

**Semi-supervised Regression Analysis with
Model Misspecification and High-dimensional Data**

Ye Tian^a, Peng Wu^b and Zhiqiang Tan^{c}*

^a*The University of Hong Kong,*

^b*Beijing Technology and Business University,*

^c*Rutgers University*

Supplementary Material

This Supplementary Material consists of Sections S1–S10. Section S1 provides a heuristic discussion on conditions for the proposed estimator to be \sqrt{N} -consistent and asymptotic normal. Section S2 provides a comparison of our paper with several related papers with regression of Y on high-dimensional $\mathbf{Z} = \mathbf{X}$. Section S3 introduces regularity conditions for the analysis of asymptotic properties of the proposed estimators. Section S4 contains technical tools used in proofs of lemmas in Section S5. Section S5 presents lemmas used in proofs of Proposition 2 and Theorem 1. Section S6 gives the technical proofs of Proposition 1 and Proposition 2. Section S7 provides the proof of Theorem 1. Section S8 includes contents of the extension to stratified sampling settings, including the proof of Proposition 3 and the variance

*Correspondence to: ztan@stat.rutgers.edu.

comparison. Sections S9 and S10 present details of the numerical implementation and application, respectively.

S1 Discussion on asymptotic normality of the proposed estimator

First, we explain why, in case (a) of $\mathbf{Z} = 1$, \sqrt{N} -consistency and asymptotic normality can be achieved using a linear OR working model, provided that either the OR model or the PS model is correctly specified. In the general case of $\mathbf{Z} \in \mathbb{R}^m$, under certain regularity conditions, we have the following expansion

$$\begin{aligned} & \tau(\mathbf{O}, \hat{\alpha}, \hat{\beta}, \hat{\gamma}) - \tau(\mathbf{O}, \bar{\alpha}, \beta^*, \bar{\gamma}) \\ &= \frac{\partial \tau(\mathbf{O}, \alpha, \beta, \gamma)}{\partial \alpha} \Big|_{(\alpha=\bar{\alpha}, \beta=\beta^*, \gamma=\bar{\gamma})}^{\text{T}} (\hat{\alpha} - \bar{\alpha}) + \frac{\partial \tau(\mathbf{O}, \alpha, \beta, \gamma)}{\partial \gamma} \Big|_{(\alpha=\bar{\alpha}, \beta=\beta^*, \gamma=\bar{\gamma})}^{\text{T}} (\hat{\gamma} - \bar{\gamma}) + o_p\left(\frac{1}{\sqrt{N}}\right). \end{aligned}$$

To obtain \sqrt{N} -consistency, we need the population estimating equation in (α, γ) :

$$\begin{cases} \mathbb{E} \left[\left\{ \frac{\partial \tau(\mathbf{O}, \alpha, \beta, \gamma)}{\partial \gamma} \right\}^{\text{T}} \Big|_{(\alpha=\bar{\alpha}, \beta=\beta^*, \gamma=\bar{\gamma})} \right] = -\mathbb{E} \left[R \frac{1 - \pi(X; \bar{\gamma})}{\pi(X; \bar{\gamma})} \{Y - \psi(\bar{\alpha}^{\text{T}} \mathbf{G})\} \mathbf{Z} \otimes \mathbf{F} \right] = 0, \\ \mathbb{E} \left[\mathbf{Z} \left\{ \frac{\partial \tau(\mathbf{O}, \alpha, \beta, \gamma)}{\partial \alpha} \right\}^{\text{T}} \Big|_{(\alpha=\bar{\alpha}, \beta=\beta^*, \gamma=\bar{\gamma})} \right] = \mathbb{E} \left[\left\{ 1 - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \right\} \psi'(\bar{\alpha}^{\text{T}} \mathbf{G}) \mathbf{Z} \otimes \mathbf{G} \right] = 0, \end{cases} \quad (\text{S1.1})$$

similarly as discussed in Tan (2020a), Section 3.2. From the system of equations (S1.1), there are a total of $m \times (p + 1) + m \times (q + 1)$ equations, but $(p + 1) + (q + 1)$ unknown parameters, i.e., (α, γ) . Therefore, α and γ cannot be identified from the equations without additional conditions.

When $\mathbf{Z} = 1$ (i.e., case (a)) and the OR model is linear (i.e., $\psi(t) = t$ is the identity function), the system of equations (S1.1) becomes

$$\begin{cases} \mathbb{E} \left\{ R \frac{1 - \pi(\mathbf{X}; \bar{\gamma})}{\pi(\mathbf{X}; \bar{\gamma})} (Y - \bar{\alpha}^\top \mathbf{G}) \mathbf{F} \right\} = 0, \\ \mathbb{E} \left\{ \left(1 - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \right) \mathbf{G} \right\} = 0. \end{cases}$$

In this case, with $\mathbf{G} = \mathbf{F}$ as in Tan (2020a), the above system is just-identified, with the same number of equations and number of parameters, and naturally serves as population estimating equations for $\bar{\alpha}$ and $\bar{\gamma}$. This is the core idea of Tan (2020a), explaining why a \sqrt{N} -consistent and asymptotically normal estimator can be achieved in the case where $\mathbf{Z} = 1$, $\mathbf{G} = \mathbf{F}$, and a linear OR model is used, provided that either the OR model or the PS model is correctly specified.

Second, we explain why, for a general choice of \mathbf{Z} , the correct specification of the propensity score model is assumed to obtain a \sqrt{N} -consistent and asymptotically normal estimator.

From the above discussion of (S1.1), there are two key challenges for a general choice of \mathbf{Z} :

- When \mathbf{Z} is multi-dimensional, the number of equations and the number of parameters in (S1.1) do not match.
- When the OR model is nonlinear or $\mathbf{F} \neq \mathbf{G}$, the two equations in (S1.1) are coupled in terms of dependency on $\bar{\alpha}$ and $\bar{\gamma}$, so that a direct approach of defining $\bar{\alpha}$ and $\bar{\gamma}$ from the two equations would not work.

In our paper, we address the first challenge by carefully specifying the form of \mathbf{G} as follows

$$\mathbf{G} = [\mathbf{F}^T, \{\mathbf{Z} \otimes \mathbf{F}\}^T]^T.$$

Then we circumvent the second challenge by assuming that the PS model is correctly specified. Alternative strategies can be considered by extending the iterative approach of Ghosh and Tan (2022). Investigating these topics would be interesting for future work.

Finally, it is worth pointing out that correct specification of the PS model is automatically satisfied in the classic SSL and stratified sampling setup. In such cases, the PS model is constant and known, as noted in Remark 1.

S2 Comparison with additional related papers

We provide a comparison of our paper with several related papers with regression of Y on high-dimensional $\mathbf{Z} = \mathbf{X}$ as mentioned in Section 2.1.

First, we compare our paper with Chakraborty et al. (2019). Chakraborty et al. (2019) consider estimating a high-dimensional parameter θ^* , which is defined as the minimizer of

$$\mathbb{L}(\theta) := \mathbb{E}[L(Y, \mathbf{X}, \theta)],$$

where $L(Y, \mathbf{X}, \theta)$ is a convex loss function. Under the setting of *semi-supervised learning with covariate shift* (i.e., the conditional distributions of Y given \mathbf{X} in the labeled and unlabeled datasets are assumed to be the same, whereas the marginal distributions of \mathbf{X} are different), Chakraborty et al. (2019) identify $\mathbb{L}(\theta)$ using the doubly robust form (their equation (2.1)):

$$\mathbb{L}(\theta) = \mathbb{E}[\phi(\mathbf{X}, \theta)] + \mathbb{E} \left[\frac{T}{\pi(\mathbf{X})} \{L(Y, \mathbf{X}, \theta) - \phi(\mathbf{X}, \theta)\} \right], \quad (\text{S2.2})$$

where $T \in \{0, 1\}$ is the indicator of Y being observed, $\phi(\mathbf{X}, \theta) = \mathbb{E}[L(Y, \mathbf{X}, \theta) | \mathbf{X}]$ and $\pi(\mathbf{X}) = \mathbb{P}(T = 1 | \mathbf{X})$. Instead of dealing with the general loss function (S2.2), they focus on a specific yet common case where $L(Y, \mathbf{X}, \theta)$ satisfies (their equation (2.3)):

$$\nabla L(Y, \mathbf{X}, \theta) \equiv \frac{\partial}{\partial \theta} L(Y, \mathbf{X}, \theta) = \mathbf{h}(\mathbf{X}) \{Y - g(\mathbf{X}, \theta)\} \quad (\text{S2.3})$$

for some functions $\mathbf{h}(\mathbf{X})$ with the same dimension as θ . In this case, $\frac{\partial}{\partial \theta} \phi(\mathbf{X}, \theta) = \mathbf{h}(\mathbf{X}) \{m(\mathbf{X}) - g(\mathbf{X}, \theta)\}$, where $m(\mathbf{X}) = \mathbb{E}(Y | \mathbf{X})$. Note that (S2.3) can be motivated from a conditional mean model $m(\mathbf{X}) = g(\mathbf{X}, \theta^*)$, but such a model is allowed

to be misspecified. The authors propose a debiased and doubly robust (DDR) estimator $\hat{\theta}^{\text{DDR}}$ for θ^* , as a minimizer of the following objective function plus a LASSO penalty on θ (after ignoring cross-fitting):

$$\hat{\mathbb{L}}(\theta) = \tilde{\mathbb{E}} \left[\hat{\phi}(\mathbf{X}, \theta) + \frac{T}{\hat{\pi}(\mathbf{X})} \{L(Y, \mathbf{X}, \theta) - \hat{\phi}(\mathbf{X}, \theta)\} \right],$$

where $\hat{\pi}(\mathbf{X})$ is an estimator of $\pi(\mathbf{X})$ through a PS model, and $\hat{\phi}(\mathbf{X}, \theta)$ is an estimator of $\phi(\mathbf{X}, \theta)$, satisfying $\frac{\partial}{\partial \theta} \hat{\phi}(\mathbf{X}, \theta) = \mathbf{h}(\mathbf{X}) \{\hat{m}(\mathbf{X}) - g(\mathbf{X}, \theta)\}$, and $\hat{m}(\mathbf{X})$ is an estimator of $m(\mathbf{X})$ through an OR model. By simple calculation, the gradient of $\hat{\mathbb{L}}(\theta)$ is

$$\begin{aligned} \frac{\partial}{\partial \theta} \hat{\mathbb{L}}(\theta) &= \tilde{\mathbb{E}} \left[\mathbf{h}(\mathbf{X}) \{\hat{m}(\mathbf{X}) - g(\mathbf{X}, \theta)\} + \frac{T}{\hat{\pi}(\mathbf{X})} \mathbf{h}(\mathbf{X}) \{Y - \hat{m}(\mathbf{X})\} \right] \\ &= \tilde{\mathbb{E}} \left[\frac{T}{\hat{\pi}(\mathbf{X})} \mathbf{h}(\mathbf{X}) \{Y - g(\mathbf{X}, \theta)\} + \left(1 - \frac{T}{\hat{\pi}(\mathbf{X})}\right) \mathbf{h}(\mathbf{X}) \{\hat{m}(\mathbf{X}) - g(\mathbf{X}, \theta)\} \right]. \end{aligned}$$

This matches the AIPW estimating function (2.4) in our paper, by linking $\hat{\pi}(\mathbf{X})$, $\hat{m}(\mathbf{X})$, $g(\mathbf{X}, \theta)$, and $\mathbf{h}(\mathbf{X})$ above with $\hat{\pi}(\mathbf{X})$, $\hat{m}(\mathbf{X})$, $\psi(\beta^\top \mathbf{Z})$, and \mathbf{Z} in our paper. Hence, Chakraborty et al. (2019) and our paper intersects in the idea of using AIPW estimating functions.

Our paper, however, differs from Chakraborty et al. (2019) in two important technical aspects. (i) Chakraborty et al. (2019) allow θ to be a high-dimensional parameter and study the L_1 and L_2 convergence rates for the DDR estimator $\hat{\theta}^{\text{DDR}}$, which are slower than $N^{-1/2}$. (ii) Chakraborty et al. (2019) then propose a despar-

sified DDR estimator $\hat{\theta}^{\text{D-DDR}}$, and prove that each *scalar* element of $\hat{\theta}^{\text{D-DDR}}$ is \sqrt{N} -consistent and asymptotically normal (called entry-wise asymptotic normality) if both OR and PS models are correct such that $\hat{\pi}(\cdot)$ and $\hat{m}(\cdot)$ converge to the *true values*, $\pi(\cdot)$ and $m(\cdot)$, at fast enough rates (specifically, the product of the estimation errors is smaller than $N^{-1/2}$). In contrast, our setting corresponds to a regression model of Y on a low-dimensional sub-vector \mathbf{Z} of \mathbf{X} with coefficient vector θ^* , but our proposed estimator is shown to be \sqrt{N} -consistent and asymptotically normal for θ^* (not just entry-wise), even when $\hat{m}(\cdot)$ is inconsistent for the true value and $\hat{\pi}(\cdot)$ converges to the true value at a rate slower than $N^{-1/2}$ (but faster than $N^{-1/4}$), with misspecified OR model and correctly specified PS model. Hence the aforementioned product of estimation errors may be greater than $N^{-1/2}$. Note that $\hat{m}(\cdot)$ is consistent for the *target value*, different from the true value, and may also converge slower than $N^{-1/2}$ (but faster than $N^{-1/4}$), with misspecified OR model; see Section 3.2 of our paper. This is the main advantage of regularized calibrated estimation in our method. Therefore, our paper is not a special case of Chakraborty et al. (2019).

Second, we compare our work with Deng et al. (2023), Zhang et al. (2023), and Chen and Zhang (2023). The settings of these works also differ from ours. We summarize the differences below, which demonstrate our unique contributions.

- Deng et al. (2023) consider estimating θ^* defined by

$$\theta^* = \operatorname{argmin}_{\theta} \mathbb{E}(Y - \theta^T \mathbf{X})^2,$$

where \mathbf{X} is a high-dimensional covariate vector and Y is the response. Under the classical semi-supervised setup without covariate shift (i.e., (\mathbf{X}, Y) have the same joint distributions in the labeled and unlabeled datasets), the authors propose several Lasso/Dantzig selector-based estimators for θ , denoted as $\hat{\theta}$, and then derive the convergence rates of $\hat{\theta}$ to the target value θ^* in terms of the L_q -norm, $\|\hat{\theta} - \theta^*\|_q$, for $q = 1, 2$. In particular, the semi-supervised Lasso estimator in Deng et al. (2023) (equation 3.6 in their paper), after ignoring the Lasso penalty on θ , can be shown to be a solution to the following equation

$$\frac{1}{n} \sum_{i=1}^n [Y_i - \hat{m}(\mathbf{X}_i)] \mathbf{X}_i + \frac{1}{N} \sum_{i=1}^N [\hat{m}(\mathbf{X}_i) - \theta^T \mathbf{X}_i] \mathbf{X}_i = 0,$$

where $\hat{m}(\mathbf{X})$ is an estimator of $m(\mathbf{X})$ as in the earlier discussion. This matches the AIPW estimating equation (5.1) in the classical semi-supervised setting in our paper, by linking $\hat{m}(\mathbf{X}_i)$ and $\theta^T \mathbf{X}_i$ with $\psi(\hat{\alpha}^T \mathbf{G}_i)$ and $\psi(\beta^T \mathbf{Z}_i)$ in our paper.

There are several differences between Deng et al. (2023) and our paper. First, Deng et al. (2023) focus on the classical semi-supervised setup without covariate shift, whereas we focus on a setup that incorporates covariate shift. Second,

Deng et al. (2023) consider a linear model $\theta^{*\text{T}} \mathbf{X}$ directly with high-dimensional \mathbf{X} , whereas we consider a generalized linear model $\psi(\beta^{*\text{T}} \mathbf{Z})$ with possibly non-linear $\psi(\cdot)$ and low-dimensional \mathbf{Z} as a subvector of \mathbf{X} . Both these models are used as approximations and allowed to be misspecified. Finally, Deng et al. (2023) concentrate on analyzing $\|\hat{\theta} - \theta^*\|_q$, which is slower than $N^{-1/2}$, without dealing with inference such as confidence intervals, whereas our theoretical analysis provides \sqrt{N} -consistency and asymptotic normality for the proposed estimator of β^* , thereby enabling valid confidence intervals under suitable sparsity conditions.

- Zhang et al. (2023) consider a high-dimensional linear regression model:

$$Y = \theta^{*\text{T}} \mathbf{X} + \epsilon,$$

where θ^* is an unknown coefficient, $\mathbb{E}(\epsilon|\mathbf{X}) = 0$, and $\mathbb{E}(\epsilon^2|\mathbf{X}) = \sigma_\epsilon^2$. Under the setting of semi-supervised learning with covariate shift, their goal is to estimate the linear regression function, $m(x) = \theta^{*\text{T}}x$, for a given query vector x . The authors propose a debiased-Lasso type estimator for $m(x)$, depending on an estimator $\hat{\pi}(\mathbf{X})$ with a correctly specified PS model, and show its \sqrt{N} -consistency and asymptotic normality under suitable sparsity conditions.

There are also notable differences between Zhang et al. (2023) and our paper. First, Zhang et al. (2023) consider a linear model $\theta^{*\text{T}} \mathbf{X}$ directly with high-

dimensional \mathbf{X} , whereas we consider a generalized linear model $\psi(\beta^{*\text{T}}\mathbf{Z})$ with possibly nonlinear $\psi(\cdot)$ and low-dimensional \mathbf{Z} as a subvector of \mathbf{X} . This is similar to the second difference between Deng et al. (2023) and our paper discussed earlier, but the linear regression model in Zhang et al. (2023) is assumed to be correctly specified along with homoscedastic noises. Second, while the goal of Zhang et al. (2023) is inference about the linear regression function $m(x) = \theta^{*\text{T}}x$ at a given point x , our goal is inference about the entire vector of regression coefficients β^* . In fact, because β^* is low-dimensional in our theory, the asymptotic normality of our proposed estimator $\hat{\beta}$ directly gives asymptotic normality of the estimator $\hat{\beta}^{\text{T}}z$ and hence enables valid confidence intervals for $\beta^{*\text{T}}z$ and $\psi(\beta^{*\text{T}}z)$ at a given point z .

- Chen and Zhang (2023) consider the following regression models

$$Y_i = \theta^{*\text{T}}\mathbf{X}_i + \epsilon_i = \phi^*V_i + \alpha^{*\text{T}}\mathbf{W}_i + \epsilon_i, \quad (\text{S2.4})$$

$$V_i = \gamma^{*\text{T}}\mathbf{W}_i + \delta_i, \quad (\text{S2.5})$$

where $\mathbf{X}_i = (V_i, \mathbf{W}_i^{\text{T}})^{\text{T}}$, with V_i being a particular covariate, and $\theta^* = (\phi^*, \alpha^{*\text{T}})^{\text{T}}$.

Model (S2.4) is allowed to be misspecified, with θ^* defined as

$$\theta^* = \operatorname{argmin}_{\theta} \mathbb{E}(Y - \theta^{\text{T}}\mathbf{X})^2.$$

Model (S2.5) is assumed to be correctly specified with $E(\delta_i|\mathbf{W}_i) = 0$ for all i ,

which is automatically satisfied in the Gaussian design (i.e., $\mathbf{X}_i = (V_i, \mathbf{W}_i^T)^T$ is jointly Gaussian). Chen and Zhang (2023) consider semi-supervised learning with covariate shift satisfying $R_i \perp\!\!\!\perp V_i | \mathbf{W}_i$, where $R_i = 1$ or 0 for i from the labeled or unlabeled dataset. The authors propose various estimators for ϕ^* and show \sqrt{N} -consistency and asymptotic normality under suitable sparsity conditions. In contrast with Chen and Zhang (2023), our paper considers semi-supervised learning with covariate shift satisfying $R_i \perp\!\!\!\perp Y_i | \mathbf{X}_i$ (Assumption 1), and studies inference for a generalized linear model $\psi(\beta^{*T} \mathbf{Z})$ for Y with possibly nonlinear $\psi(\cdot)$ and low-dimensional \mathbf{Z} as a subvector of \mathbf{X} . Therefore, Chen and Zhang (2023) and our paper deal with different semi-supervised settings and objectives of estimation.

S3 Regularity assumptions

In this section, we introduce some regularity conditions for the analysis of asymptotic properties of the proposed estimators.

The following Assumption is taken from Tan (2020a), which are plausible as discussed there.

Assumption 1 (Regularity conditions for $\hat{\gamma}$). Suppose that the following conditions are satisfied:

- (i) $\max_{j=0, \dots, p} |f_j(\mathbf{X})| \leq C_0$ a.s. for a constant $C_0 \geq 1$;

- (ii) $\bar{\gamma}^\top \mathbf{F} \geq B_0$ a.s. for a constant $B_0 \in \mathbb{R}$, that is, $\pi(\mathbf{X}; \bar{\gamma})$ is bounded from below by $\{1 + \exp(-B_0)\}^{-1}$;
- (iii) the compatibility condition holds for $\Sigma_{\bar{\gamma}}$ with the subset $S_{\bar{\gamma}} = \{0\} \cup \{j : \bar{\gamma}_j \neq 0, j = 1, \dots, p\}$ and some constants $\nu_0 > 0$ and $\xi_0 > 1$, where $\Sigma_{\bar{\gamma}} = \mathbb{E}\{Rw(\mathbf{X}; \bar{\gamma})\mathbf{F}\mathbf{F}^\top\}$ is the Hessian of $\mathbb{E}\{\ell_{\text{CAL}}(\gamma)\}$ at $\gamma = \bar{\gamma}$;
- (iv) $|S_{\bar{\gamma}}|\lambda_0 \leq \zeta_0$ for a sufficiently small constant $\zeta_0 > 0$, depending only on (A_0, C_0, ξ_0, ν_0) , where $\lambda_0 = c_\gamma \sqrt{\ln\{(1+p)/\epsilon\}/N}$, c_γ is a constant only depending on (B_0, C_0) and $A_0 > (\xi_0 + 1)/(\xi_0 - 1)$ is a constant.

