correct, $\hat{\mu}_{\text{OAD}}$

is asymptotically equivalent to $\tilde{\mu}_{\text{OAD}}$, and

$$V_{\mu, \text{OAD}}(\boldsymbol{\beta}) = n^{-1} \text{Var}_S[\{Y - g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}) - \boldsymbol{\Psi}^\top \bar{\boldsymbol{\beta}}_{\text{OAD}}\} \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}})] + N^{-1} \text{Var}_T\{g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}) + \boldsymbol{\Psi}^\top \bar{\boldsymbol{\beta}}_{\text{OAD}}\} \\ + 2\mathbf{L}^{*\top} [N^{-1} \text{Cov}_T(\mathbf{X}, \boldsymbol{\Psi}^\top \bar{\boldsymbol{\beta}}_{\text{OAD}}) + n^{-1} \text{Cov}_S\{\mathbf{X} \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}), \boldsymbol{\Psi}^\top \bar{\boldsymbol{\beta}}_{\text{OAD}} \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}})\}],$$

is the limiting function of $\widehat{V}_{\mu, \text{OAD}}(\boldsymbol{\beta})$ specified in Algorithm 1, where

$$\mathbf{L}^* = \{\mathbb{E}_S \mathbf{X} \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) \mathbf{X}^\top\}^{-1} \mathbb{E}_S \{Y - g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}})\} \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) \mathbf{X}.$$

Similar to the PAD method, when $\boldsymbol{\beta} = \mathbf{0}$, $V_{\mu, \text{OAD}}(\boldsymbol{\beta})$ reduces to the asymptotic variance of the standard DR estimator (with a constant difference invariant with $\boldsymbol{\beta}$). Thus, $\widehat{\mu}_{\text{OAD}}$ has a smaller variance than the standard DR estimator when the PS model is correct and the OR model is wrong, under which we typically have $\bar{\boldsymbol{\beta}} \neq \mathbf{0}$. On the other hand, when the OR is correctly specified, we have $\bar{\boldsymbol{\alpha}} = \boldsymbol{\alpha}_0$, $\mathbf{L}^* = \mathbf{0}$, and thus

$$\left. \frac{\partial V_{\mu, \text{OAD}}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta}=\mathbf{0}} = \text{Cov}_T(g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}), \boldsymbol{\Psi}).$$

By definition of $\boldsymbol{\Psi}$, we have $\text{Cov}_T(g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}), \boldsymbol{\Psi}) = \mathbf{0}$, implying that $\bar{\boldsymbol{\beta}}_{\text{OAD}} = \mathbf{0}$ is the (unique) population-level solution of (S2.1). Hence, similar to the analysis in Section 2.2, $\widehat{\mu}_{\text{OAD}}$ preserved the same DR property as $\widehat{\mu}_{\text{DR}}$, i.e., being root- n consistent whenever the PS or the OR model is correctly specified.

S2.3 Ensemble strategy

One might be interested in deriving an estimator more efficient than the standard DR estimator whenever the PS or OR model is misspecified (and the other one is correct). This can be achieved through a natural extension that convexly combines the PAD and OAD estimators with the optimal weights minimizing the asymptotic variance. In specific, we define the ensemble estimator

$$\widehat{\mu}_{\text{ENS}}(c) = c\widehat{\mu}_{\text{PAD}} + (1 - c)\widehat{\mu}_{\text{OAD}},$$

where $c \in [0, 1]$ is some ensemble weight to be specified. Since both $\widehat{\mu}_{\text{PAD}}$ and $\widehat{\mu}_{\text{OAD}}$ preserve the DR property, $\widehat{\mu}_{\text{ENS}}(c)$ is also an DR estimator of μ_0 with any c . Based on this, we propose to choose the c from $[0, 1]$ that minimizes

$$\widehat{\text{aVar}} \left\{ \sqrt{n}(\widehat{\mu}_{\text{ENS}}(c) - \mu_0) \right\} = \widehat{\text{aVar}} \left\{ \sqrt{n}(c\widehat{\mu}_{\text{PAD}} + (1 - c)\widehat{\mu}_{\text{OAD}} - \mu_0) \right\}. \quad (\text{S2.3})$$

$\widehat{\text{aVar}} \left\{ \sqrt{n}(\widehat{\mu}_{\text{ENS}}(c) - \mu_0) \right\}$ is a quadratic function of c and can be extracted using the asymptotic covariance matrix of $(\widehat{\mu}_{\text{PAD}}, \widehat{\mu}_{\text{OAD}})$ that can be directly estimated from the asymptotic expansions of $\widehat{\mu}_{\text{PAD}}$ and $\widehat{\mu}_{\text{OAD}}$. Let \widehat{c} be the minimizer of (S2.3), and the ensemble estimator $\widehat{\mu}_{\text{ENS}} = \widehat{\mu}_{\text{ENS}}(\widehat{c})$. It is not hard to see that the asymptotic variance of $\widehat{\mu}_{\text{ENS}}$ is always smaller or equal to those of $\widehat{\mu}_{\text{ENS}}(0) = \widehat{\mu}_{\text{OAD}}$ and $\widehat{\mu}_{\text{ENS}}(1) = \widehat{\mu}_{\text{PAD}}$ as \widehat{c} minimizes (S2.3). Thus, $\widehat{\mu}_{\text{ENS}}$ can realize bias-reduction compared to the standard DR estimator $\widehat{\mu}_{\text{DR}}$, whichever nuisance model is correct and the other one is wrong. This ensemble strategy can also be used to combine our estimator with existing

intrinsic efficient estimators like Shu and Tan (2018).

S2.4 Extension to ATE estimation

Problem setup

In this section, we extend our method to provide an PAD estimator for the average treatment effects (ATE) that attains similar DR and variance reduction properties as that for ATT. Suppose there are independent samples of (Δ_i, Y_i, X_i) for $i = 1, 2, \dots, n$, where Δ_i is the binary treatment variable. Let $Y(\delta)$ denote the counterfactual outcome under the treatment status $\Delta = \delta \in \{0, 1\}$. We are interested in estimating the ATE parameter $\mu_{\text{ATE}} = \mathbb{E}\{Y(1)\} - \mathbb{E}\{Y(0)\}$, under the standard assumptions (Hernán and Robins, 2010, e.g.) including (i) exchangeability that $\{Y(0), Y(1)\} \perp\!\!\!\perp \Delta \mid \mathbf{X}$; (ii) positivity that $\pi(\mathbf{X}) = \Pr(\Delta = 1 \mid \mathbf{X})$ stays away from 0 and 1 for all \mathbf{X} ; and (iii) consistency that $Y = Y(\delta)$ for all subjects with $\Delta = \delta$.

Following the common setup in Bang and Robins (2005), we specify the PS model $\pi(\mathbf{X}) = \text{expit}(\mathbf{X}^\top \boldsymbol{\gamma})$ with $\text{expit}(a) = e^a / (1 + e^a)$, and the OR models $\mathbb{E}(Y \mid \mathbf{X}, \Delta = \delta) = g(\mathbf{X}^\top \boldsymbol{\alpha}_\delta)$ for $\delta = 0, 1$. Similar to (2.1), we

derive their estimators $\hat{\boldsymbol{\gamma}}$ and $\hat{\boldsymbol{\alpha}}_\delta$ by solving

$$\hat{\mathbb{E}}_0 \mathbf{X} \exp(\mathbf{X}^\top \boldsymbol{\gamma}) = \hat{\mathbb{E}}_1 \mathbf{X}, \quad \hat{\mathbf{S}}_\delta(\boldsymbol{\alpha}_\delta) = \hat{\mathbb{E}}_\delta \mathbf{X} \{Y - g(\mathbf{X}^\top \boldsymbol{\alpha}_\delta)\} = \mathbf{0}, \quad \delta = 0, 1, \quad (\text{S2.4})$$

where $\hat{\mathbb{E}}_0$ and $\hat{\mathbb{E}}_1$ represent the empirical mean operator on treatment ($\Delta = 1$) and control ($\Delta = 0$) samples respectively. Then by Bang and Robins (2005), the standard DR estimator of ATE is constructed as

$$\begin{aligned} \hat{\mu}_{\text{ATE}} = & \hat{\mathbb{E}} [\Delta \{1 + \exp(\mathbf{X}^\top \hat{\boldsymbol{\gamma}})\} \{Y - g(\mathbf{X}^\top \hat{\boldsymbol{\alpha}}_1)\} + g(\mathbf{X}^\top \hat{\boldsymbol{\alpha}}_1)] \\ & - \hat{\mathbb{E}} [(1 - \Delta) \{1 + \exp(-\mathbf{X}^\top \hat{\boldsymbol{\gamma}})\} \{Y - g(\mathbf{X}^\top \hat{\boldsymbol{\alpha}}_0)\} + g(\mathbf{X}^\top \hat{\boldsymbol{\alpha}}_0)], \end{aligned}$$

where $\hat{\mathbb{E}}$ is the empirical mean operator on samples $i = 1, 2, \dots, n$.

PAD estimation of ATE

We propose the PAD estimation of ATE in Algorithm 2 and provide some heuristical justification on it.

Algorithm 2 The PAD estimation of ATE

[Step 1] Solve the estimating equations in (S2.4) to obtain $\hat{\gamma}$, $\hat{\alpha}_0$ and $\hat{\alpha}_1$, and derive the conditional variance estimators for Y on treatment and control samples as $\hat{\theta}_1$ and $\hat{\theta}_0$.

[Step 2] For $\delta = 0, 1$, specify $\Phi_\delta = \phi_\delta(\mathbf{X})$ of larger dimensionality than \mathbf{X} using any basis function $\phi_\delta(\cdot)$, and take $\hat{\Psi}_\delta = \Phi_\delta - \hat{\mathbb{E}}[\Phi_\delta v_{\hat{\theta}_\delta}(\mathbf{X})] / \hat{\mathbb{E}}v_{\hat{\theta}_\delta}(\mathbf{X})$.

[Step 3] Solve the restricted weighted least square (RWLS) problem:

$$\hat{\beta}_\delta = \operatorname{argmin}_{\beta_\delta} \hat{V}_{\mu_\delta}(\beta_\delta), \quad \text{s.t.} \quad \hat{\mathbb{E}}_\delta \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \hat{\alpha}_\delta) \hat{\Psi}_\delta \beta_\delta = \mathbf{0}, \quad (\text{S2.5})$$

where the objective function is defined as

$$\begin{aligned} \hat{V}_{\mu_\delta}(\beta_\delta) &= \hat{\mathbb{E}}I(\Delta = \delta) [\exp\{(2\delta - 1)\mathbf{X}^\top \hat{\gamma}\} + 1 + \hat{\Psi}_\delta^\top \beta_\delta]^2 v_{\hat{\theta}_\delta}(\mathbf{X}) \\ &\quad + 2\hat{\mathbb{L}}_\delta^\dagger \hat{\mathbb{E}}I(\Delta = \delta) \mathbf{X} [\exp\{(2\delta - 1)\mathbf{X}^\top \hat{\gamma}\} + 1 + \hat{\Psi}_\delta^\top \beta_\delta] v_{\hat{\theta}_\delta}(\mathbf{X}), \end{aligned}$$

and

$$\hat{\mathbf{L}}_\delta = \left\{ \hat{\mathbb{E}}_\delta \frac{\partial \hat{\mathbf{S}}_\delta(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}^\top} \Big|_{\boldsymbol{\alpha} = \hat{\boldsymbol{\alpha}}_\delta} \right\}^{-1} \left[\hat{\mathbb{E}}I(\Delta = \delta) \mathbf{X} \dot{g}(\mathbf{X}^\top \hat{\alpha}_\delta) [1 + \exp\{(2\delta - 1)\mathbf{X}^\top \hat{\gamma}\}] - \hat{\mathbb{E}}\mathbf{X} \dot{g}(\mathbf{X}^\top \hat{\alpha}_\delta) \right]$$

[Step 4] Obtain the PAD estimator for $\mu_\delta = \mathbb{E}\{Y(\delta)\}$ as

$$\hat{\mu}_\delta = \hat{\mathbb{E}} \left[I(\Delta = \delta) [1 + \exp\{(2\delta - 1)\mathbf{X}^\top \hat{\gamma}\} + \hat{\Psi}_\delta^\top \hat{\beta}_\delta] \{Y - g(\mathbf{X}^\top \hat{\alpha}_\delta)\} + g(\mathbf{X}^\top \hat{\alpha}_\delta) \right]$$

[Step 5] Obtain the ATE estimator through

$$\hat{\mu}_{\text{PAD-ATE}} = \hat{\mu}_1 - \hat{\mu}_0$$

Suppose that all estimators converge to their limiting values (holding under some reasonable regularity conditions). For $\delta = 0, 1$, let $\Psi_\delta = \Phi -$

$\mathbb{E}\Phi v_{\theta_\delta}(\mathbf{X})/\mathbb{E}v_{\theta_\delta}(\mathbf{X})$, $\bar{\boldsymbol{\beta}}_\delta$, and $\bar{\boldsymbol{\alpha}}_\delta$ and $\bar{\boldsymbol{\gamma}}$ represent the limit of $\widehat{\boldsymbol{\beta}}_\delta$, $\widehat{\boldsymbol{\alpha}}_\delta$, and $\widehat{\boldsymbol{\gamma}}$ respectively,

$$\mathbf{L}_\delta = \left\{ \mathbb{E}_\delta \frac{\partial \mathcal{S}_\delta(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}^\top} \Big|_{\bar{\boldsymbol{\alpha}}_\delta} \right\}^{-1} \{ \mathbb{E}I(\Delta = \delta) \mathbf{X} \dot{g}(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}_\delta) [1 + \exp\{(2\delta - 1)\mathbf{X}^\top \bar{\boldsymbol{\gamma}}\}] - \mathbb{E} \mathbf{X} \dot{g}(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}_\delta) \},$$

and

$$V_{\mu_\delta}(\boldsymbol{\beta}_\delta) = \mathbb{E}I(\Delta = \delta) [\exp\{(2\delta - 1)\mathbf{X}^\top \bar{\boldsymbol{\gamma}}\} + 1 + \boldsymbol{\Psi}_\delta^\top \boldsymbol{\beta}_\delta]^2 v_{\bar{\boldsymbol{\theta}}_\delta}(\mathbf{X}) + 2\mathbf{L}_\delta^\top \mathbb{E}I(\Delta = \delta) \mathbf{X} [\exp\{(2\delta - 1)\mathbf{X}^\top \bar{\boldsymbol{\gamma}}\} + 1 + \boldsymbol{\Psi}_\delta^\top \boldsymbol{\beta}_\delta] v_{\bar{\boldsymbol{\theta}}_\delta}(\mathbf{X})$$

as the limiting function of $\widehat{V}_{\mu_\delta}(\boldsymbol{\beta}_\delta)$. Next, we shall discuss the properties of $\widehat{\mu}_{\text{PAD-ATE}}$ in two scenarios including (i) correct PS model; and (ii) correct OR and wrong PS models. Similar to the ATT setting, one can show that $\widehat{\mu}_{\text{PAD-ATE}}$ can simultaneously attain the DR property and variance-reduction compared to the standard DR under wrong PS and correct OR.

Correct PS. When the PS model is correct, we have $\mathbf{L}_\delta = \mathbf{0}$ for both $\delta = 0, 1$ and, thus,

$$\frac{\partial V_{\mu_\delta}(\boldsymbol{\beta}_\delta)}{\partial \boldsymbol{\beta}_\delta} \Big|_{\boldsymbol{\beta}_\delta = \mathbf{0}} = 2\mathbb{E}I(\Delta = \delta) \boldsymbol{\Psi}_\delta [\exp\{(2\delta - 1)\mathbf{X}^\top \bar{\boldsymbol{\gamma}}\} + 1] v_{\bar{\boldsymbol{\theta}}_\delta}(\mathbf{X}) = 2\mathbb{E} \boldsymbol{\Psi}_\delta v_{\bar{\boldsymbol{\theta}}_\delta}(\mathbf{X}) = \mathbf{0}$$

by the definition of $\boldsymbol{\Psi}_\delta$ for $\delta = 0, 1$. This, combined with the strong convexity of $V_{\mu_\delta}(\boldsymbol{\beta}_\delta)$, implies $\bar{\boldsymbol{\beta}}_\delta = \mathbf{0}$ when the PS model is correct. Thus, $\widehat{\mu}_{\text{PAD-ATE}}$ attains the same doubly robust property as the standard $\widehat{\mu}_{\text{ATE}}$.