For studying the properties of $\hat{\alpha}$, we make the following Assumption 2.

Assumption 2 (Regularity conditions for $\hat{\alpha}$). Let $\psi_1(u)$ denote the derivative of $\psi(u)$. Suppose that the following conditions are satisfied:

- (i) $C_1 \leq \psi_1(\bar{\alpha}^\top \mathbf{G}) \leq C_2$ a.s. for positive constants (C_1, C_2) ;
- (ii) $\psi_1(u) \leq \psi_1(u') \exp(C_3|u - u'|)$ for any (u, u') and certain constant $C_3 \geq 0$;
- (iii) $\max_{j=0, \dots, q} |g_j(\mathbf{X})| \leq C_4$ a.s. for a constant $C_4 \geq 1$;
- (iv) $Y - \psi(\bar{\alpha}^\top \mathbf{G})$ is uniformly sub-Gaussian given \mathbf{X} : $D_0^2 \mathbb{E}(\exp[\{Y - \psi(\bar{\alpha}^\top \mathbf{G})\}^2 / D_0^2] - 1 | \mathbf{X}) \leq D_1^2$ for some positive constants (D_0, D_1) ;
- (v) the compatibility condition holds for $\Sigma_{\bar{\alpha}}$ with the subset $S_{\bar{\alpha}} = \{0\} \cup \{j : \bar{\alpha}_j \neq 0, j = 0, \dots, q\}$ and some constant $\nu_1 > 0$ and $\xi_1 > 1$, where $\Sigma_{\bar{\alpha}} =$

c_γ is defined in Section S6.1.

$$\mathbb{E}[Rw(\mathbf{X}; \bar{\gamma})\psi_1(\bar{\alpha}^\top \mathbf{G})\mathbf{G}\mathbf{G}^\top];$$

- (vi) $(1 + \xi_1)^2 \nu_1^{-2} |S_{\bar{\alpha}}| \lambda_1 \leq \zeta_1$ for a sufficiently small constant $\zeta_1 > 0$, where $\lambda_1 = \max[\lambda_0, c_\alpha \sqrt{\ln\{(1+q)/\epsilon\}/N}]$, c_α is a constant depending on $(B_0, C_2, C_4, D_0, D_1)$;
- (vii) let $A_1 > (\xi_1 + 1)/(\xi_1 - 1)$ be a constant. There exist $0 \leq \eta_2, \eta_3 < 1$, such that $\tilde{c}_\alpha |S_{\bar{\alpha}}| \lambda_1 \leq \eta_2$ and $\tilde{c}_\gamma |S_{\bar{\gamma}}| \lambda_0 \leq \eta_3$, where \tilde{c}_α and \tilde{c}_γ are both constants.

Assumptions 2(i)–(ii) are mild conditions on the smoothness of the inverse link function ψ . Commonly used functions like the identity and logit functions satisfy these requirements. Assumptions 2(iii)–(vii) are similar to those used in related analysis by Tan (2020a). A subtle difference is that the compatibility condition in Tan (2020a) [Assumption 2(ii)] is assumed for $\Sigma_{\bar{\gamma}}$ with $\mathbf{F} = \mathbf{G}$, whereas our compatibility condition is assumed for $\Sigma_{\bar{\alpha}}$. In our setting, \mathbf{G} has a higher dimension than \mathbf{F} (except in the case where $\mathbf{Z} = 1$).

Assumption 3. Suppose that the following conditions are satisfied:

- (i) $\mathbf{Z} = \{Z_0, \dots, Z_{m-1}\} \in \mathbb{R}^m$, and $\max_{j=0, \dots, m-1} |Z_j| \leq C_5$ almost surely for a constant $C_5 \geq 1$, where m is fixed as N increases;
- (ii) $\beta \in \Theta_\beta \subset \mathbb{R}^m$, and for $\forall \epsilon > 0$, $\inf_{\beta \in \Theta_\beta: \|\beta - \beta^*\|_1 \geq \epsilon} \mathbb{E}\{\|\tau(\mathbf{O}, \bar{\alpha}, \beta, \bar{\gamma})\|_1\} > 0$;
- (iii) $\mathbb{E}\{\sup_{\beta \in \Theta_\beta} \|\tau(\mathbf{O}, \bar{\alpha}, \beta, \bar{\gamma})\|_1\} < \infty$;
- (iv) $\psi_1(\beta^{*\top} \mathbf{Z}) \leq C_6$ a.e., for some constant $C_6 > 0$;

c_α is defined in Section S6.2.

\tilde{c}_α and \tilde{c}_γ depend on $(A_0, A_1, C_0, C_3, C_4, M_0, \nu_1, \xi_1, \zeta_0, \zeta_1)$ and $(A_0, A_1, C_0, C_1, C_3, C_4, D_0, D_1, M_0, \xi_1, \zeta_0)$, respectively. They are defined in Section S6.2.

(v) $|S_{\bar{\gamma}}| \sqrt{\ln(1+p)\ln(1+q)} = o(N^{1/2})$ and $|S_{\bar{\alpha}}| \ln(1+q) = o(N^{1/2})$.

Assumptions 3(ii)–(iii) are standard conditions in the asymptotic theory of estimating equations. Assumption 3(iv) is a mild condition on the smoothness of the inverse link function at the value of $\beta^{*\text{T}}\mathbf{Z}$. Assumption 3(v) is comparable to the sparsity requirement in Tan (2020a).

S4 Technical tools

We state the following concentration inequalities, to facilitate proofs of lemmas in Section S5, which can be obtained from Bühlmann and Van De Geer (2011), Lemmas 14.11, 14.16 & 14.9.

Lemma 1. Let (Y_1, \dots, Y_N) be independent variables such that $\mathbb{E}(Y_i) = 0$ for $i = 1, \dots, n$ and $\max_{i=1, \dots, n} |Y_i| \leq c_0$ for some constant c_0 . Then for any $t > 0$,

$$\mathbb{P} \left(\left| \frac{1}{N} \sum_{i=1}^N Y_i \right| > t \right) \leq 2 \exp \left(-\frac{nt^2}{2c_0^2} \right).$$

Lemma 2. Let (Y_1, \dots, Y_N) be independent variables such that $\mathbb{E}(Y_i) = 0$ for $i = 1, \dots, n$ and (Y_1, \dots, Y_N) are uniformly sub-gaussian: $\max_{i=1, \dots, n} c_1^2 \mathbb{E}\{\exp(Y_i^2/c_1^2) - 1\} \leq c_2^2$. Then for any $t > 0$,

$$\mathbb{P} \left(\left| \frac{1}{N} \sum_{i=1}^N Y_i \right| > t \right) \leq 2 \exp \left\{ -\frac{nt^2}{8(c_1^2 + c_2^2)} \right\}.$$

Lemma 3. Let (Y_1, \dots, Y_N) be independent variables such that $\mathbb{E}(Y_i) = 0$ for $i = 1, \dots, n$ and,

$$\frac{1}{N} \sum_{i=1}^N \mathbb{E}(|Y_i|^k) \leq \frac{k!}{2} c_3^{k-2} c_4^2, \quad k = 2, 3, \dots,$$

for some constants (c_3, c_4) . Then for any $t > 0$,

$$\mathbb{P} \left(\left| \frac{1}{N} \sum_{i=1}^N Y_i \right| > c_3 t + c_4 \sqrt{2t} \right) \leq 2 \exp(-nt).$$

Lemma 4. Suppose that $|X| \leq c_5$, c_5 is some constant, and Y is sub-gaussian: $c_1^2 \mathbb{E}\{\exp(X^2/c_1^2) - 1\} \leq c_2^2$ for some constants (c_1, c_2) . Then $Z = XY^2$ satisfies

$$\mathbb{E}\{|Z - \mathbb{E}(Z)|^k\} \leq \frac{k!}{2} c_6^{k-2} c_7^2, \quad k = 2, 3, \dots,$$

for $c_6 = 2c_5 c_1^2$ and $c_7 = 2c_5 c_1 c_2$.

Lemma 5. Suppose that Y is sub-gaussian: $c_1^2 \mathbb{E}\{\exp(Y^2/c_1^2) - 1\} \leq c_2^2$ for some constants (c_1, c_2) . Then

$$\mathbb{E}(|Y|^k) \leq \Gamma\left(\frac{k}{2} + 1\right) (c_1^2 + c_2^2) c_1^{k-2}, \quad k = 2, 3, \dots$$

S5 Technical lemmas

S5.1 Lemmas for the parameter in the PS model

The following Lemmas 6–8 will be used in proofs of Proposition 1 and Theorem 1.

Lemma 9 would be used in proofs of lemmas in Section S5.2 and Theorem 1.

Lemma 6. Under Assumptions 1(i) and 1(ii), the following statements hold:

(i) Denoted by Ω_{00} the event that

$$\sup_{j=0,\dots,p} \left| \tilde{\mathbb{E}} [\{-Rw(\mathbf{X}; \bar{\gamma}) + (1 - R)\} f_j(\mathbf{X})] \right| \leq \lambda_0.$$

If $\lambda_0 \geq \sqrt{2}\{1 + \exp(-B_0)\}C_0\sqrt{\ln\{(1+p)/\epsilon\}/N}$, then $\mathbb{P}(\Omega_{00}) \geq 1 - 2\epsilon$.

(ii) Denote by Ω_{01} the event that

$$\sup_{j,k=0,\dots,p} \left| (\tilde{\Sigma}_\gamma)_{jk} - (\Sigma_\gamma)_{jk} \right| \leq \lambda_0, \tag{S5.6}$$

where $\tilde{\Sigma}_\gamma$ is the empirical version of Σ_γ . If $\lambda_0 \geq 4 \exp(-B_0)C_0^2\sqrt{\ln\{(1+p)/\epsilon\}/N}$, then $\mathbb{P}(\Omega_{01}) \geq 1 - 2\epsilon^2$.

The result of Lemma 6(ii) is taken from Tan (2020b), Lemma 1(ii), the result of Lemma 6(i) can be shown similarly using Lemma 1 in Section S4 and the union bound.

Let $\Sigma_{m^1} = \mathbb{E}\{Rw(\mathbf{X}; \bar{\gamma})|Y - \psi(\bar{\alpha}^T \mathbf{G})|\mathbf{F}\mathbf{F}^T\}$, and $\tilde{\Sigma}_{m^1}$ be the sample version of Σ_{m^1} .

Lemma 7. Let Ω_{02} denote the event that

$$\sup_{j,k=0,\dots,p} |(\boldsymbol{\Sigma}_{m^1})_{jk} - (\tilde{\boldsymbol{\Sigma}}_{m^1})_{jk}| \leq \sqrt{D_0^2 + D_1^2} \lambda_0. \quad (\text{S5.7})$$

Under Assumptions 1(i), 1(ii) and 2(iv), if

$$\lambda_0 \geq 4 \exp(-B_0) C_0^2 \sqrt{\frac{\ln \left\{ \frac{(1+p)}{\epsilon} \right\}}{N}},$$

then $\mathbb{P}(\Omega_{02}) \geq 1 - 2\epsilon^2$.

Proof. Since $|Rw(\mathbf{X}; \bar{\gamma}) f_j(\mathbf{X}) f_k(\mathbf{X})| \leq \exp(-B_0) C_0^2$ for $j, k = 0, \dots, p$ by Assumptions 1(i), and 1(ii), and $|Y - \psi(\bar{\alpha}^\top \mathbf{G})|$ is uniformly sub-gaussian by Assumption 2(iv), $Rw(\mathbf{X}; \bar{\gamma}) |Y - \psi(\bar{\alpha}^\top \mathbf{G})| f_j(\mathbf{X}) f_k(\mathbf{X})$ is uniformly sub-gaussian. By Lemma 2, we have

$$\mathbb{P} \left\{ |(\boldsymbol{\Sigma}_{m^1})_{jk} - (\tilde{\boldsymbol{\Sigma}}_{m^1})_{jk}| > t \right\} \leq \frac{2\epsilon^2}{(1+p)^2},$$

for $j, k = 0, \dots, p$, where $t = \exp(-B_0) C_0^2 \sqrt{8(D_0^2 + D_1^2)} \sqrt{\ln\{(1+p)^2/\epsilon^2\}/N}$. By union bounds, (S5.7) holds. \square

Let $\boldsymbol{\Sigma}_{m^2} = \mathbb{E}[Rw(\mathbf{X}; \bar{\gamma}) \{Y - \psi(\bar{\alpha}^\top \mathbf{G})\}^2 \mathbf{F} \mathbf{F}^\top]$, the sample version of $\boldsymbol{\Sigma}_{m^2}$, $\tilde{\boldsymbol{\Sigma}}_{m^2} = \tilde{\mathbb{E}}[Rw(\mathbf{X}; \bar{\gamma}) \{Y - \psi(\bar{\alpha}^\top \mathbf{G})\}^2 \mathbf{F} \mathbf{F}^\top]$.

Lemma 8. Denote by Ω_{03} the event that

$$\sup_{j,k=0,\dots,p} |(\boldsymbol{\Sigma}_{m^2})_{jk} - (\tilde{\boldsymbol{\Sigma}}_{m^2})_{jk}| \leq (D_0^2 + D_0 D_1) \lambda_0, \quad (\text{S5.8})$$

Under Assumptions 1(i), 1(ii) and 2(iv), if

$$(D_0^2 + D_0 D_1) \lambda_0 \geq 4 \exp(-B_0) C_0^2 [D_0^2 \ln\{(1+p)/\epsilon\}/N + D_0 D_1 \sqrt{\ln\{(1+p)/\epsilon\}/N}],$$

then, $\mathbb{P}(\Omega_{03}) \geq 1 - 2\epsilon^2$. Furthermore, if we assume that $\ln\{(1+p)/\epsilon\}/N < 1$, the above condition reduces to $\lambda_0 \geq 4 \exp(-B_0) C_0^2 \sqrt{\ln\{(1+p)/\epsilon\}/N}$.

Proof. For $j, k = 0, \dots, p$, the variable $Rw(\mathbf{X}; \bar{\gamma}) \{Y - \psi(\bar{\boldsymbol{\alpha}}^T \mathbf{G})\}^2 f_j(\mathbf{X}) f_k(\mathbf{X})$ is the product of $Rw(\mathbf{X}; \bar{\gamma}) f_j(\mathbf{X}) f_k(\mathbf{X})$ and $\{Y - \psi(\bar{\boldsymbol{\alpha}}^T \mathbf{G})\}^2$, where $|Rw(\mathbf{X}; \bar{\gamma}) f_j(\mathbf{X}) f_k(\mathbf{X})| \leq \exp(-B_0) C_0^2$ by Assumptions 1(i) and 1(ii); and $\{Y - \psi(\bar{\boldsymbol{\alpha}}^T \mathbf{G})\}$ is sub-gaussian by Assumption 2(iv). By Lemmas 3 and 4 in Section S4, we have

$$\mathbb{P} \left\{ |(\tilde{\boldsymbol{\Sigma}}_{m^2})_{j,k} - \boldsymbol{\Sigma}_{m^2}_{j,k}| > 2e^{(-B_0)} C_0^2 D_0^2 t + 2e^{(-B_0)} C_0^2 D_0 D_1 \sqrt{2t} \right\} \leq 2 \frac{\epsilon^2}{(1+p)^2},$$

for $j, k = 0, \dots, p$, where $t = \ln\{(1+p)^2/\epsilon^2\}/N$. The result then follows from the union bound. \square

Lemma 9. Under the conditions of Proposition 1, in the event $\Omega_{00} \cap \Omega_{01}$, we have

$$\tilde{\mathbb{E}} [Rw(\mathbf{X}; \bar{\gamma}) \{(\hat{\boldsymbol{\gamma}} - \bar{\boldsymbol{\gamma}})^T \mathbf{F}\}^2] \leq \exp(\eta_{01}) M_0 |S_{\bar{\gamma}}| \lambda_0^2, \quad (\text{S5.9})$$

where $\eta_{01} = (A_0 - 1)^{-1}M_0C_0\zeta_0$.

Proof. By Lemma 6 and Proposition 1, in the event $\Omega_{00} \cap \Omega_{01}$, (4.1) holds, we obtain

$$\|\hat{\gamma} - \bar{\gamma}\|_1 \leq (A_0 - 1)^{-1}M_0|S_{\bar{\gamma}}|\lambda_0 \leq (A_0 - 1)^{-1}M_0\zeta_0, \quad (\text{S5.10})$$

the second inequality holds due to Assumption 1(iv). By the definition of $D_{\text{CAL}}^\dagger(\cdot, \cdot)$, we obtain

$$\begin{aligned} D_{\text{CAL}}^\dagger(\hat{\gamma}^\top \mathbf{F}, \bar{\gamma}^\top \mathbf{F}) &= -\tilde{\mathbb{E}}[R\{\exp(-\hat{\gamma}^\top \mathbf{F}) - \exp(-\bar{\gamma}^\top \mathbf{F})\}\{(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}] \\ &= \tilde{\mathbb{E}}(Rw(\mathbf{X}; \bar{\gamma}) \exp[-u\{(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}]\{(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}^2) \\ &\geq \exp(-C_0\|\hat{\gamma} - \bar{\gamma}\|_1)\tilde{\mathbb{E}}[Rw(\mathbf{X}; \bar{\gamma})\{(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}^2] \\ &\geq \exp(-\eta_{01})\tilde{\mathbb{E}}[Rw(\mathbf{X}; \bar{\gamma})\{(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}^2]. \end{aligned}$$

The second equality holds by the mean value theorem and u is some scalar in $(0, 1)$.

Combining the inequality with (4.1), we obtain

$$\tilde{\mathbb{E}}[Rw(\mathbf{X}; \bar{\gamma})\{(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}^2] \leq \exp(\eta_{01})D_{\text{CAL}}^\dagger(\hat{\gamma}^\top \mathbf{F}, \bar{\gamma}^\top \mathbf{F}) \leq \exp(\eta_{01})M_0|S_{\bar{\gamma}}|\lambda_0^2.$$

□

S5.2 Lemmas for the parameter in the OR model

The following Lemmas 10–17 will be used in the proofs of Proposition 2 and Theorem 1.

Lemma 10. Let Ω_{10} denote the event that

$$\sup_{j=0,\dots,q} \left| \tilde{\mathbb{E}} [Rw(\mathbf{X}; \bar{\gamma}) \{Y - \psi(\bar{\alpha}^\top \mathbf{G})\} g_j(\mathbf{X})] \right| \leq \lambda_1. \quad (\text{S5.11})$$

Under Assumptions 1, 2(iii) and 2(iv), if

$$\lambda_1 \geq \exp(-B_0) C_4 \sqrt{8(D_0^2 + D_1^2)} \sqrt{\ln\{(1+q)/\epsilon\}/N},$$

then $\mathbb{P}(\Omega_{10}) \geq 1 - 2\epsilon$.

Proof. Let $S_j = Rw(\mathbf{X}; \bar{\gamma}) \{Y - \psi(\bar{\alpha}^\top \mathbf{G})\} g_j(\mathbf{X})$ for $j = 0, \dots, q$. Then, $\mathbb{E}(S_j) = 0$ by the definition of $\bar{\alpha}$. Under Assumptions 1 and 2(iii), $|S_j| \leq \exp(-B_0) C_4 |R\{Y - \psi(\bar{\alpha}^\top \mathbf{G})\}|$. By Assumption 2(iv), the variables (S_0, \dots, S_q) are uniformly sub-gaussian: $\max_{j=0,\dots,q} D_2^2 \mathbb{E}\{\exp(S_j^2/D_2^2) - 1\} \leq D_3^2$, with $D_2 = \exp(-B_0) C_4 D_0$, and $D_3 = e^{(-B_0)} C_4 D_1$. Therefore, by Lemma 2 in Section S4 and the union bound, $\mathbb{P}(\Omega_{10}) \geq 1 - 2\epsilon$, if $\lambda_1 \geq \exp(-B_0) C_4 \sqrt{8(D_0^2 + D_1^2)} \sqrt{\ln\{(1+q)/\epsilon\}/N}$. \square

Lemma 11. Denote by Ω_{11} the event that

$$\sup_{j,k=0,\dots,q} |(\boldsymbol{\Sigma}_{\bar{\alpha}})_{jk} - (\tilde{\boldsymbol{\Sigma}}_{\bar{\alpha}})_{jk}| \leq \lambda_1, \quad (\text{S5.12})$$

where $\tilde{\boldsymbol{\Sigma}}_{\bar{\alpha}}$ is the empirical version of $\boldsymbol{\Sigma}_{\bar{\alpha}}$. Under Assumptions 1(ii), 2(i) and 2(iii), if $\lambda_1 \geq 4 \exp(-B_0) C_2 C_4^2 \sqrt{\ln\{(1+q)/\epsilon\}/N}$, then $\mathbb{P}(\Omega_{11}) \geq 1 - 2\epsilon^2$.