Wrong PS and Correct OR. When the OR model is correctly specified and the PS is misspecified, we have that $\widehat{\mu}_{\text{PAD-ATE}}$ is asymptotically equivalent to $\widetilde{\mu}_{\text{PAD-ATE}} = \widetilde{\mu}_1 - \widetilde{\mu}_0$, where

$$\begin{aligned} \widetilde{\mu}_\delta = & \widehat{\mathbb{E}}I(\Delta = \delta)[1 + \exp\{(2\delta - 1)\mathbf{X}^\top \bar{\boldsymbol{\gamma}}\} + \boldsymbol{\Psi}_\delta^\top \bar{\boldsymbol{\beta}}_\delta]\{Y - g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}_\delta)\} \\ & + \widehat{\mathbb{E}}g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}_\delta) + \mathbf{L}_\delta^\top \widehat{\mathbb{E}}I(\Delta = \delta)\mathbf{X}\{Y - g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}_\delta)\} \end{aligned}$$

for $\delta = 0, 1$. Based on this, have

$$\text{aVar}\{n^{1/2}(\widehat{\mu}_{\text{PAD-ATE}} - \mu_0)\} = V_{\mu_1}(\boldsymbol{\beta}_1) + V_{\mu_0}(\boldsymbol{\beta}_0) + C,$$

which, after dropping the invariant C , is the sum of $V_{\mu_1}(\boldsymbol{\beta}_1)$ and $V_{\mu_0}(\boldsymbol{\beta}_0)$.

When the PS model is wrong, $\partial V_{\mu_\delta}(\boldsymbol{\beta}_\delta)/\partial \boldsymbol{\beta}_\delta$ is typically not $\mathbf{0}$ at $\boldsymbol{\beta}_\delta = \mathbf{0}$ so its minimizer $\bar{\boldsymbol{\beta}}_\delta \neq \mathbf{0}$. Thus, similar to the reason in Section 2.3, we have $\text{aVar}\{n^{1/2}(\widehat{\mu}_{\text{PAD-ATE}} - \mu_0)\} \leq \text{aVar}\{n^{1/2}(\widehat{\mu}_{\text{ATE}} - \mu_0)\}$, and the strict “ $<$ ” will hold in general when PS is misspecified.

S3 Asymptotic justification

S3.1 Proof of Theorem 1

Proof. Proof of Theorem 1 (i).

When the OR is correctly specified, $\bar{\boldsymbol{\alpha}} = \boldsymbol{\alpha}_0$. Consider $\widetilde{\mu}_{\text{OR}}$ where

$$\widetilde{\mu}_{\text{OR}} = \widehat{\mathbb{E}}_S\{Y - g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}})\}\{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^\top \bar{\boldsymbol{\beta}}\} + \widehat{\mathbb{E}}_{\mathcal{T}}g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}})$$

$$+ \{\mathbb{E}_{\mathcal{S}} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}) \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) - \mathbb{E}_{\mathcal{T}} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \bar{\boldsymbol{\alpha}})\} \mathbb{E}_{\mathcal{S}} \left\{ \left. \frac{\partial \mathbf{S}(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}^\top} \right|_{\bar{\boldsymbol{\alpha}}} \right\}^{-1} \widehat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \{Y - g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}})\}.$$

It is obvious that $\mathbb{E} \tilde{\boldsymbol{\mu}}_{\text{OR}} = \mathbb{E}_{\mathcal{T}} g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}) = \boldsymbol{\mu}_0$. Hence, by using central limit theorem, we have that $\tilde{\boldsymbol{\mu}}_{\text{OR}} - \boldsymbol{\mu}_0 = O_p(n^{-1/2})$, $n^{1/2}(\tilde{\boldsymbol{\mu}}_{\text{OR}} - \boldsymbol{\mu}_0)$ weakly converges to gaussian distribution with mean $\mathbf{0}$. On the other hand, we have that

$$\begin{aligned} \widehat{\boldsymbol{\mu}}_{\text{PAD}} - \tilde{\boldsymbol{\mu}}_{\text{OR}} &= \widehat{\mathbb{E}}_{\mathcal{S}} \{Y - g(\mathbf{X}^\top \boldsymbol{\alpha}_0)\} \{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) \mathbf{X}^\top (\widehat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^\top (\widehat{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}}) + (\widehat{\boldsymbol{\Psi}} - \boldsymbol{\Psi})^\top \bar{\boldsymbol{\beta}}\} \\ &\quad - [\widehat{\mathbb{E}}_{\mathcal{S}} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \boldsymbol{\alpha}_0) \{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^\top \bar{\boldsymbol{\beta}}\} - \widehat{\mathbb{E}}_{\mathcal{T}} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \boldsymbol{\alpha}_0)] (\widehat{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}}) + o_p(n^{-1/2}) \\ &\quad - \{\mathbb{E}_{\mathcal{S}} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \boldsymbol{\alpha}_0) \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) - \mathbb{E}_{\mathcal{T}} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \boldsymbol{\alpha}_0)\} \mathbb{E}_{\mathcal{S}} \left\{ \left. \frac{\partial \mathbf{S}(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}^\top} \right|_{\boldsymbol{\alpha}_0} \right\}^{-1} \widehat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \{Y - g(\mathbf{X}^\top \boldsymbol{\alpha}_0)\}, \end{aligned}$$

by using central limit theorem, along with Lemma (1)-(4), we have that

$$\begin{aligned} &\widehat{\mathbb{E}}_{\mathcal{S}} \{Y - g(\mathbf{X}^\top \boldsymbol{\alpha}_0)\} \{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) \mathbf{X}^\top (\widehat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^\top (\widehat{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}}) + (\widehat{\boldsymbol{\Psi}} - \boldsymbol{\Psi})^\top \bar{\boldsymbol{\beta}}\} \\ &= [\widehat{\mathbb{E}}_{\mathcal{S}} \{Y - g(\mathbf{X}^\top \boldsymbol{\alpha}_0)\} \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) \mathbf{X}^\top] (\widehat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}}) + [\widehat{\mathbb{E}}_{\mathcal{S}} \{Y - g(\mathbf{X}^\top \boldsymbol{\alpha}_0)\} \boldsymbol{\Psi}^\top] (\widehat{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}}) \\ &\quad + [\widehat{\mathbb{E}}_{\mathcal{S}} \{Y - g(\mathbf{X}^\top \boldsymbol{\alpha}_0)\} \bar{\boldsymbol{\beta}}^\top] (\widehat{\boldsymbol{\Psi}} - \boldsymbol{\Psi}) = O_p(n^{-1/2}) o_p(1) + O_p(n^{-1/2}) o_p(1) + O_p(n^{-1/2}) o_p(1) \\ &= o_p(n^{-1/2}). \end{aligned}$$

On the other hand,

$$\begin{aligned} &- [\widehat{\mathbb{E}}_{\mathcal{S}} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \boldsymbol{\alpha}_0) \{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^\top \bar{\boldsymbol{\beta}}\} - \widehat{\mathbb{E}}_{\mathcal{T}} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \boldsymbol{\alpha}_0)] (\widehat{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}}) \\ &= [\widehat{\mathbb{E}}_{\mathcal{S}} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \boldsymbol{\alpha}_0) \{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^\top \bar{\boldsymbol{\beta}}\} - \widehat{\mathbb{E}}_{\mathcal{T}} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \boldsymbol{\alpha}_0)] \widehat{\mathbb{E}}_{\mathcal{S}} \left\{ \left. \frac{\partial \mathbf{S}(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}^\top} \right|_{\bar{\boldsymbol{\alpha}}} \right\}^{-1} \widehat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \{Y - g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}})\} \\ &= \{\mathbb{E}_{\mathcal{S}} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \boldsymbol{\alpha}_0) \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) - \mathbb{E}_{\mathcal{T}} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \boldsymbol{\alpha}_0) + O_p(n^{-1/2})\} \\ &\quad * \left[\mathbb{E}_{\mathcal{S}} \left\{ \left. \frac{\partial \mathbf{S}(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}^\top} \right|_{\boldsymbol{\alpha}_0} \right\}^{-1} + O_p(n^{-1/2}) \right] \widehat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \{Y - g(\mathbf{X}^\top \boldsymbol{\alpha}_0)\}. \end{aligned} \tag{S3.6}$$