Proof. Notice that $|Rw(\mathbf{X}; \bar{\gamma}) \psi_1(\bar{\alpha}^\top \mathbf{G}) g_j(\mathbf{X}) g_k(\mathbf{X})| \leq e^{(-B_0)} C_2 C_4^2$ for $j, k = 0, 1, \dots, q$,

by Assumptions 1(ii), 2(i), and 2(iii). Thus,

$$|Rw(\mathbf{X}; \bar{\gamma})\psi_1(\bar{\alpha}^T \mathbf{G})g_j(\mathbf{X})g_k(\mathbf{X}) - \mathbb{E}\{Rw(\mathbf{X}; \bar{\gamma})\psi_1(\bar{\alpha}^T \mathbf{G})g_j(\mathbf{X})g_k(\mathbf{X})\}| \leq 2 \exp(-B_0)C_2C_4^2.$$

By union bounds and applying Lemma 1 yields (S5.12). \square

Lemma 12. For any $\alpha \in \mathbb{R}^{q+1}$, we have

$$D_{\text{WL}}^\dagger(\hat{\alpha}^T \mathbf{G}, \alpha^T \mathbf{G}; \hat{\gamma}) + \lambda_\alpha \|\hat{\alpha}_{1:q}\|_1 \leq (\hat{\alpha} - \alpha)^T \tilde{\mathbb{E}}[Rw(\mathbf{X}; \hat{\gamma})\{Y - \psi(\alpha^T \mathbf{G})\}\mathbf{G}] + \lambda_\alpha \|\alpha_{1:q}\|_1. \quad (\text{S5.13})$$

Proof. For any $u \in (0, 1]$, by definition of $\hat{\alpha}$, we have

$$\ell_{\text{WL}}(\hat{\alpha}, \hat{\gamma}) + \lambda_\alpha \|\hat{\alpha}_{1:q}\|_1 \leq \ell_{\text{WL}}\{(1-u)\hat{\alpha} + u\alpha; \hat{\gamma}\} + \lambda_\alpha \|(1-u)\hat{\alpha}_{1:q} + u\alpha_{1:q}\|_1,$$

which implies

$$\ell_{\text{WL}}(\hat{\alpha}, \hat{\gamma}) - \ell_{\text{WL}}\{(1-u)\hat{\alpha} + u\alpha; \hat{\gamma}\} + \lambda_\alpha u \|\hat{\alpha}_{1:q}\|_1 \leq \lambda_\alpha u \|\alpha_{1:q}\|_1,$$

by the convexity of L_1 -norm. Dividing both sides of the preceding inequality by u

and letting $u \rightarrow 0_+$ leads to

$$-\tilde{\mathbb{E}}[Rw(\mathbf{X}; \hat{\gamma})\{Y - \psi(\hat{\alpha}^T \mathbf{G})\}\{(\hat{\alpha} - \alpha)^T \mathbf{G}\}] + \lambda_\alpha \|\hat{\alpha}_{1:q}\|_1 \leq \lambda_\alpha \|\alpha_{1:q}\|_1,$$

which yields (S5.13) after rearranging using (4.3). \square

Lemma 13. For any function $h(\mathbf{X})$, under the conditions of Proposition 1, in the event $\Omega_{00} \cap \Omega_{01}$,

$$D_{\text{WL}}^\dagger\{\hat{\alpha}^\text{T}\mathbf{G}, h(\mathbf{X}); \hat{\gamma}\} \geq \exp(-\eta_{01})D_{\text{WL}}^\dagger\{\hat{\alpha}^\text{T}\mathbf{G}, h(\mathbf{X}); \bar{\gamma}\}. \quad (\text{S5.14})$$

Proof. By the definition of $D_{\text{WL}}^\dagger(\cdot, \cdot; \cdot)$,

$$\begin{aligned} & D_{\text{WL}}^\dagger\{\hat{\alpha}^\text{T}\mathbf{G}, h(\mathbf{X}); \hat{\gamma}\} \\ &= \tilde{\mathbb{E}}(Rw(\mathbf{X}; \hat{\gamma})[\psi(\hat{\alpha}^\text{T}\mathbf{G}) - \psi\{h(\mathbf{X})\}]\{\hat{\alpha}^\text{T}\mathbf{G} - h(\mathbf{X})\}) \\ &= \tilde{\mathbb{E}}(Rw(\mathbf{X}; \bar{\gamma}) \exp\{-(\hat{\gamma} - \bar{\gamma})^\text{T}\mathbf{F}\}[\psi(\hat{\alpha}^\text{T}\mathbf{G}) - \psi\{h(\mathbf{X})\}]\{\hat{\alpha}^\text{T}\mathbf{G} - h(\mathbf{X})\}) \\ &\geq \tilde{\mathbb{E}}(Rw(\mathbf{X}; \bar{\gamma}) \exp(-\eta_{01})[\psi(\hat{\alpha}^\text{T}\mathbf{G}) - \psi\{h(\mathbf{X})\}]\{\hat{\alpha}^\text{T}\mathbf{G} - h(\mathbf{X})\}) \\ &= \exp(-\eta_{01})\tilde{\mathbb{E}}(Rw(\mathbf{X}; \bar{\gamma})[\psi(\hat{\alpha}^\text{T}\mathbf{G}) - \psi\{h(\mathbf{X})\}]\{\hat{\alpha}^\text{T}\mathbf{G} - h(\mathbf{X})\}) \\ &= \exp(-\eta_{01})D_{\text{WL}}^\dagger\{\hat{\alpha}^\text{T}\mathbf{G}, h(\mathbf{X}); \bar{\gamma}\}. \end{aligned}$$

The inequality holds since in the event $\Omega_{00} \cap \Omega_{01}$, (4.1) holds. \square

For functions $h(\mathbf{X})$ and $h'(\mathbf{X})$, let $Q_{\text{WL}}\{h(\mathbf{X}), h'(\mathbf{X}); \cdot\} = \tilde{\mathbb{E}}[Rw(\mathbf{X}; \cdot)\{h(\mathbf{X}) - h'(\mathbf{X})\}^2]$.

Lemma 14. Suppose Assumption 2(iv) holds, in the event $\Omega_{00} \cap \Omega_{01} \cap \Omega_{03}$, (S5.13)

implies

$$\begin{aligned}
& \exp(-\eta_{01})D_{\text{WL}}^\dagger(\hat{\alpha}^\text{T}\mathbf{G}, \bar{\alpha}^\text{T}\mathbf{G}; \bar{\gamma}) + \lambda_\alpha \|\hat{\alpha}_{1:q}\|_1 \\
& \leq (\hat{\alpha} - \bar{\alpha})^\text{T} \tilde{\mathbb{E}} [Rw(\mathbf{X}; \bar{\gamma})\{Y - \psi(\bar{\alpha}^\text{T}\mathbf{G})\}\mathbf{G}] \\
& \quad + \lambda_\alpha \|\bar{\alpha}_{1:q}\|_1 + \exp(\eta_{01})\sqrt{M_1|S_{\bar{\gamma}}|\lambda_0^2\{Q_{\text{WL}}(\hat{\alpha}^\text{T}\mathbf{G}, \bar{\alpha}^\text{T}\mathbf{G}; \bar{\gamma})\}^{1/2}},
\end{aligned}$$

where $M_1 = \left[(2D_0^2 + D_1^2 + D_0D_1) \left\{ \frac{M_0^2\zeta_0}{(A_0-1)^2} \right\} + (D_0^2 + D_1^2) \exp(\eta_{01})M_0 \right]$.

Proof. Consider the following decomposition,

$$\begin{aligned}
& (\hat{\alpha} - \bar{\alpha})^\text{T} \tilde{\mathbb{E}} [Rw(\mathbf{X}; \hat{\gamma})\{Y - \psi(\bar{\alpha}^\text{T}\mathbf{G})\}\mathbf{G}] \\
& = (\hat{\alpha} - \bar{\alpha})^\text{T} \tilde{\mathbb{E}} [Rw(\mathbf{X}; \bar{\gamma})\{Y - \psi(\bar{\alpha}^\text{T}\mathbf{G})\}\mathbf{G}] \\
& \quad + \tilde{\mathbb{E}} [R\{w(\mathbf{X}; \hat{\gamma}) - w(\mathbf{X}; \bar{\gamma})\}\{Y - \psi(\bar{\alpha}^\text{T}\mathbf{G})\}\{(\hat{\alpha} - \bar{\alpha})^\text{T}\mathbf{G}\}],
\end{aligned} \tag{S5.15}$$

denoted as $\Delta_0^0 + \Delta_1^0$. By the mean value theorem and Cauchy-Schwartz inequality,

$$\begin{aligned}
\Delta_1^0 & \leq \exp(C_0\|\hat{\gamma} - \bar{\gamma}\|_1)\tilde{\mathbb{E}}^{1/2}[Rw(\mathbf{X}; \bar{\gamma})\{(\hat{\alpha} - \bar{\alpha})^\text{T}\mathbf{G}\}^2] \\
& \quad \times \tilde{\mathbb{E}}^{1/2}[Rw(\mathbf{X}; \bar{\gamma})\{Y - \psi(\bar{\alpha}^\text{T}\mathbf{G})\}^2\{(\hat{\gamma} - \bar{\gamma})^\text{T}\mathbf{F}\}^2] \\
& \leq \exp(\eta_{01})\{Q_{\text{WL}}(\hat{\alpha}^\text{T}\mathbf{G}, \bar{\alpha}^\text{T}\mathbf{G}; \bar{\gamma})\}^{1/2} \\
& \quad \times \tilde{\mathbb{E}}^{1/2}[Rw(\mathbf{X}; \bar{\gamma})\{Y - \psi(\bar{\alpha}^\text{T}\mathbf{G})\}^2\{(\hat{\gamma} - \bar{\gamma})^\text{T}\mathbf{F}\}^2].
\end{aligned} \tag{S5.16}$$

We bound the third term in (S5.16). By Assumption 2(iv) and Lemma 5, we have

$$\mathbb{E}[\{Y - \psi(\bar{\alpha}^\text{T}\mathbf{G})\}^2|\mathbf{X}] \leq D_0^2 + D_1^2.$$

Therefore,

$$\mathbb{E}[Rw(\mathbf{X}; \bar{\gamma})\{Y - \psi(\bar{\alpha}^T \mathbf{G})\}^2\{(\hat{\gamma} - \bar{\gamma})^T \mathbf{F}\}^2] \leq (D_0^2 + D_1^2)\mathbb{E}[Rw(\mathbf{X}; \bar{\gamma})\{(\hat{\gamma} - \bar{\gamma})^T \mathbf{F}\}^2].$$

Let $(\tilde{\mathbb{E}} - \mathbb{E})(\mathbf{U})$ denote $\tilde{\mathbb{E}}\{\mathbf{U} - \mathbb{E}(\mathbf{U})\}$ for \mathbf{U} , which is a function of (\mathbf{X}, R, Y) . Then in the event Ω_{01} , by (S5.6), we have

$$(\mathbb{E} - \tilde{\mathbb{E}})[Rw(\mathbf{X}; \bar{\gamma})\{(\hat{\gamma} - \bar{\gamma})^T \mathbf{F}\}^2] \leq \lambda_0 \|\hat{\gamma} - \bar{\gamma}\|_1^2.$$

In the event Ω_{03} , by (S5.12), we have

$$(\tilde{\mathbb{E}} - \mathbb{E})[Rw(\mathbf{X}; \bar{\gamma})\{Y - \psi(\bar{\alpha}^T \mathbf{G})\}^2\{(\hat{\gamma} - \bar{\gamma})^T \mathbf{F}\}^2] \leq (D_0^2 + D_0 D_1)\lambda_0 \|\hat{\gamma} - \bar{\gamma}\|_1^2.$$

Combining preceding inequalities, we obtain in $\Omega_{00} \cap \Omega_{01} \cap \Omega_{03}$,

$$\begin{aligned} & \tilde{\mathbb{E}}[Rw(\mathbf{X}; \bar{\gamma})\{Y - \psi(\bar{\alpha}^T \mathbf{G})\}^2\{(\hat{\gamma} - \bar{\gamma})^T \mathbf{F}\}^2] \\ & \leq (D_0^2 + D_0 D_1)\lambda_0 \|\hat{\gamma} - \bar{\gamma}\|_1^2 + (D_0^2 + D_1^2) \left\{ \lambda_0 \|\hat{\gamma} - \bar{\gamma}\|_1^2 + \tilde{\mathbb{E}}[Rw(\mathbf{X}; \bar{\gamma})\{(\hat{\gamma} - \bar{\gamma})^T \mathbf{F}\}^2] \right\} \\ & \leq (2D_0^2 + D_1^2 + D_0 D_1) \left\{ \frac{M_0^2 \zeta_0}{(A_0 - 1)^2} \right\} |S_{\bar{\gamma}}| \lambda_0^2 + (D_0^2 + D_1^2) \exp(\eta_{01}) M_0 |S_{\bar{\gamma}}| \lambda_0^2 \\ & = M_1 |S_{\bar{\gamma}}| \lambda_0^2, \end{aligned} \tag{S5.17}$$

where the last inequality holds due to (4.1) and (S5.9). Combining (S5.15)–(S5.17), we obtain

$$\begin{aligned}
& (\hat{\alpha} - \bar{\alpha})^\top \tilde{\mathbb{E}} [Rw(\mathbf{X}; \hat{\gamma}) \{Y - \psi(\bar{\alpha}^\top \mathbf{G})\} \mathbf{G}] \\
& \leq (\hat{\alpha} - \bar{\alpha})^\top \tilde{\mathbb{E}} [Rw(\mathbf{X}; \bar{\gamma}) \{Y - \psi(\bar{\alpha}^\top \mathbf{G})\} \mathbf{G}] \\
& \quad + \exp(\eta_{01}) (M_1 |S_{\bar{\gamma}}| \lambda_0^2)^{1/2} \{Q_{\text{WL}}(\hat{\alpha}^\top \mathbf{G}, \bar{\alpha}^\top \mathbf{G}; \bar{\gamma})\}^{1/2}.
\end{aligned} \tag{S5.18}$$

The desired result follows by combining (S5.13), (S5.14) and (S5.18) in the event $\Omega_{00} \cap \Omega_{01} \cap \Omega_{03}$. \square

Lemma 15. Denote $b = \hat{\alpha} - \bar{\alpha}$. Suppose Assumption 2(iv) holds. In the event, $\Omega_{00} \cap \Omega_{01} \cap \Omega_{03} \cap \Omega_{10}$, we have

$$\begin{aligned}
& \exp(-\eta_{01}) D_{\text{WL}}^\dagger(\hat{\alpha}^\top \mathbf{G}, \bar{\alpha}^\top \mathbf{G}; \bar{\gamma}) + (A_1 - 1) \lambda_1 \|b\|_1 \\
& \leq 2A_1 \lambda_1 \sum_{j \in S_{\bar{\alpha}}} |b_j| + \exp(\eta_{01}) \sqrt{M_1 |S_{\bar{\gamma}}| \lambda_0^2} \{Q_{\text{WL}}(\hat{\alpha}^\top \mathbf{G}, \bar{\alpha}^\top \mathbf{G}; \bar{\gamma})\}^{1/2}.
\end{aligned} \tag{S5.19}$$

Proof. In the event Ω_{10} , we have

$$b^\top \tilde{\mathbb{E}} [Rw(\mathbf{X}; \bar{\gamma}) \{Y - \psi(\bar{\alpha}^\top \mathbf{G})\} \mathbf{G}] \leq \lambda_1 \|b\|_1,$$

by which and Lemma 14, we have in the event $\Omega_{00} \cap \Omega_{01} \cap \Omega_{03} \cap \Omega_{10}$,

$$\begin{aligned}
& \exp(-\eta_{01}) D_{\text{WL}}^\dagger(\hat{\alpha}^\top \mathbf{G}, \bar{\alpha}^\top \mathbf{G}; \bar{\gamma}) + A_1 \lambda_1 \|\hat{\alpha}_{1:q}\|_1 \\
& \leq \lambda_1 \|b\|_1 + A_1 \lambda_1 \|\bar{\alpha}_{1:q}\|_1 + \exp(\eta_{01}) \sqrt{M_1 |S_{\bar{\gamma}}| \lambda_0^2} \{Q_{\text{WL}}(\hat{\alpha}^\top \mathbf{G}, \bar{\alpha}^\top \mathbf{G}; \bar{\gamma})\}^{1/2}.
\end{aligned}$$

Applying to the preceding inequality the identity $\hat{\alpha}_j = |\hat{\alpha}_j - \bar{\alpha}_j|$ for $j \notin S_{\bar{\alpha}}$ and the triangle inequality

$$|\hat{\alpha}_j| \geq |\bar{\alpha}_j| - |\hat{\alpha}_j - \bar{\alpha}_j|, \quad j \in S_{\bar{\alpha}} \setminus \{0\},$$

and rearranging the result gives

$$\begin{aligned} & \exp(-\eta_{01})D_{\text{WL}}^\dagger(\hat{\alpha}^\text{T}\mathbf{G}, \bar{\alpha}^\text{T}\mathbf{G}; \bar{\gamma}) + (A_1 - 1)\lambda_1\|b_{1:q}\|_1 \\ & \leq \lambda_1|b_0| + 2A_1\lambda_1 \sum_{j \in S_{\bar{\alpha}} \setminus \{0\}} |b_j| + \exp(\eta_{01})\sqrt{M_1|S_{\bar{\gamma}}|\lambda_0^2}\{Q_{\text{WL}}(\hat{\alpha}^\text{T}\mathbf{G}, \bar{\alpha}^\text{T}\mathbf{G}; \bar{\gamma})\}^{1/2} \\ & \leq (A_1 + 1)\lambda_1|b_0| + 2A_1\lambda_1 \sum_{j \in S_{\bar{\alpha}} \setminus \{0\}} |b_j| + \exp(\eta_{01})\sqrt{M_1|S_{\bar{\gamma}}|\lambda_0^2}\{Q_{\text{WL}}(\hat{\alpha}^\text{T}\mathbf{G}, \bar{\alpha}^\text{T}\mathbf{G}; \bar{\gamma})\}^{1/2}. \end{aligned}$$

By adding $(A_1 - 1)\lambda_1|b_0|$ on both sides of the previous inequality, the conclusion follows. \square

Lemma 16. Suppose Assumptions 2(ii) and 2(iii) hold. Then, for any $\alpha, \tilde{\alpha} \in \mathbb{R}^{q+1}$,

$$D_{\text{WL}}^\dagger(\alpha^\text{T}\mathbf{G}, \tilde{\alpha}^\text{T}\mathbf{G}; \bar{\gamma}) \geq \{b^\text{T}\tilde{\Sigma}_\alpha(\tilde{\alpha})b\} \frac{1 - \exp(-C_{40}\|b\|_1)}{C_{40}\|b\|_1},$$

where $b = \alpha - \tilde{\alpha}$, $C_{40} = C_3C_4$ and $\tilde{\Sigma}_\alpha(\cdot) = \tilde{\mathbb{E}}[Rw(\mathbf{X}; \bar{\gamma})\psi_1\{(\cdot)^\text{T}\mathbf{G}\}\mathbf{G}\mathbf{G}^\text{T}]$. Throughout, set $\{1 - \exp(-c)\}/c = 1$, for $c = 0$.