Hence, we have that

$$\begin{aligned}
& - [\widehat{\mathbb{E}}_{\mathcal{S}} \mathbf{X}^{\top} \dot{g}(\mathbf{X}^{\top} \boldsymbol{\alpha}_0) \{\exp(\mathbf{X}^{\top} \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^{\top} \bar{\boldsymbol{\beta}}\} - \widehat{\mathbb{E}}_{\mathcal{T}} \mathbf{X}^{\top} \dot{g}(\mathbf{X}^{\top} \boldsymbol{\alpha}_0)] (\hat{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}}) \\
& - \{\mathbb{E}_{\mathcal{S}} \mathbf{X}^{\top} \dot{g}(\mathbf{X}^{\top} \boldsymbol{\alpha}_0) \exp(\mathbf{X}^{\top} \bar{\boldsymbol{\gamma}}) - \mathbb{E}_{\mathcal{T}} \mathbf{X}^{\top} \dot{g}(\mathbf{X}^{\top} \boldsymbol{\alpha}_0)\} \mathbb{E}_{\mathcal{S}} \left\{ \frac{\partial \mathbf{S}(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}^{\top}} \bigg|_{\boldsymbol{\alpha}_0} \right\}^{-1} \widehat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \{Y - g(\mathbf{X}^{\top} \boldsymbol{\alpha}_0)\} \\
& = \widehat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \{Y - g(\mathbf{X}^{\top} \boldsymbol{\alpha}_0)\} O_p(n^{-1/2}) = o_p(n^{-1/2}).
\end{aligned}$$

Thus, from previous results, we have that $\widehat{\mu}_{\text{PAD}} - \widetilde{\mu}_{\text{OR}} = o_p(n^{-1/2})$. Together with Slutsky theorem, we further have that $\widehat{\mu}_{\text{PAD}} - \mu_0 = O_p(n^{-1/2})$ and $n^{1/2}(\widehat{\mu}_{\text{PAD}} - \mu_0)$ weakly converges to gaussian distribution with mean $\mathbf{0}$.

When the PS is correctly specified, $\bar{\boldsymbol{\gamma}} = \boldsymbol{\gamma}_0$, we consider $\widetilde{\mu}_{\text{PS}}$ where

$$\begin{aligned}
\widetilde{\mu}_{\text{PS}} &= \widehat{\mathbb{E}}_{\mathcal{S}} \{Y - g(\mathbf{X}^{\top} \bar{\boldsymbol{\alpha}})\} \{\exp(\mathbf{X}^{\top} \boldsymbol{\gamma}_0) + \boldsymbol{\Psi}^{\top} \bar{\boldsymbol{\beta}}\} + \widehat{\mathbb{E}}_{\mathcal{T}} g(\mathbf{X}^{\top} \bar{\boldsymbol{\alpha}}) \\
&+ \mathbb{E}_{\mathcal{S}} \{Y - g(\mathbf{X}^{\top} \bar{\boldsymbol{\alpha}})\} \exp(\mathbf{X}^{\top} \bar{\boldsymbol{\gamma}}) \mathbf{X}^{\top} \{\mathbb{E}_{\mathcal{S}} \mathbf{X} \exp(\mathbf{X}^{\top} \boldsymbol{\gamma}_0) \mathbf{X}^{\top}\}^{-1} \{\widehat{\mathbb{E}}_{\mathcal{T}} \mathbf{X} - \widehat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \exp(\mathbf{X}^{\top} \boldsymbol{\gamma}_0)\}.
\end{aligned}$$

Together with the results from Lemma (4), we have that $\mathbb{E} \widetilde{\mu}_{\text{PS}} = \mathbb{E}_{\mathcal{S}} Y \exp(\mathbf{X}^{\top} \bar{\boldsymbol{\gamma}}) = \mathbb{E}_{\mathcal{T}} Y = \mu_0$. By using the central limit theorem, we have that $\widetilde{\mu}_{\text{PS}} - \mu_0 = O_p(n^{-1/2})$, $n^{1/2}(\widetilde{\mu}_{\text{PS}} - \mu_0)$ weakly converges to gaussian distribution with mean $\mathbf{0}$. On the other hand, we have that

$$\begin{aligned}
\widehat{\mu}_{\text{PAD}} - \widetilde{\mu}_{\text{PS}} &= \widehat{\mathbb{E}}_{\mathcal{S}} \{Y - g(\mathbf{X}^{\top} \bar{\boldsymbol{\alpha}})\} \{\exp(\mathbf{X}^{\top} \bar{\boldsymbol{\gamma}}) \mathbf{X}^{\top} (\widehat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^{\top} (\widehat{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}})\} \\
&- \mathbb{E}_{\mathcal{S}} \{Y - g(\mathbf{X}^{\top} \bar{\boldsymbol{\alpha}})\} \exp(\mathbf{X}^{\top} \bar{\boldsymbol{\gamma}}) \mathbf{X}^{\top} \{\mathbb{E}_{\mathcal{S}} \mathbf{X} \exp(\mathbf{X}^{\top} \boldsymbol{\gamma}_0) \mathbf{X}^{\top}\}^{-1} \{\widehat{\mathbb{E}}_{\mathcal{T}} \mathbf{X} - \widehat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \exp(\mathbf{X}^{\top} \boldsymbol{\gamma}_0)\} \\
&+ o_p(n^{-1/2})
\end{aligned}$$

By using the techniques from (S3.6), we would have

$$\begin{aligned} & \widehat{\mathbb{E}}_{\mathcal{S}}\{Y - g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}})\} \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) \mathbf{X}^\top (\hat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}}) \\ & - \mathbb{E}_{\mathcal{S}}\{Y - g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}})\} \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) \mathbf{X}^\top \{\mathbb{E}_{\mathcal{S}} \mathbf{X} \exp(\mathbf{X}^\top \boldsymbol{\gamma}_0) \mathbf{X}^\top\}^{-1} \{\widehat{\mathbb{E}}_{\mathcal{T}} \mathbf{X} - \widehat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \exp(\mathbf{X}^\top \boldsymbol{\gamma}_0)\} \\ & = o_p(n^{-1/2}) \end{aligned}$$

And from Lemma A4, we have $\widehat{\boldsymbol{\beta}} = O_p(n^{-1/2})$. Thus, we have $\widehat{\mu}_{\text{PAD}} - \mu_0 = O_p(n^{-1/2})$. On the other hand, it is worth noticing that $\widehat{\boldsymbol{\beta}}$ is the continuous function of $\widehat{\boldsymbol{\theta}}$, $\widehat{\boldsymbol{\gamma}}$ and $\widehat{\boldsymbol{\alpha}}$, so under central limit theorem and Slutsky theorem, we would have the asymptotic normality of $\widehat{\boldsymbol{\beta}}$. Hence, we further have that $n^{1/2}(\widehat{\mu}_{\text{PS}} - \mu_0)$ weakly converges to gaussian distribution with mean $\mathbf{0}$. \square

Proof. Proof of Theorem 1 (ii).

First we denote \mathbf{U} as

$$\mathbf{U} = \text{Var}_{\mathcal{T}}(\mathbb{E}(Y|\mathbf{X})) + \mathbf{L}^\top \mathbb{E}_{\mathcal{S}} \mathbf{X} \mathbf{X}^\top \text{Var}(Y|\mathbf{X}) \mathbf{L}$$

When the OR is correctly specified, the asymptotic variance of $\widehat{\mu}_{\text{PAD}}$, $\text{Var}\{n^{-1/2}(\widehat{\mu}_{\text{PAD}} - \mu_0)\}$ is

$$\mathbb{E}_{\mathcal{S}}\{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^\top \bar{\boldsymbol{\beta}}\}^2 v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) + 2\mathbf{L}^\top \mathbb{E}_{\mathcal{S}} \mathbf{X} \{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^\top \bar{\boldsymbol{\beta}}\} v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) + \mathbf{U},$$

and $\bar{\boldsymbol{\beta}}$ contributes to minimizing this variance. When $\bar{\boldsymbol{\beta}} = \mathbf{0}$, the function above is written as

$$\mathbb{E}_{\mathcal{S}}\{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}})\}^2 v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) + 2\mathbf{L}^\top \mathbb{E}_{\mathcal{S}} \mathbf{X} \{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}})\} v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) + \mathbf{U},$$

which is the same as the asymptotic variance of $\hat{\mu}_{\text{DR}}$, $\text{Var}\{n^{-1/2}(\hat{\mu}_{\text{DR}} - \mu_0)\}$. Hence, when $\bar{\beta} \neq \mathbf{0}$, $\hat{\mu}_{\text{PAD}}$ has the smaller asymptotic variance than standard doubly robust estimator $\hat{\mu}_{\text{DR}}$. \square

Proof. Proof of Theorem 1 (iii).

When both the PS and OR is correctly specified, consider $\tilde{\mu}_B$, where

$$\tilde{\mu}_B = \widehat{\mathbb{E}}_{\mathcal{S}}\{Y - g(\mathbf{X}^\top \boldsymbol{\alpha}_0)\} \exp(\mathbf{X}^\top \boldsymbol{\gamma}_0) + \widehat{\mathbb{E}}_{\mathcal{T}}g(\mathbf{X}^\top \boldsymbol{\alpha}_0).$$

By using central limit theorem, $n^{1/2}(\tilde{\mu}_B - \mu_0)$ weakly converges to gaussian distribution with mean $\mathbf{0}$. On the other hand, by using Taylor series expansion, we would have $\hat{\mu}_{\text{PAD}} - \tilde{\mu}_B = o_p(n^{-1/2})$ and $\hat{\mu}_{\text{DR}} - \tilde{\mu}_B = o_p(n^{-1/2})$. Hence, they have the same asymptotic variance. \square

S3.2 Technical lemma

Lemma 1. Define $\mathbf{a} := \mathbb{E}_{\mathcal{S}}\boldsymbol{\Psi}\mathbf{X}^\top \dot{g}(\mathbf{X}^\top \bar{\boldsymbol{\alpha}})$, $\mathbf{b} := \mathbb{E}_{\mathcal{S}}\boldsymbol{\Psi} \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}})v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) + \mathbb{E}_{\mathcal{S}}\boldsymbol{\Psi}\mathbf{X}^\top v_{\bar{\boldsymbol{\theta}}}(\mathbf{X})\mathbf{L}$, and $\boldsymbol{\Sigma} := \mathbb{E}_{\mathcal{S}}\boldsymbol{\Psi}\boldsymbol{\Psi}^\top v_{\bar{\boldsymbol{\theta}}}(\mathbf{X})$, under Assumption 3, the solution of the RWLS problem in Algorithm 1 is

$$\bar{\boldsymbol{\beta}} = \boldsymbol{\Sigma}^{-1}\mathbf{a}(\mathbf{a}^\top \boldsymbol{\Sigma}^{-1}\mathbf{a})^{-1}\mathbf{a}^\top \boldsymbol{\Sigma}^{-1}\mathbf{b} - \boldsymbol{\Sigma}^{-1}\mathbf{b}.$$

Proof. First we introduce Lagrange multiplier $\boldsymbol{\lambda}$ and write (2.5) as the Lagrange form:

$$\bar{\boldsymbol{\beta}} = \operatorname{argmin}_{\boldsymbol{\beta}} \mathbb{E}_{\mathcal{S}}\{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^\top \boldsymbol{\beta}\}^2 v_{\bar{\boldsymbol{\theta}}}(\mathbf{X})$$

$$+2\mathbf{L}^\top \mathbb{E}_{\mathcal{S}} \mathbf{X} \{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^\top \boldsymbol{\beta}\} v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) - \boldsymbol{\lambda}^\top \mathbb{E}_{\mathcal{S}} \mathbf{X} \dot{g}(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}) \boldsymbol{\Psi}^\top \boldsymbol{\beta}.$$

Then we have the partial derivative of $\boldsymbol{\lambda}$ and $\boldsymbol{\beta}$:

$$\mathbb{E}_{\mathcal{S}} \mathbf{X} \dot{g}(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}) \boldsymbol{\Psi}^\top \boldsymbol{\beta} = \mathbf{0}, \quad (\text{S3.7})$$

and

$$2\mathbb{E}_{\mathcal{S}} \boldsymbol{\Psi} \{\exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) + \boldsymbol{\Psi}^\top \boldsymbol{\beta}\} v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) + 2\mathbb{E}_{\mathcal{S}} \boldsymbol{\Psi} \mathbf{X}^\top v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) \mathbf{L} - \mathbb{E}_{\mathcal{S}} \boldsymbol{\Psi} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}) \boldsymbol{\lambda} = \mathbf{0}. \quad (\text{S3.8})$$

From (S3.8) we have

$$\boldsymbol{\beta} = \{2\mathbb{E}_{\mathcal{S}} \boldsymbol{\Psi} \boldsymbol{\Psi}^\top v_{\bar{\boldsymbol{\theta}}}(\mathbf{X})\}^{-1} \{\mathbb{E}_{\mathcal{S}} \boldsymbol{\Psi} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}) \boldsymbol{\lambda} - 2\mathbb{E}_{\mathcal{S}} \boldsymbol{\Psi} \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) - 2\mathbb{E}_{\mathcal{S}} \boldsymbol{\Psi} \mathbf{X}^\top v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) \mathbf{L}\},$$

together with (S3.7), we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{S}} \mathbf{X} \dot{g}(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}) \boldsymbol{\Psi}^\top \{\mathbb{E}_{\mathcal{S}} \boldsymbol{\Psi} \boldsymbol{\Psi}^\top v_{\bar{\boldsymbol{\theta}}}(\mathbf{X})\}^{-1} \\ & * \{\mathbb{E}_{\mathcal{S}} \boldsymbol{\Psi} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \bar{\boldsymbol{\alpha}}) \boldsymbol{\lambda} - 2\mathbb{E}_{\mathcal{S}} \boldsymbol{\Psi} \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) - 2\mathbb{E}_{\mathcal{S}} \boldsymbol{\Psi} \mathbf{X}^\top v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) \mathbf{L}\} = \mathbf{0}. \end{aligned}$$

this function can be simplified as

$$\mathbf{a}^\top \boldsymbol{\Sigma}^{-1} (\mathbf{a} \boldsymbol{\lambda} - 2\mathbf{b}) = \mathbf{0},$$

and we further have

$$\boldsymbol{\lambda} = 2(\mathbf{a}^\top \boldsymbol{\Sigma}^{-1} \mathbf{a})^{-1} \mathbf{a}^\top \boldsymbol{\Sigma}^{-1} \mathbf{b}.$$

Hence, we have

$$\bar{\boldsymbol{\beta}} = \boldsymbol{\Sigma}^{-1} \mathbf{a} (\mathbf{a}^\top \boldsymbol{\Sigma}^{-1} \mathbf{a})^{-1} \mathbf{a}^\top \boldsymbol{\Sigma}^{-1} \mathbf{b} - \boldsymbol{\Sigma}^{-1} \mathbf{b}.$$

□

Lemma 2. *Under Assumptions 1-4, we have that $\widehat{\Psi} - \Psi = O_p(n^{-1/2})$.*

Proof. By definition, we would have that

$$\widehat{\Psi} - \Psi = \frac{\widehat{\mathbb{E}}_{\mathcal{T}}\{\Phi v_{\widehat{\theta}}(\mathbf{X})\}}{\widehat{\mathbb{E}}_{\mathcal{T}}v_{\widehat{\theta}}(\mathbf{X})} - \frac{\mathbb{E}_{\mathcal{T}}\{\Phi v_{\bar{\theta}}(\mathbf{X})\}}{\mathbb{E}_{\mathcal{T}}v_{\bar{\theta}}(\mathbf{X})}.$$