Proof. By the definition of $D_{\text{WL}}^\dagger(\cdot, \cdot; \cdot)$,

$$D_{\text{WL}}^\dagger(\alpha^\text{T}\mathbf{G}, \tilde{\alpha}^\text{T}\mathbf{G}; \bar{\gamma})$$

$$\begin{aligned}
&= \tilde{\mathbb{E}} [Rw(\mathbf{X}; \bar{\gamma}) \{\psi(\alpha^\top \mathbf{G}) - \psi(\tilde{\alpha}^\top \mathbf{G})\} (\alpha^\top \mathbf{G} - \tilde{\alpha}^\top \mathbf{G})] \\
&= \tilde{\mathbb{E}} \left(Rw(\mathbf{X}; \bar{\gamma}) \left[\int_0^1 \psi_1 \{\tilde{\alpha}^\top \mathbf{G} + u(\alpha^\top \mathbf{G} - \tilde{\alpha}^\top \mathbf{G})\} du \right] \{\alpha^\top \mathbf{G} - \tilde{\alpha}^\top \mathbf{G}\}^2 \right).
\end{aligned}$$

By Assumptions 2(ii), 2(iii), and the fact that $|\alpha^\top \mathbf{G} - \tilde{\alpha}^\top \mathbf{G}| \leq \{\sup_{j=0, \dots, q} |g_j(\mathbf{X})|\} \|\alpha - \tilde{\alpha}\|_1 \leq C_4 \|\alpha - \tilde{\alpha}\|_1$, it follows that

$$\begin{aligned}
&D_{\text{WL}}^\dagger(\alpha^\top \mathbf{G}, \tilde{\alpha}^\top \mathbf{G}; \bar{\gamma}) \\
&\geq \tilde{\mathbb{E}} \left[Rw(\mathbf{X}; \bar{\gamma}) \left\{ \int_0^1 \psi_1(\tilde{\alpha}^\top \mathbf{G}) \exp(-uC_3|\alpha^\top \mathbf{G} - \tilde{\alpha}^\top \mathbf{G}|) du \right\} (\alpha^\top \mathbf{G} - \tilde{\alpha}^\top \mathbf{G})^2 \right] \\
&\geq \tilde{\mathbb{E}} \{Rw(\mathbf{X}; \bar{\gamma}) \psi_1(\tilde{\alpha}^\top \mathbf{G}) (\alpha^\top \mathbf{G} - \tilde{\alpha}^\top \mathbf{G})^2\} \left\{ \int_0^1 \exp(-uC_{40}\|\tilde{\alpha} - \alpha\|_1) du \right\},
\end{aligned}$$

which gives the desired result since $\int_0^1 \exp(-cu) du = \{1 - \exp(-c)\}/c$ for $c \geq 0$. \square

Lemma 17. Suppose that Assumption 2(v) holds. In the event Ω_{11} , Assumption 2(vi) implies a compatibility condition for $\tilde{\Sigma}_{\tilde{\alpha}}$: for any vector $b = (b_0, \dots, b_q)^\top$ such that $\sum_{j \notin S_{\tilde{\alpha}}} |b_j| \leq \xi_1 \sum_{j \in S_{\tilde{\alpha}}} |b_j|$, we have

$$(1 - \zeta_1) \nu_1^2 \left(\sum_{j \in S_{\tilde{\alpha}}} |b_j| \right)^2 \leq |S_{\tilde{\alpha}}| (b^\top \tilde{\Sigma}_{\tilde{\alpha}} b). \quad (\text{S5.20})$$

Proof. In the event Ω_{11} , we have $|b^\top (\tilde{\Sigma}_{\tilde{\alpha}} - \Sigma_{\tilde{\alpha}}) b| \leq \lambda_1 \|b\|_1^2$. Then Assumption 2(v) implies that for any $b = (b_0, \dots, b_q)^\top$ satisfying $\sum_{j \notin S_{\tilde{\alpha}}} |b_j| \leq \xi_1 \sum_{j \in S_{\tilde{\alpha}}} |b_j|$,

$$\begin{aligned}
\nu_1^2 \|b_{S_{\tilde{\alpha}}}\|_1^2 &\leq |S_{\tilde{\alpha}}| (b^\top \Sigma_{\tilde{\alpha}} b) \leq |S_{\tilde{\alpha}}| (b^\top \tilde{\Sigma}_{\tilde{\alpha}} b + \lambda_1 \|b\|_1^2) \\
&\leq |S_{\tilde{\alpha}}| (b^\top \tilde{\Sigma}_{\tilde{\alpha}} b) + |S_{\tilde{\alpha}}| \lambda_1 (1 + \xi_1)^2 \|b_{S_{\tilde{\alpha}}}\|_1^2
\end{aligned}$$

$$\leq |S_{\bar{\alpha}}|(b^T \tilde{\Sigma}_{\bar{\alpha}} b) + \zeta_1 \nu_1^2 \|b_{S_{\bar{\alpha}}}\|_1^2$$

where $\|b_{S_{\bar{\alpha}}}\|_1 = \sum_{j \in S_{\bar{\alpha}}} |b_j|$; and the last inequality holds due to Assumption 2(vi), $(1 + \xi_1)^2 \nu_1^{-2} |S_{\bar{\alpha}}| \lambda_1 \leq \zeta_1$. Thus (S5.20) follows by rearrangement. \square

S6 Proofs of Propositions 1 and 2

S6.1 Proof of Proposition 1

Proof. Let c_γ in Assumption 2(iv) $= \max[\sqrt{2}\{1 + \exp(-B_0)\}C_0, 4 \exp(-B_0)C_0^2]$. This can be shown similarly to Tan (2020a), Theorem 1 by Lemmas 6–8 and Lemmas similar to 10–17. The small difference in probability is due to extra constraints of λ_0 on Ω_{02} and Ω_{03} from the sequential estimate, which is also demonstrated in Tan (2020a), Theorem 5. \square

S6.2 Proof of Proposition 2

Proof. To facilitate the proof, we first define some constants. Let $\nu_2 = \nu_1(1 - \zeta_1)^{1/2}$, $\xi_2 = 1 - 2A_1/\{(\xi_1 + 1)(A_1 - 1)\}$, $\xi_3 = (\xi_1 + 1)(A_1 - 1)$, c_α in Assumption 2(vi) equals to $\max\{C_4\sqrt{8(D_0^2 + D_1^2)}, 4C_2C_4^2\} \exp(-B_0)$, \tilde{c}_α and \tilde{c}_γ in Assumption 2(vii) equal to $C_3C_4 \exp(\eta_{01})(A_1 - 1)^{-1} \xi_3^2 \nu_2^{-2}$ and $C_1^{-1}C_3C_4(A_1 - 1)^{-1} \xi_2^{-2} \exp(3\eta_{01})M_1$, respectively.

Denote $b = \hat{\alpha} - \bar{\alpha}$, $D_{\text{WL}}^\dagger = D_{\text{WL}}^\dagger(\hat{\alpha}^T \mathbf{G}, \bar{\alpha}^T \mathbf{G}; \bar{\gamma})$, $Q_{\text{WL}} = Q_{\text{WL}}(\hat{\alpha}^T \mathbf{G}, \bar{\alpha}^T \mathbf{G}; \bar{\gamma})$ and $D_{\text{WL}}^\dagger = \exp(-\eta_{01})D_{\text{WL}}^\dagger + (A_1 - 1)\lambda_1 \|b\|_1$.

In the event $\Omega_{00} \cap \Omega_{01} \cap \Omega_{03} \cap \Omega_{10} \cap \Omega_{11}$, (S5.19) in Lemma 15 leads to two possible cases: either

$$\xi_2 D_{\text{WL}}^\ddagger \leq \exp(\eta_{01})(M_1 |S_{\bar{\gamma}}| \lambda_0^2 Q_{\text{WL}})^{1/2}, \quad (\text{S6.21})$$

or $(1 - \xi_2) D_{\text{WL}}^\ddagger \leq 2A_1 \lambda_1 \sum_{j \in S_{\bar{\alpha}}} |b_j|$, i.e.,

$$D_{\text{WL}}^\ddagger \leq (\xi_1 + 1)(A_1 - 1) \lambda_1 \sum_{j \in S_{\bar{\alpha}}} |b_j| = \xi_3 \lambda_1 \sum_{j \in S_{\bar{\alpha}}} |b_j|. \quad (\text{S6.22})$$

By Lemma 16, we have

$$D_{\text{WL}}^\dagger \geq \{b^\top \tilde{\Sigma}_\alpha(\bar{\alpha})b\} \frac{1 - \exp(-C_{40}\|b\|_1)}{C_{40}\|b\|_1} = (b^\top \tilde{\Sigma}_{\bar{\alpha}}b) \frac{1 - \exp(-C_{40}\|b\|_1)}{C_{40}\|b\|_1}. \quad (\text{S6.23})$$

If (S6.21) holds, notice that $D_{\text{WL}}^\dagger \leq \exp(\eta_{01}) D_{\text{WL}}^\ddagger$ and by Assumption 2(i), $Q_{\text{WL}} < C_1^{-1}(b^\top \tilde{\Sigma}_{\bar{\alpha}}b)$, which together with (S6.23) yields

$$D_{\text{WL}}^\dagger \leq \exp(3\eta_{01}) \xi_2^{-2} C_1^{-1}(M_1 |S_{\bar{\gamma}}| \lambda_0^2) \frac{C_{40}\|b\|_1}{1 - \exp(-C_{40}\|b\|_1)}. \quad (\text{S6.24})$$

Since $(A_1 - 1) \lambda_1 \|b\|_1 \leq D_{\text{WL}}^\ddagger$ and Assumption 2(vii) holds, (S6.24) implies that

$$1 - \exp(-C_{40}\|b\|_1) \leq (A_1 - 1)^{-1} \exp(3\eta_{01}) \xi_2^{-2} C_1^{-1}(M_1 |S_{\bar{\gamma}}| \lambda_0) C_{40} \leq \eta_3 < 1.$$

As a result, $C_{40}\|b\|_1 \leq -\ln(1 - \eta_3)$, which leads to

$$\frac{1 - \exp(-C_{40}\|b\|_1)}{C_{40}\|b\|_1} = \int_0^1 \exp(-C_{40}\|b\|_1 u) du \geq \exp(-C_{40}\|b\|_1) \geq 1 - \eta_3. \quad (\text{S6.25})$$

Combining the inequality with (S6.24), we obtain

$$D_{\text{WL}}^\dagger \leq \exp(3\eta_{01})\xi_2^{-2}\{C_1(1-\eta_3)\}^{-1}(M_1|S_{\bar{\gamma}}|\lambda_0^2)$$

If (S6.22) holds, then $\sum_{j \notin S_{\bar{\alpha}}} |b_j| \leq \xi_1 \sum_{j \in S_{\bar{\alpha}}} |b_j|$, which, together with Assumptions 2(v)–2(vi), implies (S5.20) in Lemma 17, that is,

$$\sum_{j \in S_{\bar{\alpha}}} |b_j| \leq (1 - \zeta_1)^{-1/2} \nu_1^{-1} |S_{\bar{\alpha}}|^{1/2} \left(b^\top \tilde{\Sigma}_{\bar{\alpha}} b \right)^{1/2}. \quad (\text{S6.26})$$

Since $D_{\text{WL}}^\dagger \leq \exp(\eta_{01}) D_{\text{WL}}^\dagger$, combining (S6.22), (S6.23) and (S6.26) yields

$$D_{\text{WL}}^\dagger \leq \exp(\eta_{01}) \xi_3^2 (1 - \zeta_1)^{-1} \nu_1^{-2} |S_{\bar{\alpha}}| \lambda_1^2 \frac{C_{40} \|b\|_1}{1 - \exp(-C_{40} \|b\|_1)}. \quad (\text{S6.27})$$

Since $(A_1 - 1)\lambda_1 \|b\|_1 \leq D_{\text{WL}}^\dagger$ and Assumption 2(ii) holds, (S6.27) implies that

$$1 - \exp(-C_{40} \|b\|_1) \leq (A_1 - 1)^{-1} \exp(\eta_{01}) \xi_3^2 (1 - \zeta_1)^{-1} \nu_1^{-2} |S_{\bar{\alpha}}| \lambda_1 C_{40} \leq \eta_2 < 1.$$

As a result, $C_{40} \|b\|_1 \leq -\ln(1 - \eta_2)$, which leads to

$$\frac{1 - \exp(-C_{40} \|b\|_1)}{C_{40} \|b\|_1} = \int_0^1 \exp(-C_{40} \|b\|_1 u) du \geq \exp(-C_{40} \|b\|_1) \geq 1 - \eta_2.$$

Combining the inequality with (S6.27), we obtain $D_{\text{WL}}^\dagger \leq \exp(\eta_{01}) \xi_3^2 \nu_2^{-2} (1 - \eta_2)^{-2} |S_{\bar{\alpha}}| \lambda_1^2$.

Therefore, we obtain

$$\begin{aligned} & D_{\text{WL}}^\dagger(\hat{\alpha}^\top \mathbf{G}, \bar{\alpha}^\top \mathbf{G}, \bar{\gamma}) + \exp(\eta_{01})(A_1 - 1)\lambda_1 \|\hat{\alpha} - \bar{\alpha}\|_1 \\ & \leq \exp(4\eta_{01})\xi_2^{-2} \{C_1(1 - \eta_3)\}^{-1} (M_1 |S_{\bar{\gamma}}| \lambda_0^2) + \exp(2\eta_{01})\xi_3^2 \nu_2^{-2} (1 - \eta_2)^{-2} |S_{\bar{\alpha}}| \lambda_1^2, \end{aligned}$$

Let $M_{11} = \exp(4\eta_{01})\xi_2^{-2} \{C_1(1 - \eta_3)\}^{-1} M_1$ and $M_{12} = \exp(2\eta_{01})\xi_3^2 \nu_2^{-2} (1 - \eta_2)^{-2}$, then (4.2) holds. □

S7 Proof of Theorem 1

S7.1 Lemmas for the proposed estimator

Lemma 18. Suppose that Assumptions 1(i), 2(i), 2(ii), and 2(vi) hold, if

$$\lambda_0 \geq \sqrt{2} \{1 + \exp(-B_0)\} C_0 \sqrt{\frac{\ln(\frac{1+P}{\epsilon})}{N}},$$

then, $\mathbb{P}(\Omega_{20}) \geq 1 - 2\epsilon$ for any $r \geq 0$, where Ω_{20} denotes the event

$$\sup_{\|\alpha - \bar{\alpha}\|_1 \leq r; j=0, \dots, d} \left| (\tilde{\mathbb{E}} - \mathbb{E}) \left[\left\{ \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} - 1 \right\} \{ \psi(\hat{\alpha}^\top \mathbf{G}) - \psi(\bar{\alpha}^\top \mathbf{G}) \} Z_j \right] \right| \leq B_1 \lambda_0 r,$$

where B_1 is a positive constant, depending on (C_0, C_2, C_3, C_5) .

Proof. This can be shown similarly to Lemma 13 in the Supplement of Tan (2020a).

□

Lemma 19. Suppose Assumptions 1, 2 and 3 hold, if a function $h(\cdot)$ on a set of samples satisfying $h[\{\mathbf{X}_i\}_{i=1}^N] \leq M\{|S_{\bar{\gamma}}|\lambda_0(\epsilon) + |S_{\bar{\alpha}}|\lambda_1(\epsilon)\}$ for some constant M with probability $1 - c\epsilon$ for some constant $c > 0$ and any $\epsilon > 0$, then, $h[\{\mathbf{X}_i\}_{i=1}^N] = o_p(1)$.

Proof. For $\forall \epsilon > 0$, let Ω_ϵ be the event where $h[\{\mathbf{X}_i\}_{i=1}^N] \leq M\{|S_{\bar{\gamma}}|\lambda_0(\epsilon/c) + |S_{\bar{\alpha}}|\lambda_1(\epsilon/c)\}$ holds. Suppose $|S_{\bar{\gamma}}|\lambda_0(\epsilon/c) < |S_{\bar{\alpha}}|\lambda_1(\epsilon/c)$, then, on Ω_ϵ ,

$$\frac{h[\{\mathbf{X}_i\}_{i=1}^N]\sqrt{N}}{|S_{\bar{\alpha}}|\sqrt{\ln\{e(q+1)\}}} \leq 2M\sqrt{\frac{\ln(\frac{q+1}{\epsilon/c})}{\ln\{e(q+1)\}}} \leq 2M\sqrt{\frac{\ln(q+1)}{\ln(q+1)} + \ln\left(\frac{c}{\epsilon}\right)} \leq 2M\sqrt{1 + \ln\left(\frac{c}{\epsilon}\right)}, \quad (\text{S7.28})$$

which implies that

$$\mathbb{P}\left[\frac{h[\{\mathbf{X}_i\}_{i=1}^N]\sqrt{N}}{|S_{\bar{\alpha}}|\sqrt{\ln\{e(q+1)\}}} \leq 2M\sqrt{1 + \ln\left(\frac{c}{\epsilon}\right)}\right] \geq \mathbb{P}(\Omega_\epsilon) = 1 - \epsilon.$$

Similarly, suppose $|S_{\bar{\alpha}}|\lambda_1 < |S_{\bar{\gamma}}|\lambda_0$, then, on Ω_ϵ ,

$$\mathbb{P}\left[\frac{h[\{\mathbf{X}_i\}_{i=1}^N]\sqrt{N}}{|S_{\bar{\gamma}}|\sqrt{\ln\{e(p+1)\}}} \leq 2M\sqrt{1 + \ln\left(\frac{c}{\epsilon}\right)}\right] \geq \mathbb{P}(\Omega_\epsilon) = 1 - \epsilon.$$

For $\forall \epsilon > 0$, we have

$$\mathbb{P}\left[\frac{h[\{\mathbf{X}_i\}_{i=1}^N]\sqrt{N}}{|S_{\bar{\alpha}}|\sqrt{\ln\{e(q+1)\}} + |S_{\bar{\gamma}}|\sqrt{\ln\{e(p+1)\}}} \leq 2M\sqrt{1 + \ln\left(\frac{c}{\epsilon}\right)}\right] \geq 1 - \epsilon,$$

We write λ_0, λ_1 as $\lambda_0(\epsilon), \lambda_1(\epsilon)$, since we treat them as functions of ϵ .

thus, $\frac{h[\{\mathbf{X}_i\}_{i=1}^N]\sqrt{N}}{|S_{\bar{\alpha}}|\sqrt{\ln\{e(q+1)\}}+|S_{\bar{\gamma}}|\sqrt{\ln\{e(p+1)\}}} = O_p(1)$, and by Assumption 3(vi), it follows that

$$h[\{\mathbf{X}_i\}_{i=1}^N] = o_p\left[\frac{|S_{\bar{\alpha}}|\sqrt{\ln\{e(q+1)\}} + |S_{\bar{\gamma}}|\sqrt{\ln\{e(p+1)\}}}{\sqrt{N}}\right] = o_p(1).$$

□

Lemma 20. Suppose Assumptions 1, 2 and 3 hold, if a function $h(\cdot)$ on a set of samples satisfying $h[\{\mathbf{X}_i\}_{i=1}^N] \leq M(|S_{\bar{\gamma}}| + |S_{\bar{\alpha}}|)\lambda_0(\epsilon)\lambda_1(\epsilon)$ for some constant M with probability $1 - c\epsilon$ for some constant c and any $\epsilon > 0$, then, $h[\{\mathbf{X}_i\}_{i=1}^N] = o_p(1/\sqrt{N})$.

Proof. By similar trick used in the proof of Lemma 19, it can be easily shown that

$$h[\{\mathbf{X}_i\}_{i=1}^N] = o_p\left[\frac{(|S_{\bar{\gamma}}|+|S_{\bar{\alpha}}|)\sqrt{\ln\{e(q+1)\}}\sqrt{\ln\{e(p+1)\}}}{N}\right] = o_p(1/\sqrt{N}) \text{ by Assumption 3(vi).}$$

□

Lemma 21. Under Assumptions 1, 2 and 3, suppose either the PS model (3.1) is correct or the OR model (3.2) is correct, the AIPW estimator $\hat{\beta} \xrightarrow{\mathbb{P}} \beta^*$.

Proof. First, we notice that when $\pi(\mathbf{X}; \gamma)$ is correct or $\phi(\mathbf{X}; \alpha) = \psi(\alpha^T \mathbf{G})$ is correct, β^* is the unique solution to $\mathbb{E}\{\tau(\mathbf{O}, \bar{\alpha}, \beta, \bar{\gamma})\} = 0$. If $\pi(\mathbf{X}; \gamma)$ is correct, we obtain,

$$\begin{aligned} & \mathbb{E}\{\tau(\mathbf{O}, \bar{\alpha}, \beta^*, \bar{\gamma})\} \\ &= \mathbb{E}\left[\frac{R}{\pi(\mathbf{X}; \bar{\gamma})}\{Y - \psi(\beta^{*\top} \mathbf{Z})\} \mathbf{Z} + \left\{1 - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})}\right\} \{\psi(\bar{\alpha}^T \mathbf{G}) - \psi(\beta^{*\top} \mathbf{Z})\} \mathbf{Z}\right] \\ &= \mathbb{E}_{\mathbf{X}}\left(\mathbb{E}\left[\frac{R\{Y - \psi(\beta^{*\top} \mathbf{Z})\}}{\pi(\mathbf{X}; \bar{\gamma})} + \left\{1 - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})}\right\} \{\psi(\bar{\alpha}^T \mathbf{G}) - \psi(\beta^{*\top} \mathbf{Z})\}\right] \mathbf{Z} \middle| \mathbf{X}\right) \end{aligned}$$

$$\begin{aligned}
&= \mathbb{E} [\{Y - \psi(\beta^{*\text{T}} \mathbf{Z})\} \mathbf{Z}] \\
&= 0.
\end{aligned}$$

If $\psi(\alpha^{\text{T}} \mathbf{G})$ is correct, we obtain,

$$\begin{aligned}
&\mathbb{E}\{\tau(\mathbf{O}, \bar{\alpha}, \beta^*, \bar{\gamma})\} \\
&= \mathbb{E} \left[\frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \{Y - \psi(\beta^{*\text{T}} \mathbf{Z})\} \mathbf{Z} + \left\{ 1 - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \right\} \{\psi(\bar{\alpha}^{\text{T}} \mathbf{G}) - \psi(\beta^{*\text{T}} \mathbf{Z})\} \mathbf{Z} \right] \\
&= \mathbb{E} [\{\psi(\bar{\alpha}^{\text{T}} \mathbf{G}) - \psi(\beta^{*\text{T}} \mathbf{Z})\} \mathbf{Z}] \\
&= \mathbb{E} [\{Y - \psi(\beta^{*\text{T}} \mathbf{Z})\} \mathbf{Z}] \\
&= 0.
\end{aligned}$$

The uniqueness is determined by the uniqueness of β^* .