Under Assumption 4, we have that

$$\widehat{\mathbb{E}}_{\mathcal{T}}v_{\widehat{\theta}}(\mathbf{X}) - \mathbb{E}_{\mathcal{T}}v_{\bar{\theta}}(\mathbf{X}) = \widehat{\mathbb{E}}_{\mathcal{T}}v_{\bar{\theta}}(\mathbf{X}) + \widehat{\mathbb{E}}_{\mathcal{T}} \frac{\partial v_{\theta}(\mathbf{X})}{\partial \theta} \Big|_{\bar{\theta}} (\widehat{\theta} - \bar{\theta}) - \mathbb{E}_{\mathcal{T}}v_{\bar{\theta}}(\mathbf{X}) = O_p(n^{-1/2}) \quad (\text{S3.9})$$

for $\tilde{\theta}$ between $\widehat{\theta}$ and $\bar{\theta}$. By using the same techniques, we have that

$\widehat{\mathbb{E}}_{\mathcal{T}}\{\Phi v_{\widehat{\theta}}(\mathbf{X})\} - \mathbb{E}_{\mathcal{T}}\{\Phi v_{\bar{\theta}}(\mathbf{X})\} = O_p(n^{-1/2})$. And we have

$$\begin{aligned} \widehat{\Psi} - \Psi &= \frac{\widehat{\mathbb{E}}_{\mathcal{T}}\{\Phi v_{\widehat{\theta}}(\mathbf{X})\} \mathbb{E}_{\mathcal{T}}v_{\bar{\theta}}(\mathbf{X}) - \mathbb{E}_{\mathcal{T}}\{\Phi v_{\bar{\theta}}(\mathbf{X})\} \widehat{\mathbb{E}}_{\mathcal{T}}v_{\widehat{\theta}}(\mathbf{X})}{\widehat{\mathbb{E}}_{\mathcal{T}}v_{\widehat{\theta}}(\mathbf{X}) \mathbb{E}_{\mathcal{T}}v_{\bar{\theta}}(\mathbf{X})} \\ &= \frac{O_p(n^{-1/2}) \mathbb{E}_{\mathcal{T}}v_{\bar{\theta}}(\mathbf{X}) - \mathbb{E}_{\mathcal{T}}\{\Phi v_{\bar{\theta}}(\mathbf{X})\} O_p(n^{-1/2})}{\{\mathbb{E}_{\mathcal{T}}v_{\bar{\theta}}(\mathbf{X}) + O_p(n^{-1/2})\} \mathbb{E}_{\mathcal{T}}v_{\bar{\theta}}(\mathbf{X})} = O_p(n^{-1/2}). \end{aligned}$$

□

Lemma 3. *Under Assumptions 1 and 2, we have that $\widehat{\gamma} - \bar{\gamma} = O_p(n^{-1/2})$*

and $\widehat{\alpha} - \bar{\alpha} = O_p(n^{-1/2})$.

Proof. The estimation of γ has been given as

$$\widehat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \exp(\mathbf{X}^{\top} \widehat{\gamma}) = \widehat{\mathbb{E}}_{\mathcal{T}} \mathbf{X},$$

by applying Taylor series expansion, we have

$$n^{-1} \sum_{i=1}^n \mathbf{X}_i \exp(\mathbf{X}_i^\top \bar{\boldsymbol{\gamma}}) + n^{-1} \sum_{i=1}^n \mathbf{X}_i \exp(\mathbf{X}_i^\top \tilde{\boldsymbol{\gamma}}) \mathbf{X}_i^\top (\hat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}}) = N^{-1} \sum_{i=n+1}^{n+N} \mathbf{X}_i,$$

where $\tilde{\boldsymbol{\gamma}}$ is some vector between $\hat{\boldsymbol{\gamma}}$ and $\bar{\boldsymbol{\gamma}}$. According to (Van der Vaart,

2000, Chapter 5), we have $\hat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}} = o_p(1)$. Let \mathbf{J} represent matrix $n^{-1} \sum_{i=1}^n \mathbf{X}_i \exp(\mathbf{X}_i^\top \tilde{\boldsymbol{\gamma}}) \mathbf{X}_i^\top$,

and we have that

$$\mathbf{J} = n^{-1} \sum_{i=1}^n \mathbf{X}_i \exp(\mathbf{X}_i^\top \tilde{\boldsymbol{\gamma}}) \mathbf{X}_i^\top + n^{-1} \sum_{i=1}^n \mathbf{X}_i \exp(\mathbf{X}_i^\top \boldsymbol{\gamma}^*) \mathbf{X}_i^\top \mathbf{X}_i (\tilde{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}}) = \mathbb{E}_S \mathbf{X} \exp(\mathbf{X}^\top \tilde{\boldsymbol{\gamma}}) \mathbf{X}^\top + o_p(1)$$

for $\boldsymbol{\gamma}^*$ between $\tilde{\boldsymbol{\gamma}}$ and $\bar{\boldsymbol{\gamma}}$. Hence, by central limit theorem and Slutsky

theorem, we have that,

$$\begin{aligned} \hat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}} &= \mathbf{J}^{-1} \left\{ N^{-1} \sum_{i=n+1}^{n+N} \mathbf{X}_i - n^{-1} \sum_{i=1}^n \mathbf{X}_i \exp(\mathbf{X}_i^\top \tilde{\boldsymbol{\gamma}}) \right\} \\ &= \mathbf{J}^{-1} \left\{ N^{-1} \sum_{i=n+1}^{n+N} \mathbf{X}_i - \mathbb{E}_{\mathcal{T}} \mathbf{X} + \mathbb{E}_S \mathbf{X} \exp(\mathbf{X}^\top \tilde{\boldsymbol{\gamma}}) - n^{-1} \sum_{i=1}^n \mathbf{X}_i \exp(\mathbf{X}_i^\top \tilde{\boldsymbol{\gamma}}) \right\} = O_p(n^{-1/2}). \end{aligned}$$

Furthermore, The estimation equation of $\hat{\boldsymbol{\alpha}}$ is given by

$$\widehat{\mathbb{E}}_S \mathbf{S}(\hat{\boldsymbol{\alpha}}) = \widehat{\mathbb{E}}_S \mathbf{X} \{Y - g(\mathbf{X}^\top \hat{\boldsymbol{\alpha}})\} = \mathbf{0},$$

by using Taylor series expansion, we have that

$$\widehat{\mathbb{E}}_S \mathbf{X} \{Y - g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}})\} + \widehat{\mathbb{E}}_S \left. \frac{\partial \mathbf{S}(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}^\top} \right|_{\bar{\boldsymbol{\alpha}}} (\hat{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}}) = \mathbf{0}$$

for $\bar{\boldsymbol{\alpha}}$ between $\hat{\boldsymbol{\alpha}}$ and $\bar{\boldsymbol{\alpha}}$, and we have

$$\hat{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}} = -\widehat{\mathbb{E}}_S \left\{ \left. \frac{\partial \mathbf{S}(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}^\top} \right|_{\bar{\boldsymbol{\alpha}}} \right\}^{-1} \widehat{\mathbb{E}}_S \mathbf{X} \{Y - g(\mathbf{X}^\top \bar{\boldsymbol{\alpha}})\}.$$

By using the same techniques as those for obtaining the asymptotic properties of $\widehat{\gamma}$, under Assumptions 1 and 2, we have $\widehat{\boldsymbol{\alpha}} - \bar{\boldsymbol{\alpha}} = O_p(n^{-1/2})$. \square

Lemma 4. *Under Assumptions 1–4 and by Lemmas 1–3, we can obtain that $\widehat{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}} = O_p(n^{-1/2})$. In addition, when the PS is correctly specified, we further have $\bar{\boldsymbol{\beta}} = \mathbf{0}$ and $\widehat{\boldsymbol{\beta}} = O_p(n^{-1/2})$.*

Proof. By using the same techniques as (S3.9), under Condition 2-4, we first have that

$$\begin{aligned}\widehat{\mathbf{a}} - \mathbf{a} &= \widehat{\mathbb{E}}_{\mathcal{S}} \widehat{\boldsymbol{\Psi}} \mathbf{X}^{\top} \dot{g}(\mathbf{X}^{\top} \widehat{\boldsymbol{\alpha}}) - \mathbb{E}_{\mathcal{S}} \boldsymbol{\Psi} \mathbf{X}^{\top} \dot{g}(\mathbf{X}^{\top} \bar{\boldsymbol{\alpha}}) \\ &= \widehat{\mathbb{E}}_{\mathcal{S}} \boldsymbol{\Psi} \mathbf{X}^{\top} \dot{g}(\mathbf{X}^{\top} \bar{\boldsymbol{\alpha}}) - \mathbb{E}_{\mathcal{S}} \boldsymbol{\Psi} \mathbf{X}^{\top} \dot{g}(\mathbf{X}^{\top} \bar{\boldsymbol{\alpha}}) + O_p(n^{-1/2}) = O_p(n^{-1/2}).\end{aligned}$$

In addition, we can have that $\widehat{\mathbf{b}} - \mathbf{b} = O_p(n^{-1/2})$ and $\widehat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma} = O_p(n^{-1/2})$.