Since Assumptions 3(ii) and 3(iii) hold, by standard argument of consistency (e.g. Van der Vaart (2000)), it suffices to show that

$$\tilde{\mathbb{E}}\{\tau(\mathbf{O}, \bar{\alpha}, \hat{\beta}, \bar{\gamma})\} = o_p(1). \quad (\text{S7.29})$$

We consider the following decomposition,

$$\begin{aligned}
&\tilde{\mathbb{E}}\{\tau(\mathbf{O}, \bar{\alpha}, \hat{\beta}, \bar{\gamma})\} = \tilde{\mathbb{E}}\{\tau(\mathbf{O}, \bar{\alpha}, \hat{\beta}, \bar{\gamma})\} - \tilde{\mathbb{E}}\{\tau(\mathbf{O}, \hat{\alpha}, \hat{\beta}, \hat{\gamma})\} \\
&= \tilde{\mathbb{E}}\{\tau(\mathbf{O}, \bar{\alpha}, \hat{\beta}, \hat{\gamma})\} - \tilde{\mathbb{E}}\{\tau(\mathbf{O}, \hat{\alpha}, \hat{\beta}, \hat{\gamma})\} + \tilde{\mathbb{E}}\{\tau(\mathbf{O}, \bar{\alpha}, \hat{\beta}, \bar{\gamma})\} - \tilde{\mathbb{E}}\{\tau(\mathbf{O}, \bar{\alpha}, \hat{\beta}, \hat{\gamma})\}
\end{aligned}$$

$$=\Delta_0^1 + \Delta_1^1,$$

where

$$\begin{aligned}\Delta_0^1 &= \tilde{\mathbb{E}}\{\tau(\mathbf{O}, \bar{\alpha}, \hat{\beta}, \hat{\gamma})\} - \tilde{\mathbb{E}}\{\tau(\mathbf{O}, \hat{\alpha}, \hat{\beta}, \hat{\gamma})\} = \tilde{\mathbb{E}} \left[\left\{ \frac{R}{\pi(\mathbf{X}; \hat{\gamma})} - 1 \right\} \{\psi(\hat{\alpha}^\top \mathbf{G}) - \psi(\bar{\alpha}^\top \mathbf{G})\} \mathbf{Z} \right], \\ \Delta_1^1 &= \tilde{\mathbb{E}}\{\tau(\mathbf{O}, \bar{\alpha}, \hat{\beta}, \bar{\gamma})\} - \tilde{\mathbb{E}}\{\tau(\mathbf{O}, \bar{\alpha}, \hat{\beta}, \hat{\gamma})\} = \tilde{\mathbb{E}} \left[R \left\{ \frac{1}{\pi(\mathbf{X}; \bar{\gamma})} - \frac{1}{\pi(\mathbf{X}; \hat{\gamma})} \right\} \{Y - \psi(\bar{\alpha}^\top \mathbf{G})\} \mathbf{Z} \right].\end{aligned}$$

By Assumptions 1(iv) and 2(vi), we know that $\exists \zeta_{10} < 0$ a constant, such that $|S_{\bar{\gamma}}|\lambda_0 + |S_{\bar{\alpha}}|\lambda_1 \leq \zeta_{10}$. By (4.2) in Proposition 2, we know that $\exists M_0^0 > 0$ a constant, such that in the event $\Omega_{00} \cap \Omega_{01} \cap \Omega_{03} \cap \Omega_{10} \cap \Omega_{11}$, $\|\hat{\alpha} - \bar{\alpha}\|_1 \leq M_0^0(|S_{\bar{\gamma}}|\lambda_0 + |S_{\bar{\alpha}}|\lambda_1)$.

Consider the j -th coordinate of Δ_0^1 ,

$$\begin{aligned}|\Delta_{0,j}^1| &= \left| \tilde{\mathbb{E}} \left[\left\{ \frac{R}{\pi(\mathbf{X}; \hat{\gamma})} - 1 \right\} \{\psi(\hat{\alpha}^\top \mathbf{G}) - \psi(\bar{\alpha}^\top \mathbf{G})\} Z_j \right] \right| \\ &\leq \tilde{\mathbb{E}}^{1/2} \left[\left\{ \frac{R}{\pi(\mathbf{X}; \hat{\gamma})} - 1 \right\}^2 Z_j^2 \right] \times \tilde{\mathbb{E}}^{1/2} \{\psi(\hat{\alpha}^\top \mathbf{G}) - \psi(\bar{\alpha}^\top \mathbf{G})\}^2 \\ &\leq \tilde{\mathbb{E}}^{1/2} \left([R - 1 + R w(\mathbf{X}; \bar{\gamma}) \exp\{-(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}]^2 Z_j^2 \right) \\ &\quad \times \sqrt{2} \tilde{\mathbb{E}}^{1/2} [\{\psi(\tilde{\alpha}^\top \mathbf{G}) - \psi(\bar{\alpha}^\top \mathbf{G})\} \psi_1(\tilde{\alpha}^\top \mathbf{G}) (\hat{\alpha} - \bar{\alpha})^\top \mathbf{G}] \\ &\leq \sqrt{2} C_5 \tilde{\mathbb{E}}^{1/2} [(R - 1)^2 + R w^2(\mathbf{X}; \bar{\gamma}) \exp\{-2(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}] \\ &\quad \times \tilde{\mathbb{E}}^{1/2} \left[\int_0^1 \psi_1\{\bar{\alpha}^\top \mathbf{G} + u(\tilde{\alpha} - \bar{\alpha})^\top \mathbf{G}\} du \psi_1(\tilde{\alpha}^\top \mathbf{G}) \{(\hat{\alpha} - \bar{\alpha})^\top \mathbf{G}\}^2 \right] \\ &\leq \sqrt{2} C_5 \{1 + \exp(-2B_0 + 2C_0 \|\hat{\gamma} - \bar{\gamma}\|_1)\}^{\frac{1}{2}} \\ &\quad \times \tilde{\mathbb{E}}^{1/2} [\psi_1^2(\bar{\alpha}^\top \mathbf{G}) \exp(2C_3 |\tilde{\alpha}^\top \mathbf{G} - \bar{\alpha}^\top \mathbf{G}|) \{(\hat{\alpha} - \bar{\alpha})^\top \mathbf{G}\}^2]\end{aligned}$$

$$\begin{aligned}
&\leq 2C_5\{1 + \exp(-2B_0 + 2\eta_{01})\}^{\frac{1}{2}} \times C_2C_4 \exp(C_3C_4\|\hat{\alpha} - \bar{\alpha}\|_1)\|\hat{\alpha} - \bar{\alpha}\|_1 \\
&\leq \sqrt{2}C_2C_4C_5\{1 + \exp(-2B_0 + 2\eta_{01})\}^{\frac{1}{2}} \times \exp(C_3C_4M_0^0\zeta_{10})\|\hat{\alpha} - \bar{\alpha}\|_1 \\
&= M_0^1(|S_{\bar{\gamma}}|\lambda_0 + |S_{\bar{\alpha}}|\lambda_1) \\
&= o_p(1)
\end{aligned}$$

where M_0^1 is a constant, $\tilde{\alpha} = u\hat{\alpha} + (1-u)\bar{\alpha}$ for some constant $u \in (0, 1)$ and the last equality holds by Lemma 19.

We consider the j -th coordinate of Δ_1^1 . By Cauchy–Schwarz inequality, Lemma 8, and Lemma 5, since Assumption 2(iv) holds, we obtain

$$\begin{aligned}
|\Delta_{1,j}^1| &= \left| \tilde{\mathbb{E}} \left[R \left\{ \frac{1}{\pi(\mathbf{X}; \bar{\gamma})} - \frac{1}{\pi(\mathbf{X}; \hat{\gamma})} \right\} \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} Z_j \right] \right| \\
&= \left| \tilde{\mathbb{E}} [Rw(\mathbf{X}; \bar{\gamma}) \exp\{-u(\hat{\gamma} - \bar{\gamma})^T \mathbf{F}\} (\bar{\gamma} - \hat{\gamma})^T \mathbf{F} \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} Z_j] \right| \\
&\leq \tilde{\mathbb{E}}^{1/2} [Rw(\mathbf{X}; \bar{\gamma}) \exp\{-2u(\hat{\gamma} - \bar{\gamma})^T \mathbf{F}\} \{(\bar{\gamma} - \hat{\gamma})^T \mathbf{F}\}^2 Z_j^2] \\
&\quad \times \tilde{\mathbb{E}}^{1/2} [Rw(\mathbf{X}; \bar{\gamma}) \{Y - \psi(\bar{\alpha}^T \mathbf{G})\}^2] \\
&\leq \exp(-B_0/2 + C_0\|\hat{\gamma} - \bar{\gamma}\|_1) C_0 C_5 \|\hat{\gamma} - \bar{\gamma}\|_1 \\
&\quad \times \sqrt{\mathbb{E}[Rw(\mathbf{X}; \bar{\gamma}) \{Y - \psi(\bar{\alpha}^T \mathbf{G})\}^2] + (D_0^2 + D_0 D_1) \lambda_0} \\
&\leq \exp(-B_0/2 + \eta_{01}) C_0 C_5 \sqrt{(D_0^2 + D_1^2) \{\lambda_0 + \exp(-B_0)\}} \|\hat{\gamma} - \bar{\gamma}\|_1 \\
&\leq M_1^0 (|S_{\bar{\gamma}}| \lambda_0 + |S_{\bar{\alpha}}| \lambda_1) \\
&= o_p(1),
\end{aligned}$$

where $u \in (0, 1)$ and M_1^0 are both constants. Therefore, $\|\tilde{\mathbb{E}}\{\tau(\mathbf{O}, \bar{\alpha}, \hat{\beta}, \bar{\gamma})\}\|_\infty \leq \|\Delta_0^1\|_\infty + \|\Delta_1^1\|_\infty = o_p(1)$. Hence, $\hat{\beta} \xrightarrow{\mathbb{P}} \beta^*$. \square

S7.2 Proof of Theorem 1(i)

We show the asymptotic normality of $\sqrt{N}(\hat{\beta} - \beta^*)$. First, we consider the following decomposition,

$$\begin{aligned} & \tau(\mathbf{O}, \hat{\alpha}, \hat{\beta}, \hat{\gamma}) - \tau(\mathbf{O}, \bar{\alpha}, \beta^*, \bar{\gamma}) \\ &= \left[\frac{R}{\pi(\mathbf{X}; \hat{\gamma})} \{Y - \psi(\hat{\beta}^\top \mathbf{Z})\} - \left\{ \frac{R}{\pi(\mathbf{X}; \hat{\gamma})} - 1 \right\} \{\psi(\hat{\alpha}^\top \mathbf{G}) - \psi(\hat{\beta}^\top \mathbf{Z})\} \right] \mathbf{Z} \\ & \quad - \left[\frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \{Y - \psi(\beta^{*\top} \mathbf{Z})\} - \left\{ \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} - 1 \right\} \{\psi(\bar{\alpha}^\top \mathbf{G}) - \psi(\beta^{*\top} \mathbf{Z})\} \right] \mathbf{Z} \\ &= \{\psi(\beta^{*\top} \mathbf{Z}) - \psi(\hat{\beta}^\top \mathbf{Z})\} \mathbf{Z} + R \left\{ \frac{1}{\pi(\mathbf{X}; \hat{\gamma})} - \frac{1}{\pi(\mathbf{X}; \bar{\gamma})} \right\} \{Y - \psi(\bar{\alpha}^\top \mathbf{G})\} \mathbf{Z} \\ & \quad + \left\{ 1 - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \right\} \{\psi(\hat{\alpha}^\top \mathbf{G}) - \psi(\bar{\alpha}^\top \mathbf{G})\} \mathbf{Z} \\ & \quad + \left\{ \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} - \frac{R}{\pi(\mathbf{X}; \hat{\gamma})} \right\} \{\psi(\hat{\alpha}^\top \mathbf{G}) - \psi(\bar{\alpha}^\top \mathbf{G})\} \mathbf{Z}, \end{aligned}$$

denoted as $\delta_0^0 + \delta_1^0 + \delta_2^0 + \delta_3^0$. Then, $-\tilde{\mathbb{E}}\tau(\mathbf{O}, \bar{\alpha}, \beta^*, \bar{\gamma}) = \Delta_0^2 + \Delta_1^2 + \Delta_2^2 + \Delta_3^2$, with $\Delta_i^2 = \tilde{\mathbb{E}}(\delta_i^0)$, $i = 0, 1, 2, 3$. First, we show that $\Delta_1^2 + \Delta_2^2 + \Delta_3^2 = o_p(1/\sqrt{N})$. To upper-bound Δ_1^2 , consider $\Delta_{1,j}^2$, the j -th coordinate of Δ_1^2 . By Taylor expansion in a neighborhood of $\bar{\gamma}$,

$$\begin{aligned} \Delta_{1,j}^2 &= \tilde{\mathbb{E}} \left[\left\{ \frac{R}{\pi(\mathbf{X}; \hat{\gamma})} - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \right\} \{Y - \psi(\bar{\alpha}^\top \mathbf{G})\} Z_j \right] \\ &= -(\hat{\gamma} - \bar{\gamma})^\top \tilde{\mathbb{E}} [\mathbf{F}Rw(\mathbf{X}; \bar{\gamma}) \{Y - \psi(\bar{\alpha}^\top \mathbf{G})\} Z_j] \end{aligned}$$

$$+ \frac{1}{2}(\hat{\gamma} - \bar{\gamma})^T \tilde{\mathbb{E}} [\mathbf{F} R w(\mathbf{X}, \tilde{\gamma}_j) \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} Z_j \mathbf{F}^T] (\hat{\gamma} - \bar{\gamma}),$$

denoted as $\Delta_{10,j}^2 + \Delta_{11,j}^2$ where $\tilde{\gamma}_j = u_j \hat{\gamma} + (1 - u_j) \bar{\gamma}$ for some $u_j \in (0, 1)$.

In the event $\Omega_{00} \cap \Omega_{01} \cap \Omega_{10}$, by (4.1) and Lemma 10, we obtain

$$|\Delta_{10,j}^2| \leq \|\hat{\gamma} - \bar{\gamma}\|_1 \|\tilde{\mathbb{E}} [\mathbf{F} R w(\mathbf{X}; \bar{\gamma}) \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} Z_j] \|_\infty \leq M_{10,j}^1 |S_{\bar{\gamma}}| \lambda_0 \lambda_1,$$

for some constant $M_{10,j}^1 > 0$. In the event $\Omega_{00} \cap \Omega_{01}$, by (S5.9), we obtain

$$\begin{aligned} |\Delta_{11,j}^2| &= \frac{1}{2} |(\hat{\gamma} - \bar{\gamma})^T \tilde{\mathbb{E}} [\mathbf{F} R w(\mathbf{X}; \tilde{\gamma}_j) \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} Z_j \mathbf{F}^T] (\hat{\gamma} - \bar{\gamma})| \\ &\leq \frac{1}{2} C_5 \exp(\|\hat{\gamma} - \bar{\gamma}\|_1 C_0) \tilde{\mathbb{E}} \{R w(\mathbf{X}; \bar{\gamma}) |Y - \psi(\bar{\alpha}^T \mathbf{G})| |(\hat{\gamma} - \bar{\gamma})^T \mathbf{F}|^2\}. \end{aligned}$$

We bound $\tilde{\mathbb{E}} \{R w(\mathbf{X}; \bar{\gamma}) |Y - \psi(\bar{\alpha}^T \mathbf{G})| |(\hat{\gamma} - \bar{\gamma})^T \mathbf{F}|^2\}$ by following steps. First, by Lemma 7, in the event Ω_{02} ,

$$(\tilde{\mathbb{E}} - \mathbb{E}) [R w(\mathbf{X}; \bar{\gamma}) |Y - \psi(\bar{\alpha}^T \mathbf{G})| \{(\hat{\gamma} - \bar{\gamma})^T \mathbf{F}\}^2] \leq \sqrt{D_0^2 + D_1^2} \|\hat{\gamma} - \bar{\gamma}\|_1^2 \lambda_0.$$

Second, by Assumption 2(iv) and Lemma 5, we have

$$\mathbb{E}[\{Y - \psi(\bar{\alpha}^T \mathbf{G})\}^2 | \mathbf{X}] \leq D_0^2 + D_1^2.$$

Therefore,

$$\begin{aligned}
& \mathbb{E} [Rw(\mathbf{X}; \bar{\gamma}) | Y - \psi(\bar{\alpha}^\top \mathbf{G}) | \{(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}^2] \\
& \leq \mathbb{E}^{1/2} [Rw(\mathbf{X}; \bar{\gamma}) | Y - \psi(\bar{\alpha}^\top \mathbf{G})|^2 \{(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}^2] \times \mathbb{E}^{1/2} [Rw(\mathbf{X}; \bar{\gamma}) \{(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}^2] \\
& \leq \sqrt{D_0^2 + D_1^2} \mathbb{E} [Rw(\mathbf{X}; \bar{\gamma}) \{(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}^2].
\end{aligned}$$

Third, in the event Ω_{01} , by (S5.6),

$$(\mathbb{E} - \tilde{\mathbb{E}}) [Rw(\mathbf{X}; \bar{\gamma}) \{(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}^2] \leq \lambda_0 \|\bar{\gamma} - \hat{\gamma}\|_1^2.$$

Combining preceding inequalities and (S5.9) in Lemma 9, in the event $\Omega_{00} \cap \Omega_{01} \cap \Omega_{03} \cap \Omega_{02}$,

$$\begin{aligned}
& \tilde{\mathbb{E}} \{Rw(\mathbf{X}; \bar{\gamma}) | Y - \psi(\bar{\alpha}^\top \mathbf{G}) | \{(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}^2\} \\
& \leq \sqrt{D_0^2 + D_1^2} \|\hat{\gamma} - \bar{\gamma}\|_1^2 \lambda_0 + \sqrt{D_0^2 + D_1^2} \{\lambda_0 \|\bar{\gamma} - \hat{\gamma}\|_1^2 + \exp(\eta_{01}) M_0 |S_{\bar{\gamma}}| \lambda_0^2\} \\
& \leq 2\sqrt{D_0^2 + D_1^2} \left(\frac{M_0}{A_0 - 1}\right)^2 |S_{\bar{\gamma}}|^2 \lambda_0^3 + \sqrt{D_0^2 + D_1^2} \exp(\eta_{01}) M_0 |S_{\bar{\gamma}}| \lambda_0^2 \\
& \leq 2\zeta_0 \sqrt{D_0^2 + D_1^2} \left(\frac{M_0}{A_0 - 1}\right)^2 |S_{\bar{\gamma}}| \lambda_0^2 + \sqrt{D_0^2 + D_1^2} \exp(\eta_{01}) M_0 |S_{\bar{\gamma}}| \lambda_0^2.
\end{aligned}$$

Therefore,

$$\begin{aligned}
|\Delta_{11,j}^2| & \leq \frac{1}{2} C_5 \exp(\|\hat{\gamma} - \bar{\gamma}\|_1 C_0) \tilde{\mathbb{E}} \{Rw(\mathbf{X}; \bar{\gamma}) | Y - \psi(\bar{\alpha}^\top \mathbf{G}) | \{(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}^2\} \\
& \leq \frac{1}{2} C_5 \exp(\eta_{01}) \times \left\{ 2\zeta_0 \sqrt{D_0^2 + D_1^2} \left(\frac{M_0}{A_0 - 1}\right)^2 |S_{\bar{\gamma}}| \lambda_0^2 + \sqrt{D_0^2 + D_1^2} \exp(\eta_{01}) M_0 |S_{\bar{\gamma}}| \lambda_0^2 \right\}
\end{aligned}$$

$$=M_{11}^1|S_{\bar{\gamma}}|\lambda_0^2,$$

for some constant $M_{11}^1 > 0$. Hence, in the event $\Omega_{00} \cap \Omega_{01} \cap \Omega_{03} \cap \Omega_{02} \cap \Omega_{10}$,

$$\begin{aligned} & \sup_{j=0,\dots,m-1} |\Delta_{1,j}^2| \\ & \leq \sup_{j=0,\dots,m-1} (|\Delta_{10,j}^2| + |\Delta_{11,j}^2|) \\ & \leq \sup_{j=0,\dots,m-1} M_{10}^1|S_{\bar{\gamma}}|\lambda_0\lambda_1 + \sup_{j=0,\dots,m-1} M_{11}^1|S_{\bar{\gamma}}|\lambda_0^2 \\ & = M_{10}^1|S_{\bar{\gamma}}|\lambda_0\lambda_1 + M_{11}^1|S_{\bar{\gamma}}|\lambda_0^2 \\ & \leq M_1^1(|S_{\bar{\gamma}}|\lambda_0\lambda_1 + |S_{\bar{\gamma}}|\lambda_0^2), \end{aligned}$$

where $M_1^1 = \max(M_{10}^1, M_{11}^1)$.

To bound Δ_2^2 , consider $\Delta_{2,j}^2$, the j -th coordinate of Δ_2^2 , $\Delta_{2,j}^2$ can be decomposed as

$$\begin{aligned} \Delta_{2,j}^2 &= (\tilde{\mathbb{E}} - \mathbb{E}) \left[\left\{ 1 - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \right\} \{ \psi(\hat{\alpha}^T \mathbf{G}) - \psi(\bar{\alpha}^T \mathbf{G}) \} Z_j \right] \\ &+ \mathbb{E} \left[\left\{ 1 - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \right\} \{ \psi(\hat{\alpha}^T \mathbf{G}) - \psi(\bar{\alpha}^T \mathbf{G}) \} Z_j \right], \end{aligned}$$

denoted as $\Delta_{20,j}^2 + \Delta_{21,j}^2$. Since $\pi(\mathbf{X}; \bar{\gamma})$ is correct, $\Delta_{21,j}^2 = 0$. By (4.2), $\exists M_\alpha > 0$, such that $\|\hat{\alpha} - \bar{\alpha}\|_1 \leq M_\alpha(|S_{\bar{\gamma}}|\lambda_0 + |S_{\bar{\alpha}}|\lambda_1)$. Take $r_\alpha = M_\alpha(|S_{\bar{\gamma}}|\lambda_0 + |S_{\bar{\alpha}}|\lambda_1)$ in Lemma 18, then in the event $\Omega_{00} \cap \Omega_{01} \cap \Omega_{03} \cap \Omega_{10} \cap \Omega_{11} \cap \Omega_{20}$, we have $\|\hat{\alpha} - \bar{\alpha}\|_1 \leq r_\alpha$ and

Note that $\lambda_0 \leq \lambda_1$.

hence

$$|\Delta_{20,j}^2| \leq B_1 M_\alpha (|S_{\bar{\gamma}}| \lambda_0 + |S_{\bar{\alpha}}| \lambda_1) \lambda_0. \quad (\text{S7.30})$$

Thus, in the event $\Omega_{00} \cap \Omega_{01} \cap \Omega_{03} \cap \Omega_{10} \cap \Omega_{11} \cap \Omega_{20}$,

$$\sup_{j=0,\dots,m-1} |\Delta_{2,j}^2| \leq B_1 M_\alpha (|S_{\bar{\gamma}}| \lambda_0 + |S_{\bar{\alpha}}| \lambda_1) \lambda_0 = M_2^1 (|S_{\bar{\gamma}}| \lambda_0 + |S_{\bar{\alpha}}| \lambda_1) \lambda_0, \quad (\text{S7.31})$$

for some positive constant M_2^1 .

To deal with Δ_3^2 , first, by mean value theorem, we obtain for some $u \in (0, 1)$,

$$\begin{aligned} \frac{1}{\pi(\mathbf{X}; \hat{\gamma})} - \frac{1}{\pi(\mathbf{X}; \bar{\gamma})} &= -\exp\{-u\hat{\gamma}^\top \mathbf{F} - (1-u)\bar{\gamma}^\top \mathbf{F}\} (\hat{\gamma} - \bar{\gamma})^\top \mathbf{F} \\ &= -w(\mathbf{X}; \bar{\gamma}) \exp\{-u(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\} (\hat{\gamma} - \bar{\gamma})^\top \mathbf{F} \end{aligned} \quad (\text{S7.32})$$

and for some $\tilde{\alpha}$ lies between $\bar{\alpha}$ and $\hat{\alpha}$,

$$\psi(\bar{\alpha}^\top \mathbf{G}) - \psi(\hat{\alpha}^\top \mathbf{G}) = -\psi_1(\tilde{\alpha}^\top \mathbf{G})(\hat{\alpha} - \bar{\alpha})^\top \mathbf{G}. \quad (\text{S7.33})$$

Combining (S7.32) and (S7.33) and applying Cauchy-Schwartz inequality to j -th

coordinate of Δ_3^2 in the event $\Omega_{00} \cap \Omega_{01} \cap \Omega_{03} \cap \Omega_{10} \cap \Omega_{11} \cap \Omega_{20}$, we get

$$\begin{aligned} |\Delta_{3,j}^2| &= \left| \tilde{\mathbb{E}} \left[\left\{ \frac{R}{\pi(\mathbf{X}; \hat{\gamma})} - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \right\} \{\psi(\bar{\alpha}^\top \mathbf{G}) - \psi(\hat{\alpha}^\top \mathbf{G})\} Z_j \right] \right| \\ &= |\tilde{\mathbb{E}}([Rw(\mathbf{X}; \bar{\gamma}) \exp\{-u_j(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\} (\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}] \{\psi_1(\tilde{\alpha}_j^\top \mathbf{G})(\hat{\alpha} - \bar{\alpha})^\top \mathbf{G}\}) Z_j| \\ &\leq C_5 \exp(\eta_{01}) \tilde{\mathbb{E}}^{1/2} |Rw(\mathbf{X}; \bar{\gamma}) \{(\hat{\gamma} - \bar{\gamma})^\top \mathbf{F}\}^2| \times \tilde{\mathbb{E}}^{1/2} |Rw(\mathbf{X}; \bar{\gamma}) \psi_1^2(\tilde{\alpha}_j^\top \mathbf{G}) \{(\hat{\alpha} - \bar{\alpha})^\top \mathbf{G}\}^2| \end{aligned}$$

$$\begin{aligned}
&\leq C_5 \exp(\eta_{01}) \{ \exp(\eta_{01}) M_0 |S_{\bar{\gamma}}| \lambda_0^2 \}^{1/2} \\
&\quad \times \tilde{\mathbb{E}}^{1/2} \left(R w(\mathbf{X}; \bar{\gamma}) [\psi_1(\bar{\alpha}_j^T \mathbf{G}) \exp\{C_3 |(\tilde{\alpha} - \bar{\alpha})^T \mathbf{G}\}|]^2 \{(\hat{\alpha} - \bar{\alpha})^T \mathbf{G}\}^2 \right) \\
&\leq C_5 \exp(\eta_{01}) \{ \exp(\eta_{01}) M_0 |S_{\bar{\gamma}}| \lambda_0^2 \}^{1/2} C_2^{1/2} \exp(C_3 C_4 r_\alpha) \left\{ \frac{M^\dagger (|S_{\bar{\gamma}}| \lambda_0^2 + |S_{\bar{\alpha}}| \lambda_1^2)}{1 - \eta_3} \right\}^{1/2} \\
&\leq M_3^1 \{ |S_{\bar{\gamma}}| \lambda_0^2 + (|S_{\bar{\gamma}}| |S_{\bar{\alpha}}|)^{1/2} \lambda_0 \lambda_1 \}.
\end{aligned}$$

The second inequality holds due to (S5.9) in Lemma 9. The third inequality holds by Assumption 2(ii), (S6.25), the facts that $\tilde{\mathbb{E}} [R w(\mathbf{X}; \bar{\gamma}) \psi_1(\bar{\alpha}_j^T \mathbf{G}) \{(\hat{\alpha} - \bar{\alpha})^T \mathbf{G}\}^2] = b^T \tilde{\Sigma}_\alpha b$ in (S6.23) and that by (4.2) in Proposition 2, $\exists M^\dagger > 0$ a constant, such that $D_{\text{WL}}^\dagger \leq M^\dagger (|S_{\bar{\gamma}}| \lambda_0^2 + |S_{\bar{\alpha}}| \lambda_1^2)$. Therefore,

$$\sup_{j=0, \dots, m-1} |\Delta_{3,j}^2| \leq M_3^1 \{ |S_{\bar{\gamma}}| \lambda_0^2 + (|S_{\bar{\gamma}}| |S_{\bar{\alpha}}|)^{1/2} \lambda_0 \lambda_1 \}, \quad (\text{S7.34})$$

for some constant M_3^1 . Thus, on the event $\Omega_{00} \cap \Omega_{01} \cap \Omega_{03} \cap \Omega_{02} \cap \Omega_{10} \cap \Omega_{11} \cap \Omega_{20}$,

$$\begin{aligned}
\|\Delta_1^2 + \Delta_2^2 + \Delta_3^2\|_\infty &\leq \sup_{j=0, \dots, m-1} |\Delta_{1,j}^2| + \sup_{j=0, \dots, m-1} |\Delta_{2,j}^2| + \sup_{j=0, \dots, m-1} |\Delta_{3,j}^2| \\
&\leq M_1^1 (|S_{\bar{\gamma}}| \lambda_0 \lambda_1 + |S_{\bar{\gamma}}| \lambda_0^2) + M_2^1 (|S_{\bar{\gamma}}| \lambda_0 + |S_{\bar{\alpha}}| \lambda_1) \lambda_0 \\
&\quad + M_3^1 \{ |S_{\bar{\gamma}}| \lambda_0^2 + (|S_{\bar{\gamma}}| |S_{\bar{\alpha}}|)^{1/2} \lambda_0 \lambda_1 \} \\
&\leq M^1 (|S_{\bar{\gamma}}| + |S_{\bar{\alpha}}|) \lambda_0 \lambda_1.
\end{aligned}$$

By Lemma 20, $\|\Delta_1^2 + \Delta_2^2 + \Delta_3^2\|_\infty = o_p(1/\sqrt{N})$.

Then, we deal with Δ_0^2 . For the j -th coordinate of Δ_0^2 , $\Delta_{0,j}^2$, by mean value

theorem,

$$\Delta_{0,j}^2 = \tilde{\mathbb{E}}[\{\psi(\beta^{*\text{T}} \mathbf{Z}) - \psi(\hat{\beta}^{\text{T}} \mathbf{Z})\} Z_j] = -\tilde{\mathbb{E}}\{\psi_1(\tilde{\beta}_j^{\text{T}} \mathbf{Z}) Z_j \mathbf{Z}^{\text{T}} (\hat{\beta} - \beta^*)\},$$

where $\tilde{\beta}_j = (1-u_j)\beta^* + u_j\hat{\beta}$ for some $u_j \in (0, 1)$. We first show that $\tilde{\mathbb{E}}\{\psi_1(\beta^{*\text{T}} \mathbf{Z}) Z_j \mathbf{Z}\} - \tilde{\mathbb{E}}\{\psi_1(\tilde{\beta}_j^{\text{T}} \mathbf{Z}) Z_j \mathbf{Z}\} \xrightarrow{\mathbb{P}} 0$. By Assumption 2(ii), we know that for $\forall u, u'$, if $\psi_1(u) > \psi_1(u')$, since $\psi_1(u') \geq \psi_1(u) \exp(-C_3|u - u'|)$, then $|\psi_1(u) - \psi_1(u')| \leq \psi_1(u)\{1 - \exp(-C_3|u - u'|)\}$; if $\psi_1(u) < \psi_1(u')$, $\psi_1(u') \leq \psi_1(u) \exp(C_3|u - u'|)$, $|\psi_1(u) - \psi_1(u')| \leq \psi_1(u)\{\exp(C_3|u - u'|) - 1\}$; therefore, $|\psi_1(u) - \psi_1(u')| \leq \psi_1(u) \max\{1 - \exp(-C_3|u - u'|), \exp(C_3|u - u'|) - 1\}$. Consider the i -th element of the difference, if $C_3 = 0$, ψ_1 is a constant, then

$$|\tilde{\mathbb{E}}\{\psi_1(\beta^{*\text{T}} X) Z_i Z_j\} - \tilde{\mathbb{E}}\{\psi_1(\tilde{\beta}_j^{\text{T}} \mathbf{Z}) Z_i Z_j\}| = 0.$$

Otherwise,

$$\begin{aligned} & |\tilde{\mathbb{E}}\{\psi_1(\beta^{*\text{T}} X) Z_i Z_j\} - \tilde{\mathbb{E}}\{\psi_1(\tilde{\beta}_j^{\text{T}} \mathbf{Z}) Z_i Z_j\}| \leq \tilde{\mathbb{E}}\{|\psi_1(\beta^{*\text{T}} \mathbf{Z}) - \psi_1(\tilde{\beta}_j^{\text{T}} \mathbf{Z})| |Z_i Z_j|\} \\ & \leq C_5^2 \tilde{\mathbb{E}}[\psi_1(\beta^{*\text{T}} \mathbf{Z}) \max\{1 - \exp(-C_3|(\hat{\beta} - \beta^*)^{\text{T}} \mathbf{Z}|), \exp(C_3|(\hat{\beta} - \beta^*)^{\text{T}} \mathbf{Z}) - 1\}] \\ & = C_5^2 C_6 \tilde{\mathbb{E}}\{C_3|(\hat{\beta} - \beta^*)^{\text{T}} \mathbf{Z}| + o_p(|(\hat{\beta} - \beta^*)^{\text{T}} \mathbf{Z}|)\} \\ & \leq C_3 C_5^3 C_6 \|\hat{\beta} - \beta^*\|_1 + o_p(\|\hat{\beta} - \beta^*\|_1) \\ & = o_p(1), \end{aligned}$$

which leads to $-\tilde{\mathbb{E}}\{\psi_1(\tilde{\beta}_j^T \mathbf{Z})Z_j \mathbf{Z}^T\} + \tilde{\mathbb{E}}\{\psi_1(\beta^{*T} \mathbf{Z})Z_j \mathbf{Z}^T\} \xrightarrow{\mathbb{P}} 0$. Therefore, we consider the following decomposition,

$$\begin{aligned} & -\tilde{\mathbb{E}}\{\psi_1(\tilde{\beta}_j^T \mathbf{Z})Z_j \mathbf{Z}\} + \mathbb{E}\{\psi_1(\beta^{*T} \mathbf{Z})Z_j \mathbf{Z}\} \\ &= -\tilde{\mathbb{E}}\{\psi_1(\tilde{\beta}_j^T \mathbf{Z})Z_j \mathbf{Z}\} + \tilde{\mathbb{E}}\{\psi_1(\beta^{*T} X)Z_j \mathbf{Z}\} - \tilde{\mathbb{E}}\{\psi_1(\beta^{*T} \mathbf{Z})Z_j \mathbf{Z}\} + E\{\psi_1(\beta^{*T} \mathbf{Z})Z_j \mathbf{Z}\} \\ & \xrightarrow{\mathbb{P}} 0. \end{aligned}$$

Hence,

$$\tilde{\mathbb{E}}\{\psi_1(\tilde{\beta}_j^T \mathbf{Z})Z_j \mathbf{Z}\} \xrightarrow{\mathbb{P}} \mathbb{E}\{\psi_1(\beta^{*T} \mathbf{Z})Z_j \mathbf{Z}\} = \mathbf{\Gamma}_j, \quad (\text{S7.35})$$

where $\mathbf{\Gamma}_j$ is the j -th row of $\mathbf{\Gamma}$, and (S7.35) holds for $j = 0, \dots, m-1$. Hence, $\Delta_0^2 = \tilde{\mathbb{E}}(\tilde{\mathbf{Z}} \mathbf{Z}^T)(\hat{\beta} - \beta^*)$, where $\tilde{\mathbf{Z}}_j = -\psi_1(\tilde{\beta}_j^T \mathbf{Z})Z_j$, and $\tilde{\mathbb{E}}(\tilde{\mathbf{Z}} \mathbf{Z}^T) \xrightarrow{\mathbb{P}} -\mathbf{\Gamma}$. Suppose $\sqrt{N}(\hat{\beta} - \beta^*) \xrightarrow{d} G_2$, by continuous mapping theorem,

$$\sqrt{N}\Delta_0^2 \xrightarrow{d} -\mathbf{\Gamma}G_2. \quad (\text{S7.36})$$

Besides, by central limit theorem,

$$\sqrt{N}\Delta_0^2 = \sqrt{N}\tilde{\mathbb{E}}\{\tau(\mathbf{O}, \bar{\alpha}, \beta^*, \bar{\gamma})\} + o_p(1) \xrightarrow{d} N(0, \mathbf{\Lambda}), \quad (\text{S7.37})$$

where

$$\mathbf{\Lambda} = \mathbb{E}\{\tau(\mathbf{O}, \bar{\alpha}, \beta^*, \bar{\gamma})\tau(\mathbf{O}, \bar{\alpha}, \beta^*, \bar{\gamma})^T\}$$

$$\begin{aligned}
&= \mathbb{E} \left(\left[\frac{1}{\pi(\mathbf{X}; \bar{\gamma})} \{Y - \psi(\bar{\alpha}^\top \mathbf{G})\}^2 + \{\psi(\bar{\alpha}^\top \mathbf{G}) - \psi(\beta^{*\top} \mathbf{Z})\}^2 \right] \mathbf{Z} \mathbf{Z}^\top \right) \\
&\quad + 2\mathbb{E}[\{Y - \psi(\bar{\alpha}^\top \mathbf{G})\} \{\psi(\bar{\alpha}^\top \mathbf{G}) - \psi(\beta^{*\top} \mathbf{Z})\} \mathbf{Z} \mathbf{Z}^\top]. \tag{S7.38}
\end{aligned}$$

Therefore,

$$\sqrt{N}(\hat{\beta} - \beta^*) \xrightarrow{d} G_2 \sim -\mathbf{\Gamma}^{-1} \mathbf{N}(0, \mathbf{\Lambda}) \sim \mathbf{N}(0, \mathbf{\Sigma}),$$

where \sim denotes “distributed as”, i.e., for any two distributions G_0 and G_1 , $G_0 \sim G_1$ means the two distributions are the same.

S7.3 Proof of Theorem 1(ii)

We show the consistency of $\hat{\Sigma}$. First, if we let $\hat{\mathbf{\Gamma}}_j = \tilde{\mathbb{E}}\{\psi_1(\hat{\beta}^\top \mathbf{Z}) Z_j \mathbf{Z}\}$, then

$$\hat{\mathbf{\Gamma}}_j \xrightarrow{\mathbb{P}} \mathbb{E}\{\psi_1(\beta^{*\top} \mathbf{Z}) Z_j \mathbf{Z}\},$$

i.e., $\hat{\mathbf{\Gamma}}_j \xrightarrow{\mathbb{P}} \mathbf{\Gamma}_j$, which can be shown in the way similar to the proof of (S7.35). Then we get $\hat{\mathbf{\Gamma}} \xrightarrow{\mathbb{P}} \mathbf{\Gamma}$.

Next, we want to show that $\hat{\mathbf{\Lambda}} \xrightarrow{\mathbb{P}} \mathbf{\Lambda}$. Since $\mathbf{\Lambda} = \mathbb{E}\{\tau(\mathbf{O}, \bar{\alpha}, \beta^*, \bar{\gamma}) \tau(\mathbf{O}, \bar{\alpha}, \beta^*, \bar{\gamma})^\top\}$ and $\hat{\mathbf{\Lambda}} = \tilde{\mathbb{E}}\{\tau(\mathbf{O}, \hat{\alpha}, \hat{\beta}, \hat{\gamma}) \tau(\mathbf{O}, \hat{\alpha}, \hat{\beta}, \hat{\gamma})^\top\}$, it suffices to show that

$$\tilde{\mathbb{E}}\{\tau(\mathbf{O}, \hat{\alpha}, \hat{\beta}, \hat{\gamma}) \tau(\mathbf{O}, \hat{\alpha}, \hat{\beta}, \hat{\gamma})^\top\} - \tilde{\mathbb{E}}\{\tau(\mathbf{O}, \bar{\alpha}, \beta^*, \bar{\gamma}) \tau(\mathbf{O}, \bar{\alpha}, \beta^*, \bar{\gamma})^\top\} = o_p(1).$$

We consider the i, j -th element of the difference above:

$$\begin{aligned}
& \left| \tilde{\mathbb{E}}\{\tau(\mathbf{O}, \hat{\alpha}, \hat{\beta}, \hat{\gamma})\tau(\mathbf{O}, \hat{\alpha}, \hat{\beta}, \hat{\gamma})^T\} - \tilde{\mathbb{E}}\{\tau(\mathbf{O}, \bar{\alpha}, \beta^*, \bar{\gamma})\tau(\mathbf{O}, \bar{\alpha}, \beta^*, \bar{\gamma})^T\} \right|_{ij} \\
&= \left| \tilde{\mathbb{E}} \left(\left[\frac{R}{\pi(\mathbf{X}; \hat{\gamma})} \{Y - \psi(\hat{\alpha}^T \mathbf{G})\} + \{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\hat{\beta}^T \mathbf{Z})\} \right]^2 Z_i Z_j \right. \right. \\
&\quad \left. \left. - \left[\frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} + \{\psi(\bar{\alpha}^T \mathbf{G}) - \psi(\beta^{*T} \mathbf{Z})\} \right]^2 Z_i Z_j \right) \right| \\
&\leq \tilde{\mathbb{E}} \left(\left| \left[\frac{R}{\pi(\mathbf{X}; \hat{\gamma})} \{Y - \psi(\hat{\alpha}^T \mathbf{G})\} + \{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\hat{\beta}^T \mathbf{Z})\} \right]^2 \right. \right. \\
&\quad \left. \left. - \left[\frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} + \{\psi(\bar{\alpha}^T \mathbf{G}) - \psi(\beta^{*T} \mathbf{Z})\} \right]^2 \right| \left| Z_i Z_j \right| \right) \\
&\leq C_5^2 \tilde{\mathbb{E}} \left(\left| \left[\frac{R}{\pi(\mathbf{X}; \hat{\gamma})} \{Y - \psi(\hat{\alpha}^T \mathbf{G})\} + \{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\hat{\beta}^T \mathbf{Z})\} \right]^2 \right. \right. \\
&\quad \left. \left. - \left[\frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} + \{\psi(\bar{\alpha}^T \mathbf{G}) - \psi(\beta^{*T} \mathbf{Z})\} \right]^2 \right| \right) \\
&\leq C_5^2 \tilde{\mathbb{E}} \left(\left[\frac{R}{\pi(\mathbf{X}; \hat{\gamma})} \{Y - \psi(\hat{\alpha}^T \mathbf{G})\} + \{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\hat{\beta}^T \mathbf{Z})\} \right. \right. \\
&\quad \left. \left. - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} - \{\psi(\bar{\alpha}^T \mathbf{G}) - \psi(\beta^{*T} \mathbf{Z})\} \right]^2 \right) \\
&\quad + 2C_5^2 \tilde{\mathbb{E}} \left(\left[\frac{R}{\pi(\mathbf{X}; \hat{\gamma})} \{Y - \psi(\hat{\alpha}^T \mathbf{G})\} + \{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\hat{\beta}^T \mathbf{Z})\} \right. \right. \\
&\quad \left. \left. - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} - \{\psi(\bar{\alpha}^T \mathbf{G}) - \psi(\beta^{*T} \mathbf{Z})\} \right] \right. \\
&\quad \left. \times \left[\frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} + \{\psi(\bar{\alpha}^T \mathbf{G}) - \psi(\beta^{*T} \mathbf{Z})\} \right] \right) \\
&\leq C_5^2 \tilde{\mathbb{E}} \left(\left[\frac{R}{\pi(\mathbf{X}; \hat{\gamma})} \{Y - \psi(\hat{\alpha}^T \mathbf{G})\} + \{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\hat{\beta}^T \mathbf{Z})\} \right. \right. \\
&\quad \left. \left. - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} - \{\psi(\bar{\alpha}^T \mathbf{G}) - \psi(\beta^{*T} \mathbf{Z})\} \right]^2 \right) \\
&\quad + 2C_5^2 \tilde{\mathbb{E}}^{\frac{1}{2}} \left(\left[\frac{R}{\pi(\mathbf{X}; \hat{\gamma})} \{Y - \psi(\hat{\alpha}^T \mathbf{G})\} + \{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\hat{\beta}^T \mathbf{Z})\} \right. \right.
\end{aligned}$$

$$\begin{aligned}
& - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} - \{\psi(\bar{\alpha}^T \mathbf{G}) - \psi(\beta^{*\top} \mathbf{Z})\} \Big]^2 \Big) \\
& \times \tilde{\mathbb{E}}^{\frac{1}{2}} \left(\left[\frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} + \{\psi(\bar{\alpha}^T \mathbf{G}) - \psi(\beta^{*\top} \mathbf{Z})\} \right]^2 \right),
\end{aligned}$$

therefore, we only need to show that

$$\begin{aligned}
& \tilde{\mathbb{E}} \left(\left[\frac{R}{\pi(\mathbf{X}; \hat{\gamma})} \{Y - \psi(\hat{\alpha}^T \mathbf{G})\} + \{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\hat{\beta}^T \mathbf{Z})\} \right. \right. \\
& \quad \left. \left. - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} - \{\psi(\bar{\alpha}^T \mathbf{G}) - \psi(\beta^{*\top} \mathbf{Z})\} \right]^2 \right) \\
& = o_p(1).
\end{aligned}$$

Consider the following decomposition:

$$\begin{aligned}
& \frac{R}{\pi(\mathbf{X}; \hat{\gamma})} \{Y - \psi(\hat{\alpha}^T \mathbf{G})\} + \{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\hat{\beta}^T \mathbf{Z})\} \\
& - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} - \{\psi(\bar{\alpha}^T \mathbf{G}) - \psi(\beta^{*\top} \mathbf{Z})\} \\
& = \{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\bar{\alpha}^T \mathbf{G})\} \left\{ 1 - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \right\} \\
& + R \{Y - \psi(\bar{\alpha}^T \mathbf{G})\} \left\{ \frac{1}{\pi(\mathbf{X}; \hat{\gamma})} - \frac{1}{\pi(\mathbf{X}; \bar{\gamma})} \right\} \\
& + R \{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\bar{\alpha}^T \mathbf{G})\} \left\{ \frac{1}{\pi(\mathbf{X}; \bar{\gamma})} - \frac{1}{\pi(\mathbf{X}; \hat{\gamma})} \right\} \\
& + \psi(\beta^{*\top} \mathbf{Z}) - \psi(\hat{\beta}^T \mathbf{Z}),
\end{aligned}$$

denoted as $\delta_0^1 + \delta_1^1 + \delta_2^1 + \delta_3^1$. Let $\Delta_i^3 = \tilde{\mathbb{E}}\{(\delta_i^1)^2\}$, $i = 0, \dots, 3$, we only need to show that $\Delta_i^3 = o_p(1)$, $i = 0, \dots, 3$.