Furthermore, we can easily have that

$$\begin{aligned}\widehat{\boldsymbol{\Sigma}}^{-1} - \boldsymbol{\Sigma}^{-1} &= \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma} \{ \boldsymbol{\Sigma} + O_p(n^{-1/2}) \}^{-1} - \boldsymbol{\Sigma}^{-1} \\ &= \boldsymbol{\Sigma}^{-1} [\boldsymbol{\Sigma} \{ \boldsymbol{\Sigma} + O_p(n^{-1/2}) \}^{-1} - \{ \boldsymbol{\Sigma} + O_p(n^{-1/2}) \} \{ \boldsymbol{\Sigma} + O_p(n^{-1/2}) \}^{-1}] = O_p(n^{-1/2}),\end{aligned}$$

based on which we can have $(\widehat{\mathbf{a}}^{\top} \widehat{\boldsymbol{\Sigma}}^{-1} \widehat{\mathbf{a}})^{-1} - (\mathbf{a}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{a})^{-1} = O_p(n^{-1/2})$. More-

over, let $\widehat{\boldsymbol{\Omega}}$ denote $\widehat{\boldsymbol{\Sigma}}^{-1} \widehat{\mathbf{a}} (\widehat{\mathbf{a}}^{\top} \widehat{\boldsymbol{\Sigma}}^{-1} \widehat{\mathbf{a}})^{-1} \widehat{\mathbf{a}}^{\top} \widehat{\boldsymbol{\Sigma}}^{-1}$ and $\boldsymbol{\Omega}$ denote $\boldsymbol{\Sigma}^{-1} \mathbf{a} (\mathbf{a}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{a})^{-1} \mathbf{a}^{\top} \boldsymbol{\Sigma}^{-1}$,

and we can have that $\widehat{\boldsymbol{\Omega}} - \boldsymbol{\Omega} = O_p(n^{-1/2})$. Hence, we have that $\widehat{\boldsymbol{\beta}} - \bar{\boldsymbol{\beta}} =$

$$\widehat{\boldsymbol{\Omega}} \widehat{\mathbf{b}} - \boldsymbol{\Omega} \mathbf{b} = O_p(n^{-1/2}).$$

On the other hand, when the PS is correctly specified, $\mathbf{L} = \mathbf{0}$ and

$\mathbb{E}_{\mathcal{S}} \Psi \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) = \mathbb{E}_{\mathcal{T}} \Psi v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) = \mathbf{0}$, which means

$$\bar{\boldsymbol{\beta}} = \boldsymbol{\Omega} \mathbf{b} = \boldsymbol{\Omega} \{ \mathbb{E}_{\mathcal{S}} \Psi \exp(\mathbf{X}^\top \bar{\boldsymbol{\gamma}}) v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) + \mathbb{E}_{\mathcal{S}} \Psi \mathbf{X}^\top v_{\bar{\boldsymbol{\theta}}}(\mathbf{X}) \mathbf{L} \} = \boldsymbol{\Omega} \mathbf{0} = \mathbf{0}.$$

And at the same time, we have $\hat{\boldsymbol{\beta}} = O_p(n^{-1/2})$. □

S4 The mPAD estimator

S4.1 Technical details

To address the range issue of the PS model in aPAD, we propose a multiplicative PAD (mPAD) estimator in Algorithm 3, in which the PS model is taken as $\exp(\mathbf{X}^\top \hat{\boldsymbol{\gamma}} + \hat{\boldsymbol{\Psi}}^\top \hat{\boldsymbol{\beta}}')$ with a multiplicative augmentation term $\exp(\hat{\boldsymbol{\Psi}}^\top \hat{\boldsymbol{\beta}}')$ added to the original PS model.

Algorithm 3 Multiplicative PAD estimation

[Step 1] Solve the estimating equations in (2.1) to obtain $\hat{\gamma}$ and $\hat{\alpha}$, and obtain the conditional variance estimator as $\hat{\theta}$.

[Step 2] Specify $\Phi = \phi(\mathbf{X})$ of larger dimensionality than \mathbf{X} using any basis function $\phi(\cdot)$, and take $\hat{\Psi}' = \Phi - \hat{\mathbb{E}}_{\mathcal{T}}[\Phi \exp(\mathbf{X}^\top \hat{\gamma}) v_{\hat{\theta}}(\mathbf{X})] / \hat{\mathbb{E}}_{\mathcal{T}} \exp(\mathbf{X}^\top \hat{\gamma}) v_{\hat{\theta}}(\mathbf{X})$.

[Step 3] Solve the optimization problem:

$$\hat{\beta}' = \operatorname{argmin}_{\beta'} \hat{V}'_{\mu}(\beta'), \quad \text{s.t.} \quad \hat{\mathbb{E}}_{\mathcal{S}} \mathbf{X}^\top \dot{g}(\mathbf{X}^\top \hat{\alpha}) \exp(\mathbf{X}^\top \hat{\gamma}) \{\exp(\hat{\Psi}'^\top \beta') - 1\} = \mathbf{0}, \quad (\text{S4.10})$$

where

$$\hat{V}'_{\mu}(\beta') = \hat{\mathbb{E}}_{\mathcal{S}} \exp\{2(\mathbf{X}^\top \hat{\gamma} + \hat{\Psi}'^\top \beta')\} v_{\hat{\theta}}(\mathbf{X}) + 2\hat{\mathbf{L}}^\top \hat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \exp(\mathbf{X}^\top \hat{\gamma} + \hat{\Psi}'^\top \beta') v_{\hat{\theta}}(\mathbf{X}), \quad (\text{S4.11})$$

and $\hat{\mathbf{L}}' = -\hat{\mathbf{H}}^{-1} \left\{ \hat{\mathbb{E}}_{\mathcal{S}} \mathbf{X} \dot{g}(\mathbf{X}^\top \hat{\alpha}) \exp(\mathbf{X}^\top \hat{\gamma} + \hat{\Psi}'^\top \beta') - \hat{\mathbb{E}}_{\mathcal{T}} \mathbf{X} \dot{g}(\mathbf{X}^\top \hat{\alpha}) \right\}$.

[Step 4] Obtain the mPAD estimator through

$$\hat{\mu}_{\text{mPAD}} = \hat{\mathbb{E}}_{\mathcal{S}} \{Y - g(\mathbf{X}^\top \hat{\alpha})\} \exp(\mathbf{X}^\top \hat{\gamma} + \hat{\Psi}'^\top \hat{\beta}') + \hat{\mathbb{E}}_{\mathcal{T}} g(\mathbf{X}^\top \hat{\alpha}).$$

Similar to our discussion in Section 2.3, suppose that all estimators converge to their limiting values. Let $\Psi' = \Phi - \mathbb{E}_{\mathcal{T}} \Phi \exp(\mathbf{X}^\top \bar{\gamma}) v_{\bar{\theta}}(\mathbf{X}) / \mathbb{E}_{\mathcal{T}} \exp(\mathbf{X}^\top \bar{\gamma}) v_{\bar{\theta}}(\mathbf{X})$,

$$\mathbf{L}' = -\mathbf{H}^{-1} \left\{ \mathbb{E}_{\mathcal{S}} \mathbf{X} \dot{g}(\mathbf{X}^\top \bar{\alpha}) \exp(\mathbf{X}^\top \bar{\gamma} + \Psi'^\top \bar{\beta}') - \mathbb{E}_{\mathcal{T}} \mathbf{X} \dot{g}(\mathbf{X}^\top \bar{\alpha}) \right\},$$