By mean value theorem,

$$\begin{aligned}
\Delta_0^3 &= \tilde{\mathbb{E}} \left[\{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\bar{\alpha}^T \mathbf{G})\}^2 \left\{ 1 - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \right\}^2 \right] \\
&= 2\tilde{\mathbb{E}} \left[\psi_1(\tilde{\alpha}^T \mathbf{G}) \{\psi(\tilde{\alpha}^T \mathbf{G}) - \psi(\bar{\alpha}^T \mathbf{G})\} (\hat{\alpha} - \bar{\alpha})^T \mathbf{G} \left\{ 1 - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \right\}^2 \right] \\
&= 2\tilde{\mathbb{E}} \left(\psi_1(\tilde{\alpha}^T \mathbf{G}) \left[\int_0^1 \psi_1\{\bar{\alpha}^T \mathbf{G} + u(\tilde{\alpha} - \bar{\alpha})^T \mathbf{G}\} du \right] \{(\hat{\alpha} - \bar{\alpha})^T \mathbf{G}\}^2 \left\{ 1 - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \right\}^2 \right) \\
&\leq 2\tilde{\mathbb{E}} \left[\psi_1(\bar{\alpha}^T \mathbf{G}) \exp(C_3 C_4 r_\alpha) \left\{ \int_0^1 \psi_1(\bar{\alpha}^T \mathbf{G}) \exp(C_3 C_4 r_\alpha) du \right\} \{(\hat{\alpha} - \bar{\alpha})^T \mathbf{G}\}^2 \left\{ 1 - \frac{R}{\pi(\mathbf{X}; \bar{\gamma})} \right\}^2 \right] \\
&\leq 2C_2^2 C_4^2 \exp(2C_3 C_4 r_\alpha) \{1 + \exp(-B_0)\}^2 \|\hat{\alpha} - \bar{\alpha}\|_1^2 \\
&= o_p(1).
\end{aligned}$$

$$\begin{aligned}
\Delta_1^3 &= \tilde{\mathbb{E}} \left[R \{Y - \psi(\bar{\alpha}^T \mathbf{G})\}^2 \left\{ \frac{1}{\pi(\mathbf{X}; \hat{\gamma})} - \frac{1}{\pi(\mathbf{X}; \bar{\gamma})} \right\}^2 \right] \\
&= -2\tilde{\mathbb{E}} \left(R w(\mathbf{X}; \bar{\gamma})^2 [\exp\{-(\tilde{\gamma} - \bar{\gamma})^T \mathbf{F}\} - 1] \exp\{-(\tilde{\gamma} - \bar{\gamma})^T \mathbf{F}\} \times \{Y - \psi(\bar{\alpha}^T \mathbf{G})\}^2 (\hat{\gamma} - \bar{\gamma})^T \mathbf{G} \right) \\
&\leq 2\{1 + \exp(\eta_{01})\}^2 \exp(-B_0 + \eta_{01}) C_4 \|\hat{\gamma} - \bar{\gamma}\|_1 \tilde{\mathbb{E}} \left[R w(\mathbf{X}; \bar{\gamma}) \{Y - \psi(\bar{\alpha}^T \mathbf{G})\}^2 \right] \\
&\leq 2\{1 + \exp(\eta_{01})\}^2 \exp(-B_0 + \eta_{01}) C_4 \|\hat{\gamma} - \bar{\gamma}\|_1 \left(\mathbb{E} \left[R w(\mathbf{X}; \bar{\gamma}) \{Y - \psi(\bar{\alpha}^T \mathbf{G})\}^2 \right] + (D_0^2 + D_0 D_1) \lambda_0 \right) \\
&= o_p(1).
\end{aligned}$$

By mean value theorem,

$$\begin{aligned}
\Delta_2^3 &= \tilde{\mathbb{E}} \left[R\{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\bar{\alpha}^T \mathbf{G})\}^2 \left\{ \frac{1}{\pi(\mathbf{X}; \bar{\gamma})} - \frac{1}{\pi(\mathbf{X}; \hat{\gamma})} \right\}^2 \right] \\
&= \tilde{\mathbb{E}} \left(R\{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\bar{\alpha}^T \mathbf{G})\}^2 w(\mathbf{X}; \bar{\gamma})^2 [1 - \exp\{-(\hat{\gamma} - \bar{\gamma})^T \mathbf{F}\}]^2 \right) \\
&\leq \{1 + \exp(\eta_{01})\}^2 \exp(-2B_0) \tilde{\mathbb{E}} [\{\psi(\hat{\alpha}^T \mathbf{G}) - \psi(\bar{\alpha}^T \mathbf{G})\}^2] \\
&= 2\{1 + \exp(\eta_{01})\}^2 \exp(-2B_0) \tilde{\mathbb{E}} [\{\psi(\tilde{\alpha}^T \mathbf{G}) - \psi(\bar{\alpha}^T \mathbf{G})\} \psi_1(\tilde{\alpha}^T \mathbf{G}) (\hat{\alpha} - \bar{\alpha})^T \mathbf{G}] \\
&= 2\{1 + \exp(\eta_{01})\}^2 \exp(-2B_0) \times \tilde{\mathbb{E}} \left(\psi_1(\tilde{\alpha}^T \mathbf{G}) \left[\int_0^1 \psi_1\{\tilde{\alpha}^T \mathbf{G} + u(\tilde{\alpha} - \bar{\alpha})^T \mathbf{G}\} du \right] \{(\hat{\alpha} - \bar{\alpha})^T \mathbf{G}\}^2 \right) \\
&\leq 2C_2^2 C_4^2 \{1 + \exp(\eta_{01})\}^2 \exp(2C_3 C_4 r_\alpha - 2B_0) \|\hat{\alpha} - \bar{\alpha}\|_1^2 \\
&= o_p(1).
\end{aligned}$$

$$\begin{aligned}
\Delta_3^3 &= \tilde{\mathbb{E}} [\{\psi(\beta^{*\top} \mathbf{Z}) - \psi(\hat{\beta}^T \mathbf{Z})\}^2] \\
&= 2\tilde{\mathbb{E}} \left[\{\psi(\tilde{\beta}^T \mathbf{Z}) - \psi(\beta^{*\top} \mathbf{Z})\} \psi_1(\tilde{\beta}^T \mathbf{Z}) (\hat{\beta} - \beta^*)^T \mathbf{Z} \right] \\
&= 2\tilde{\mathbb{E}} \left(\psi_1(\tilde{\beta}^T \mathbf{Z}) \left[\int_0^1 \psi_1\{\beta^{*\top} \mathbf{Z} + u(\tilde{\beta} - \beta^*)^T \mathbf{Z}\} du \right] \{(\hat{\beta} - \beta^*)^T \mathbf{Z}\}^2 \right) \\
&\leq 2C_5^2 C_6^2 \{1 + \exp(\eta_{01})\}^2 \exp(2C_3 C_5 \|\hat{\beta} - \beta^*\|_1) \|\hat{\beta} - \beta^*\|_1^2 \\
&= o_p(1).
\end{aligned}$$

Therefore, $\hat{\Lambda} \xrightarrow{\mathbb{P}} \Lambda$. Then, by continuous mapping theorem, $\hat{\Gamma}^{-1} \xrightarrow{\mathbb{P}} \Gamma^{-1}$. Thus, by continuous mapping theorem again, $\hat{\Sigma} \xrightarrow{\mathbb{P}} \Sigma$.

S8 Extension to the setting of stratified sampling

S8.1 Proof of Proposition 3

Let

$$\tau^s(\mathbf{O}, \alpha, \beta, \cdot) = \frac{1}{N} \sum_{i=1}^n \frac{1}{(\cdot)} \{Y_i - \psi(\alpha^T \mathbf{G}_i)\} \mathbf{Z}_i + \frac{1}{N} \sum_{i=1}^N \{\psi(\alpha^T \mathbf{G}_i) - \psi(\beta^T \mathbf{Z}_i)\} \mathbf{Z}_i. \quad (\text{S8.39})$$

We have

$$\begin{aligned} & -\tau^s\{\mathbf{O}, \bar{\alpha}, \beta^*, \pi(\mathbf{X}; \bar{\gamma})\} = -\tau^s\{\mathbf{O}, \bar{\alpha}, \beta^*, \pi^*(\mathbf{X})\} = \tau^s\{\mathbf{O}, \hat{\alpha}, \hat{\beta}^s, \hat{\pi}(\mathbf{X})\} - \tau^s\{\mathbf{O}, \bar{\alpha}, \beta^*, \pi^*(\mathbf{X})\} \\ &= \frac{1}{N} \sum_{i=1}^N [\{\psi(\hat{\alpha}^T \mathbf{G}_i) - \psi(\hat{\beta}^{sT} \mathbf{Z}_i)\} \mathbf{Z}_i] - \frac{1}{N} \sum_{i=1}^N [\{\psi(\hat{\alpha}^T \mathbf{G}_i) - \psi(\beta^{*T} \mathbf{Z}_i)\} \mathbf{Z}_i] \\ &+ \frac{1}{N} \sum_{i=1}^N [\{\psi(\hat{\alpha}^T \mathbf{G}_i) - \psi(\beta^{*T} \mathbf{Z}_i)\} \mathbf{Z}_i] - \frac{1}{N} \sum_{i=1}^N [\{\psi(\bar{\alpha}^T \mathbf{G}_i) - \psi(\beta^{*T} \mathbf{Z}_i)\} \mathbf{Z}_i] \\ &+ \frac{1}{n} \sum_{i=1}^n \{Y_i - \psi(\hat{\alpha}^T \mathbf{G}_i)\} \mathbf{Z}_i - \frac{1}{n} \sum_{i=1}^n \{Y_i - \psi(\bar{\alpha}^T \mathbf{G}_i)\} \mathbf{Z}_i, \end{aligned}$$

denoted as $\Delta_5^0 + \Delta_5^1 + \Delta_5^2$.

We first deal with Δ_5^0 . Let Z_{ij} denote the j -th co-ordinate of the i -th sample \mathbf{Z}_i .

Then, we consider the j -th coordinate of Δ_5^0 , $\Delta_{5,j}^0$

$$\begin{aligned} \Delta_{5,j}^0 &= \frac{1}{N} \sum_{i=1}^N [\{\psi(\hat{\alpha}^T \mathbf{G}_i) - \psi(\hat{\beta}^{sT} \mathbf{Z}_i)\} Z_{ij}] - \frac{1}{N} \sum_{i=1}^N [\{\psi(\hat{\alpha}^T \mathbf{G}_i) - \psi(\beta^{*T} \mathbf{Z}_i)\} Z_{ij}] \\ &= \frac{1}{N} \sum_{i=1}^N [\{\psi(\beta^{*T} \mathbf{Z}_i) - \psi(\hat{\beta}^{sT} \mathbf{Z}_i)\} Z_{ij}] \end{aligned}$$

$$\begin{aligned}
 &= \frac{1}{N} \sum_{i=1}^N \{\psi_1(\tilde{\beta}_j^T \mathbf{Z}_i) Z_{ij} (\beta^* - \hat{\beta}^s)^T \mathbf{Z}_i\} \\
 &= \frac{1}{N} \sum_{i=1}^N \{\psi_1(\beta^{*T} \mathbf{Z}_i) Z_{ij} (\beta^* - \hat{\beta}^s)^T \mathbf{Z}_i\} + O_p(\|\beta^* - \hat{\beta}^s\|_1^2),
 \end{aligned}$$

where $\tilde{\beta}_j = u_j \beta^* + (1 - u_j) \hat{\beta}^s$ for some $u_j \in (0, 1)$. The last equality holds since $\psi_1(\beta^{*T} \mathbf{Z}_i) \exp(-C_3 C_5 \|\beta^* - \hat{\beta}^s\|_1) \leq \psi_1(\tilde{\beta}_j^T \mathbf{Z}_i) \leq \psi_1(\beta^{*T} \mathbf{Z}_i) \exp(C_3 C_5 \|\beta^* - \hat{\beta}^s\|_1)$, which leads to $\psi_1(\tilde{\beta}_j^T \mathbf{Z}_i) = \psi_1(\beta^{*T} \mathbf{Z}_i) + O_p(\|\beta^* - \hat{\beta}^s\|_1)$.

Then we deal with Δ_5^1 and Δ_5^2 together. Consider the j -th coordinate of $\Delta_5^1 + \Delta_5^2$:

$$\begin{aligned}
 &|\Delta_{5,j}^1 + \Delta_{5,j}^2| \\
 &= \left| \frac{1}{n} \sum_{i=1}^n [\{\psi(\tilde{\alpha}^T \mathbf{G}_i) - \psi(\hat{\alpha}^T \mathbf{G}_i)\} Z_{ij}] + \frac{1}{N} \sum_{i=1}^N [\{\psi(\hat{\alpha}^T \mathbf{G}_i) - \psi(\tilde{\alpha}^T \mathbf{G}_i)\} Z_{ij}] \right| \\
 &= \left| \frac{1}{n} \sum_{i=1}^n \{\psi_1(\tilde{\alpha}_j^T \mathbf{G}_i) (\tilde{\alpha} - \hat{\alpha})^T \mathbf{G}_i Z_{ij}\} + \frac{1}{N} \sum_{i=1}^N \{\psi_1(\tilde{\alpha}_j^T \mathbf{G}_i) (\hat{\alpha} - \tilde{\alpha})^T \mathbf{G}_i Z_{ij}\} \right| \\
 &\leq \|\hat{\alpha} - \tilde{\alpha}\|_1 \left\| \frac{1}{n} \sum_{i=1}^n \{\psi_1(\tilde{\alpha}_j^T \mathbf{G}_i) \mathbf{G}_i Z_{ij}\} - \frac{1}{N} \sum_{i=1}^N \{\psi_1(\tilde{\alpha}_j^T \mathbf{G}_i) \mathbf{G}_i Z_{ij}\} \right\|_{\infty},
 \end{aligned}$$

where $\tilde{\alpha}_j = u_j \hat{\alpha} + (1 - u_j) \tilde{\alpha}$, for some $0 < u_j < 1$, $j = 0, \dots, m-1$. Consider the k -th coordinate of \mathbf{G}_i . For technical convenience, we assume that n is divisible by N . Let

$$V_{ijk} \triangleq \left\{ \frac{n}{N} - 1 \right\} \psi_1(\tilde{\alpha}_j^T \mathbf{G}_i) G_{ik} Z_{ij} + \frac{n}{N} \sum_{s=n+(N/n-1)(i+1)+1}^{n+i(N/n-1)} \psi_1(\tilde{\alpha}_j^T \mathbf{G}_s) G_{sk} Z_{sj},$$

then, \exists constant $C_7 > 0$, such that $|V_{ijk}| \leq C_7$. Moreover, we have $\mathbb{E}(V_{ijk}) = 0$;

therefore, $\mathbb{E}(V_{ijk}^2) = \text{Var}(V_{ijk}) = \{(n-N)^2/N^2 + n^2/N^2(N/n-1)\} \text{Var}\{\psi_1(\tilde{\alpha}_j^T \mathbf{G}) \mathbf{G}_k Z_j\} \leq C_8$ for some constant $C_8 > 0$. Let $t = \sqrt{C_8 \ln\{(q+1)/\epsilon\}/n}$, since $\ln(q+1) = o(n)$, $\exists n$ large enough, such that $t^2/C_8 \leq 3t/C_7$, then by Bernstein's inequality, we have

$$\mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n V_{ijk}\right| \geq t\right) = \mathbb{P}\left(\left|\sum_{i=1}^n V_{ijk}\right| \geq nt\right) \leq \exp\left(-\frac{nt^2}{C_8}\right) = \frac{\epsilon}{q+1}. \quad (\text{S8.40})$$

Let V_j denote $\sup_{k=0, \dots, q} |1/n \sum_{i=1}^n V_{ijk}|$, and

$$V_j < t \Rightarrow \frac{V_j \sqrt{n}}{\sqrt{C_8 \ln\{e(q+1)\}}} \leq \sqrt{1 - \ln(\epsilon)}.$$

It follows that

$$\begin{aligned} & \mathbb{P}\left[\frac{V_j \sqrt{n}}{\sqrt{C_8 \ln\{e(q+1)\}}} \leq \sqrt{1 - \ln(\epsilon)}\right] \geq \mathbb{P}(V_j < t) \\ & = 1 - \mathbb{P}(V_j \geq t) \geq 1 - \sum_{k=0}^q \mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n V_{ijk}\right| \geq t\right) \\ & = 1 - \epsilon. \end{aligned}$$

Therefore, $\frac{V_j \sqrt{n}}{\sqrt{C_8 \ln\{e(q+1)\}}} = O_p(1)$, for $j = 0, \dots, m-1$, it follows that

$$\left\| \frac{1}{n} \sum_{i=1}^n \left\{ \frac{1}{\pi^s(\tilde{\gamma})} \psi_1(\tilde{\alpha}_j^T \mathbf{G}_i) \mathbf{G}_i Z_{ij} \right\} - \frac{1}{N} \sum_{i=1}^N \left\{ \psi_1(\tilde{\alpha}_j^T \mathbf{G}_i) \mathbf{G}_i Z_{ij} \right\} \right\|_{\infty} = V_j = O_p(\sqrt{\ln(q+1)/n}).$$

Therefore,

$$|\Delta_{5,j}^1 + \Delta_{5,j}^2| \leq O_p(\sqrt{\ln(q+1)/n})O_p(|S_{\bar{\alpha}}|\sqrt{\ln(q+1)/n}) = o_p(1/\sqrt{n}),$$

it follows that $\Delta_5^1 + \Delta_5^2 = o_p(1/\sqrt{n})$. Hence, by central limit theorem and continuous mapping theorem, we have

$$\sqrt{n}(\hat{\beta}^s - \beta^*) \xrightarrow{d} \sqrt{n}\Gamma^{-1}\tau^s\{\mathbf{O}, \bar{\alpha}, \beta^*, \pi^*(\mathbf{X})\} \sim N(0, \Sigma^s).$$

S8.2 Variance comparison

Because the following relationship holds,

$$\begin{aligned} \frac{1}{n}\mathbf{\Lambda}^s &= \frac{1}{n}\mathbb{E} \left[\left\{ Y - \frac{N-n}{N}\psi(\bar{\alpha}^T\mathbf{G}) - \frac{n}{N}\psi(\beta^{*T}\mathbf{Z}) \right\}^2 \mathbf{Z}\mathbf{Z}^T \right] + \frac{(N-n)}{N^2}\mathbb{E} [\{\psi(\bar{\alpha}^T\mathbf{G}) - \psi(\beta^{*T}\mathbf{Z})\}^2 \mathbf{Z}\mathbf{Z}^T] \\ &= \frac{1}{n}\mathbb{E} \left[\{Y - \psi(\bar{\alpha}^T\mathbf{G})\}^2 \mathbf{Z}\mathbf{Z}^T + \left(\frac{n}{N}\right)^2 \{\psi(\bar{\alpha}^T\mathbf{G}) - \psi(\beta^{*T}\mathbf{Z})\}^2 \mathbf{Z}\mathbf{Z}^T \right] \\ &\quad + \frac{2}{N}\mathbb{E} [\{Y - \psi(\bar{\alpha}^T\mathbf{G})\}\{\psi(\bar{\alpha}^T\mathbf{G}) - \psi(\beta^{*T}\mathbf{Z})\} \mathbf{Z}\mathbf{Z}^T] + \frac{N-n}{N^2}\mathbb{E} [\{\psi(\bar{\alpha}^T\mathbf{G}) - \psi(\beta^{*T}\mathbf{Z})\}^2 \mathbf{Z}\mathbf{Z}^T] \\ &= \frac{1}{n}\mathbb{E} [\{Y - \psi(\bar{\alpha}^T\mathbf{G})\}^2 \mathbf{Z}\mathbf{Z}^T] + \frac{2}{N}\mathbb{E} [\{Y - \psi(\bar{\alpha}^T\mathbf{G})\}\{\psi(\bar{\alpha}^T\mathbf{G}) - \psi(\beta^{*T}\mathbf{Z})\} \mathbf{Z}\mathbf{Z}^T] \\ &\quad + \frac{1}{N}\mathbb{E} [\{\psi(\bar{\alpha}^T\mathbf{G}) - \psi(\beta^{*T}\mathbf{Z})\}^2 \mathbf{Z}\mathbf{Z}^T] \\ &= \frac{1}{N}\mathbf{\Lambda}, \end{aligned}$$

we obtain

$$\boldsymbol{\Sigma}^s/n = \boldsymbol{\Gamma}^{-1} \left(\frac{1}{n} \boldsymbol{\Lambda}^s \right) \boldsymbol{\Gamma}^{-1} = \boldsymbol{\Gamma}^{-1} \left(\frac{1}{N} \boldsymbol{\Lambda} \right) \boldsymbol{\Gamma}^{-1} = \boldsymbol{\Sigma}/N.$$

S9 Details of numerical implementation and simulation

In this section, we provide details of the numerical implementation and simulation.