$\bar{\beta}'$ be the limits of $\hat{\beta}'$, and

$$V'_{\mu}(\beta') = \mathbb{E}_{\mathcal{S}} \exp\{2(\mathbf{X}^\top \bar{\gamma} + \Psi'^\top \beta')\} v_{\bar{\theta}}(\mathbf{X}) + 2\mathbf{L}'^\top \mathbb{E}_{\mathcal{S}} \mathbf{X} \exp(\mathbf{X}^\top \bar{\gamma} + \Psi'^\top \beta') v_{\bar{\theta}}(\mathbf{X})$$

be the limiting function of $\hat{V}'_{\mu}(\beta')$ specified in Algorithm 3. Similarly to

aPAD, we shall comment on two scenarios: (i) correct PS model; and (ii) correct OR and wrong PS, to demonstrate the properties of the proposed mPAD estimator. In specific, similar to aPAD, the mPAD estimator can attain the DR property and a smaller variance than the standard DR under wrong PS. However, one should note that the optimization problem (S4.10) in Algorithm 3 is non-convex and not ensured to converge to the global solution. In addition, as will be discussed in Appendix S4.2, mPAD tends to have worse and less stable finite-sample performance than aPAD due to the presence of squared exponential terms like $\exp\{2(\mathbf{X}^\top \hat{\boldsymbol{\gamma}} + \hat{\boldsymbol{\Psi}}^\top \boldsymbol{\beta}')\}$ in mPAD. Therefore, we still recommend aPAD as the default choice when implementing PAD.

Correct PS model. When the PS model is correct, we have that $\mathbf{L}' = \mathbf{0}$ using the linear constraint in (S4.10) and

$$\frac{\partial V'_\mu(\boldsymbol{\beta}')}{\partial \boldsymbol{\beta}'} \Big|_{\boldsymbol{\beta}'=\mathbf{0}} = 2\mathbb{E}_{\mathcal{S}} \boldsymbol{\Psi}' \exp\{2(\mathbf{X}^\top \hat{\boldsymbol{\gamma}})\} v_{\hat{\boldsymbol{\theta}}}(\mathbf{X}) = 2\mathbb{E}_{\mathcal{T}} \boldsymbol{\Psi}' \exp(\mathbf{X}^\top \hat{\boldsymbol{\gamma}}) v_{\hat{\boldsymbol{\theta}}}(\mathbf{X}) = \mathbf{0}$$

by the definition of $\boldsymbol{\Psi}'$. Thus, $\boldsymbol{\beta}' = \mathbf{0}$ is a (local) minimizer of $V'_\mu(\boldsymbol{\beta}')$, which implies that $\hat{\boldsymbol{\beta}}'$ converges to $\mathbf{0}$ when PS model is correct. Meanwhile, the PS augmentation does not affect the OR model. Therefore, $\hat{\boldsymbol{\mu}}_{\text{mPAD}}$ have the same double robustness as $\hat{\boldsymbol{\mu}}_{\text{DR}}$.

Wrong PS and correct OR. When the OR model is correctly specified and

PS is misspecified, $\hat{\mu}_{\text{mPAD}}$ is asymptotically equivalent with

$$\tilde{\mu}_{\text{mPAD}} = \widehat{\mathbb{E}}_{\mathcal{S}}\{Y - g(\mathbf{X}^{\top}\bar{\boldsymbol{\alpha}})\} \exp(\mathbf{X}^{\top}\bar{\boldsymbol{\gamma}} + \boldsymbol{\Psi}^{\top}\bar{\boldsymbol{\beta}}') + \widehat{\mathbb{E}}_{\mathcal{T}}g(\mathbf{X}^{\top}\bar{\boldsymbol{\alpha}}) + \mathbf{L}^{\top}\widehat{\mathbb{E}}_{\mathcal{S}}\mathbf{X}\{Y - g(\mathbf{X}^{\top}\bar{\boldsymbol{\alpha}})\}.$$

This implies that $\text{aVar}\{n^{1/2}(\hat{\mu}_{\text{mPAD}} - \mu_0)\}$ is equal to

$$\mathbb{E}_{\mathcal{S}} \exp^2(\mathbf{X}^{\top}\bar{\boldsymbol{\gamma}} + \boldsymbol{\Psi}^{\top}\bar{\boldsymbol{\beta}}')v(\mathbf{X}) + 2\mathbf{L}^{\top}\mathbb{E}_{\mathcal{S}}\mathbf{X} \exp(\mathbf{X}^{\top}\bar{\boldsymbol{\gamma}} + \boldsymbol{\Psi}^{\top}\bar{\boldsymbol{\beta}}')v(\mathbf{X}) + C, \tag{S4.12}$$

which, after dropping the invariant C , is equal to $V'_{\mu}(\boldsymbol{\beta}')$. When the PS model is misspecified, $\partial V_{\mu}(\boldsymbol{\beta}')/\partial\boldsymbol{\beta}'$ is typically not $\mathbf{0}$ at $\boldsymbol{\beta}' = \mathbf{0}$ so the minimizer is non-zero. Thus, when the OR model is correct, we have $\text{aVar}\{n^{1/2}(\hat{\mu}_{\text{mPAD}} - \mu_0)\} \leq \text{aVar}\{n^{1/2}(\hat{\mu}_{\text{DR}} - \mu_0)\}$, and the strict “ $<$ ” will hold in general when PS is wrong.

S4.2 Simulation study

We consider the same settings (L1)–(L3) as introduced in Section 4 and evaluate the performance of mPAD. To solve for the non-convex problem (S4.10), we adopt an iterative Newton-type optimization procedure with a quadratic approximation to the objective function $\widehat{V}'_{\mu}(\boldsymbol{\beta}')$ and a linear approximation to the constraint in (S4.10) at each iteration. Also, we impose an early stop whenever the objective \widehat{V}'_{μ} will not decrease in the next iteration.

Table 1: The absolute average bias (Bias), standard error (SE), and coverage probability (CP) of the 95% confidence intervals of the mPAD estimators, Relative efficiency (RE) between DR and mPAD, i.e., $\text{Var}(\hat{\mu}_{\text{DR}})/\text{Var}(\hat{\mu}_{\text{mPAD}})$ under the settings described in Section 4. All results are based on 1000 repetitions.

Setting	$n = N = 500$				$n = N = 1000$			
	Bias	SE	CP	RE	bias	SE	CP	RE
(L1)	0.001	0.054	0.95	1.09	0.002	0.039	0.94	1.07
(L2)	0.005	0.054	0.96	1.01	0.001	0.041	0.94	0.96
(L3)	0.002	0.051	0.94	1.23	0.001	0.035	0.94	1.18

In Table 1, we report mPAD’s bias, standard error (SE), and coverage probability (CP) of the 95% CI, and relative efficiency (RE) to the standard DR. Similar to the aPAD estimator studied in Section 4, mPAD performs closely to the standard DR estimator when the PS model is correct, and attains better estimation efficiency than DR when the OR model is correct and PS model is misspecified. Also, mPAD shows similar performance to aPAD when comparing Table 1 with Table 1.

Meanwhile, mainly due to the non-convexity issue and the outlying exponential terms $\exp\{2(\mathbf{X}^\top \hat{\boldsymbol{\gamma}} + \hat{\boldsymbol{\Psi}}^\top \boldsymbol{\beta}')\}$, mPAD performs poorly in the Gaussian linear settings (G1)–(G3) supposed to be reasonable for the aPAD and standard DR methods. For example, under misspecified OR models,

the mPAD estimator has a bias more than 0.06 in all settings (G1)–(G3) under $N = n = 1000$, which is substantially larger than aPAD.

Bibliography

Bang, H. and Robins, J. M. (2005). Doubly robust estimation in missing data and causal inference models. *Biometrics*, 61(4):962–973.

Hernán, M. A. and Robins, J. M. (2010). Causal inference.

Shu, H. and Tan, Z. (2018). Improved estimation of average treatment effects on the treated: Local efficiency, double robustness, and beyond. *arXiv preprint arXiv:1808.01408*.

Van der Vaart, A. W. (2000). *Asymptotic statistics*, volume 3. Cambridge university press.