We consider the estimators of population mean for $\mathbf{Z} = 1$, regression coefficients in the mean model for $\mathbf{Z} = X_1$ and $\mathbf{Z} = \mathbf{X}$, respectively, and $\psi(\cdot)$ is assumed to be the identity function.

S9.1 Data generating process

Throughout the simulation, we generate the covariates \mathbf{X} as follows. We first generate a random vector from $N(0, \boldsymbol{\Sigma})$, where the variance matrix $\boldsymbol{\Sigma} \in \mathbb{R}^{3 \times 3}$ has elements $\boldsymbol{\Sigma}_{i,j}$ defined as $2^{-|i-j|}$ for $i, j = 1, 2, 3$. Then we clamp each of its coordinates within $[-3, 3]$ to obtain (X_1, X_2, X_3) and $\mathbf{X} = (1, X_1, X_2, X_3)$. In addition, the data source indicator R follows a Bernoulli distribution with success probability $\pi(\mathbf{X})$, where $\pi(\mathbf{X}) = \{1 + \exp(-\gamma^\top \mathbf{F})\}^{-1}$, the parameter $\gamma = (-1.5, -0.8, -0.2, 0.3, 0, \dots, 0)^\top$ and the basis functions \mathbf{F} are described in Section 6.

Study I. We first focus on the estimation of the population mean and consider two data-generating mechanisms:

- **Case 1.** The outcome $Y = -0.2 + 0.1\tilde{X}_1 + 0.4\tilde{X}_2 + 0.7\tilde{X}_3 + \epsilon$, where $\epsilon \sim \mathcal{N}(0, 0.1)$

and $\tilde{X}_j = X_j \cdot |X_j|^{0.1} + X_j \cdot |X_j|^{0.3} + X_j \cdot |X_j|^{0.5}$, for $j = 1, 2, 3$. We set $\mathbf{Z} = 1$.

- **Case 2.** The outcome $Y = -0.2 + 0.1\tilde{X}_1 + 0.4\tilde{X}_2 + 0.7\tilde{X}_3 + \epsilon$, where $\epsilon \sim \mathcal{N}(0, 0.1)$ and $\tilde{X}_j = |X_j| \exp(|X_j|^{0.1} + |X_j|^{0.3})$ for $j = 1, 2, 3$. We set $\mathbf{Z} = 1$.

In many scenarios, estimating the conditional mean given a subset of variables in \mathbf{X} garners statistical interest. Accordingly, we design experiments in Study II to evaluate the performance of the proposed estimator in such setups.

Study II. We further consider three additional cases for estimating regression coefficients in the conditional mean outcome model.

- **Case 3.** The outcome $Y = 0.4\tilde{X}_1 + 0.2\tilde{X}_2 + \epsilon$, where $\epsilon \sim \mathcal{N}(0, 0.1)$, $\tilde{X}_1 = X_1 + X_1^2$ and $\tilde{X}_2 = \cos(\pi/9 \cdot X_1 \cdot X_3)$. We set $\mathbf{Z} = X_1$.
- **Case 4.** The outcome $Y = 0.4\tilde{X}_1 + 0.2\tilde{X}_2 + \epsilon$, where $\epsilon \sim \mathcal{N}(0, 0.1)$, $\tilde{X}_1 = X_1 \cdot \mathbb{I}\{X_1 > 0\} \sqrt{|X_1|}$ and $\tilde{X}_2 = X_1 \cdot X_2$. We set $\mathbf{Z} = X_1$.
- **Case 5.** The outcome $Y = -0.2 + 0.1\tilde{X}_1 + 0.4\tilde{X}_2 + 0.7\tilde{X}_3 + \epsilon$, where $\epsilon \sim \mathcal{N}(0, 0.1)$, $\tilde{X}_1 = X_1 \cdot X_2$, $\tilde{X}_2 = X_2 \cdot X_3$ and $\tilde{X}_3 = X_1 \cdot X_3$. We set $\mathbf{Z} = \mathbf{X}$.

Cases 3 and 4 involve \mathbf{Z} as a specific covariate X_1 , while Case 5 involves \mathbf{Z} as the full set of covariates \mathbf{X} . In addition, for all Cases 1–5, OR models are misspecified.

S9.2 Implementation details

We compare the proposed method $\text{AIPW}_{\text{RCAL}}$ with IPW, AIPW_{RML} and AIPW_{CF} introduced in Section 6. Both the Lasso-regularized calibrated and maximum likeli-

hood estimators for the PS and OR models can be implemented using the R package RCAL (Tan and Sun, 2020). We employ 5-fold cross-fitting to select the optimal tuning parameters. In addition, by equation (2.2), $\beta^* = \mathbb{E}(\mathbf{Z}\mathbf{Z}^T)^{-1}\mathbb{E}(Y\mathbf{Z})$. Thus the true value of β^* is calculated as $\tilde{\mathbb{E}}(\mathbf{Z}\mathbf{Z}^T)^{-1}\tilde{\mathbb{E}}(Y\mathbf{Z})$ through a simulation with a sample size of 100,000. For $\mathbf{Z} = X_1$ and $\mathbf{Z} = \mathbf{X}$, we denote $\beta^* = \beta_1$ and $\beta^* = (\beta_0, \beta_1, \beta_2, \beta_3)^T$, respectively.

S9.3 Summary of results

We present and analyze the simulation results for Study I and Study II. In the following tables, we compare various methods in terms of Bias, $\sqrt{\text{Var}}$, $\sqrt{\text{EVar}}$, CP90, and CP95. As discussed in the paragraph below Theorem 1, under high-dimensional settings, the IPW estimator is not \sqrt{N} -consistent, and its asymptotic normality is not well established. Therefore, we do not report its numerical results for $\sqrt{\text{EVar}}$, CP90, and CP95.

Results for Study I. Table 1 shows the numerical results for the estimation of population mean $\mathbb{E}(Y)$. From the table, the proposed method AIPW_{RCAL} has the smallest $\sqrt{\text{Var}}$ and $\sqrt{\text{EVar}}$, and Bias. Moreover, CP90 and CP95 of the proposed method are more aligned with their nominal values of 0.90 and 0.95, respectively. This indicates the effectiveness of the proposed method in terms of estimating the population mean.

Figure 1 depicts the box plots of the estimates for Case 1 and Case 2, where the

Table 1: Summary of estimates of population mean for Study I.

	Case 1			Case 2		
	AIPW _{RCAL}	AIPW _{RML}	AIPW _{CF}	AIPW _{RCAL}	AIPW _{RML}	AIPW _{CF}
Bias	0.004	0.004	0.006	0.002	-0.006	0.031
$\sqrt{\text{Var}}$	0.078	0.078	0.079	0.139	0.143	0.166
$\sqrt{\text{EVar}}$	0.079	0.079	0.082	0.140	0.143	0.161
CP90	0.904	0.908	0.918	0.898	0.886	0.896
CP95	0.948	0.954	0.956	0.948	0.946	0.956

horizontal line indicates the true value. In both cases, our method AIPW_{RCAL} exhibits the smallest biases, interquartile ranges, and whiskers, indicating the smallest variances compared to the other methods. In addition, AIPW_{CF} shows more outliers than the other two methods, which is more apparent in the results of Study II, suggesting that cross-fitting may cause instability for the estimates.

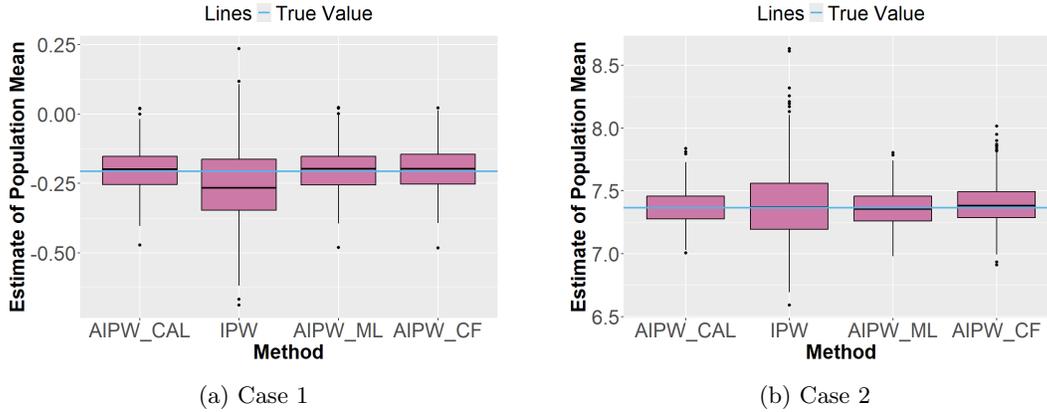


Figure 1: Box plots for estimates of population mean.

Results for Study II. The simulation results for Cases 3–5 are presented in Tables 2–3, and the corresponding box plots are displayed in Figures 2–3. We observe similar patterns as those in Cases 1 and 2: the proposed method performs well in terms of all metrics.

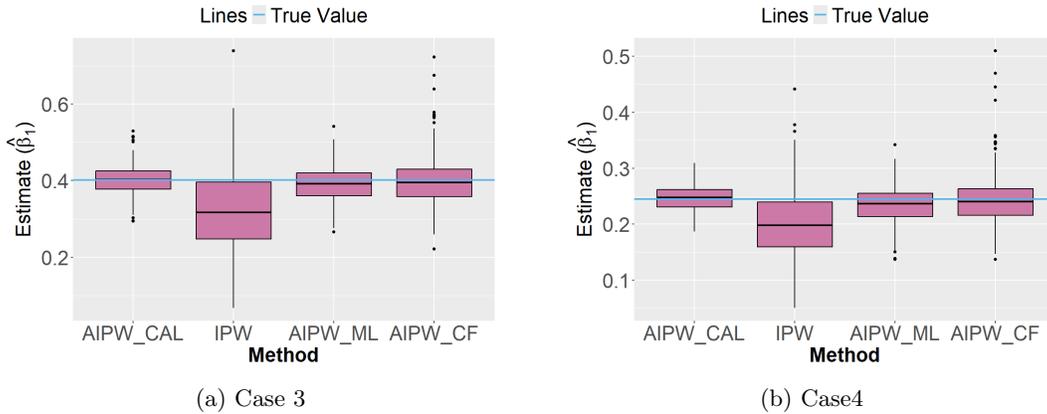
Table 2: Summary of estimates of β_1 in Cases 3 and 4.

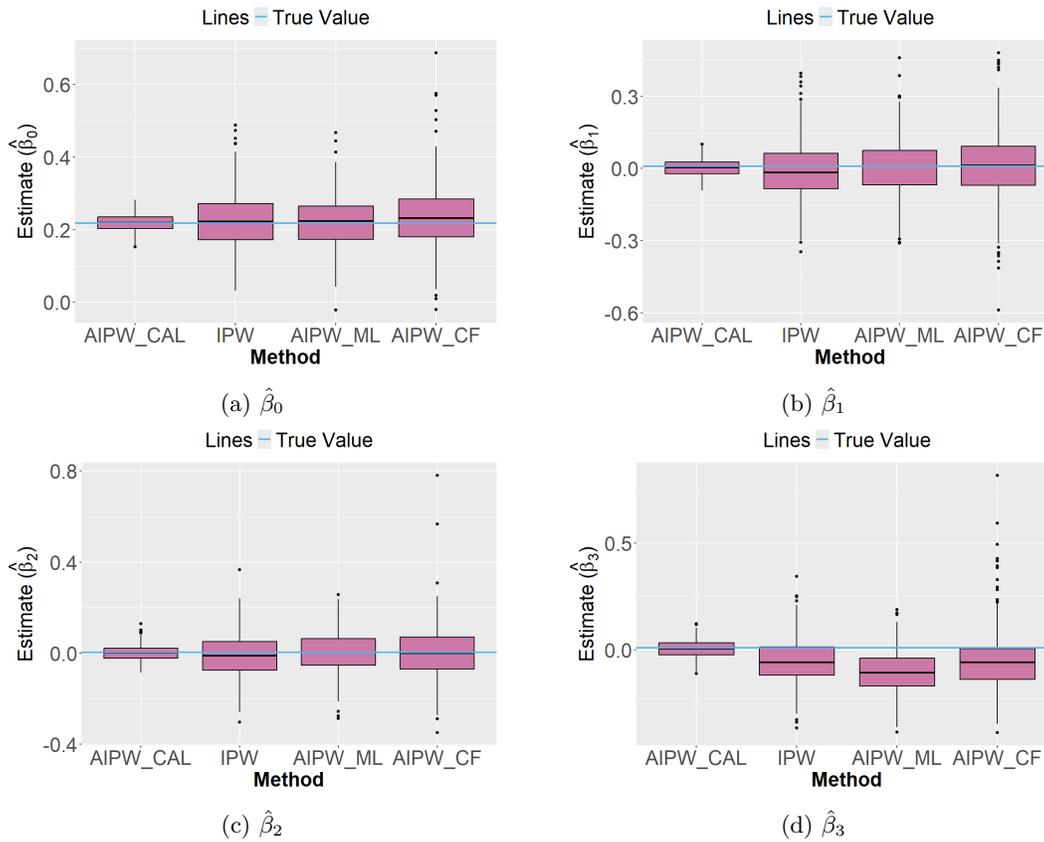
	Case 3			Case 4		
	AIPW _{RCAL}	AIPW _{RML}	AIPW _{CF}	AIPW _{RCAL}	AIPW _{RML}	AIPW _{CF}
Bias	0.000	-0.011	-0.004	0.002	-0.010	-0.002
$\sqrt{\text{Var}}$	0.036	0.043	0.058	0.023	0.031	0.042
$\sqrt{\text{EVar}}$	0.036	0.037	0.056	0.021	0.025	0.042
CP90	0.884	0.824	0.840	0.866	0.774	0.838
CP95	0.946	0.886	0.900	0.930	0.834	0.906

Table 3: Summary of estimates of $\beta_0, \beta_1, \beta_2, \beta_3$ in Case 5

	β_0			β_1		
	AIPW _{RCAL}	AIPW _{RML}	AIPW _{CF}	AIPW _{RCAL}	AIPW _{RML}	AIPW _{CF}
Bias	0.001	0.001	0.013	-0.007	-0.004	0.000
$\sqrt{\text{Var}}$	0.024	0.068	0.084	0.036	0.114	0.139
$\sqrt{\text{EVar}}$	0.025	0.055	0.077	0.036	0.089	0.131
CP90	0.910	0.800	0.846	0.886	0.774	0.830
CP95	0.958	0.874	0.914	0.942	0.842	0.900

	β_2			β_3		
	AIPW _{RCAL}	AIPW _{RML}	AIPW _{CF}	AIPW _{RCAL}	AIPW _{RML}	AIPW _{CF}
Bias	-0.001	0.002	0.001	-0.006	0.114	-0.060
$\sqrt{\text{Var}}$	0.033	0.086	0.111	0.040	0.095	0.132
$\sqrt{\text{EVar}}$	0.034	0.075	0.108	0.041	0.095	0.133
CP90	0.920	0.818	0.866	0.898	0.592	0.774
CP95	0.958	0.902	0.938	0.940	0.682	0.842

Figure 2: Box plots of estimates of β_1 in Cases 3 and 4.

Figure 3: Box plots of estimates of $\beta_0, \beta_1, \beta_2, \beta_3$ in Case 5

S10 Details of the application

S10.1 Pre-processing details of the community and crime dataset

We pre-process the data in following steps:

Step 1. remove 22 covariates missing 84% of data and 2 variables missing roughly 59% of data;

Step 2. remove covariates with weak linear relationships to the response `ViolentCrimesPerPop`

based on their correlation coefficients.

Step 3. remove covariates that exhibit multi-collinearity based on their values of variance inflation factors.

After the process, we obtain 1993 observations of 26 covariates.

S10.2 Test results of the covariate shift

Kernel two-sample test with maximum mean discrepancy

- Kernel: $\exp(-\|\cdot\|_2^2)$
- MMD: 0.39227
- P-value: 0.001

Bootstrap KS-tests for univariate covariates

Covariate	Bootstrap-KS P-value	KS-test Statistic	KS-test Approximate P-value
racePctHisp	0.000	0.335	0.000
pctWwage	0.000	0.230	0.000
pctWInvInc	0.000	0.330	0.000
blackPerCap	0.000	0.421	0.000
PctLess9thGrade	0.010	0.119	0.010
PctUnemployed	0.000	0.231	0.000
PctOccupManu	0.000	0.205	0.000

MalePctDivorce	0.000	0.373	0.000
MalePctNevMarr	0.000	0.202	0.000
PctTeen2Par	0.000	0.289	0.000
PctIlleg	0.000	0.200	0.000
NumImmig	0.000	0.267	0.000
PctImmigRec10	0.001	0.141	0.001
PctHousLess3BR	0.000	0.218	0.000
MedNumBR	0.000	0.138	0.000
HousVacant	0.000	0.184	0.000
PctHousOccup	0.000	0.195	0.000
PctHousOwnOcc	0.000	0.285	0.000
PctVacantBoarded	0.797	0.040	0.923
PctHousNoPhone	0.000	0.373	0.000
PctWofFullPlumb	0.000	0.146	0.000
RentLowQ	0.000	0.529	0.000
MedRentPctHousInc	0.062	0.089	0.099
NumInShelters	0.022	0.087	0.113
NumStreet	0.001	0.101	0.043
PopDens	0.000	0.266	0.000

Table 4: Bootstrap KS-tests for univariate covariates

S10.3 Design of basis functions

In this application, we design basis functions in the following way: Given $\{X_{ij}\}_{j=1}^N$, i.e., N samples of the i -th coordinate of \mathbf{X} , let $\{\xi_{ij}\}_{j=1}^{n_k}$ be the n_k points equally spaced within the $[-a_i, b_i]$, where $a_i = \min_{j=1, \dots, N} X_{ij}$ and $b_i = \max_{j=1, \dots, N} X_{ij}$. Let $f_{ij}(\mathbf{X})$ denote $(X_i - \xi_{ij})_+$, $i = 1, \dots, d$; $j = 1, \dots, n_k$. Let $\mathbf{F} = \{1, f_{11}(\mathbf{X}), \dots, f_{1n_k}(\mathbf{X}), \dots, f_{d1}(\mathbf{X}), \dots, f_{dn_k}(\mathbf{X})\}^T$ be the basis functions of the PS model, and let $\mathbf{G} = \{\mathbf{F}^T, (\mathbf{Z} \otimes \mathbf{F})^T\}^T$. We choose $n_k = 4$ in this application.

References

- Bühlmann, P. and S. Van De Geer (2011). *Statistics for High-dimensional Data: Methods, Theory and Applications*. Springer Science & Business Media.
- Chakraborty, A., J. Lu, T. T. Cai, and H. Li (2019). High dimensional m-estimation with missing outcomes: A semi-parametric framework. *arXiv preprint arXiv:1911.11345*.
- Chen, K. and Y. Zhang (2023). Enhancing efficiency and robustness in high-dimensional linear regression with additional unlabeled data. *arXiv preprint arXiv:2311.17685*.
- Deng, S., Y. Ning, J. Zhao, and H. Zhang (2023). Optimal and safe estimation for high-dimensional semi-supervised learning. *Journal of the American Statistical Association*, 1–12.
- Ghosh, S. and Z. Tan (2022). Doubly robust semiparametric inference using regularized calibrated estimation with high-dimensional data. *Bernoulli* 28, 1675–1703.
- Tan, Z. (2020a). Model-assisted inference for treatment effects using regularized calibrated estimation with

- high-dimensional data. *The Annals of Statistics* 48, 811–837.
- Tan, Z. (2020b). Regularized calibrated estimation of propensity scores with model misspecification and high-dimensional data. *Biometrika* 107, 137–158.
- Tan, Z. and B. Sun (2020). *RCAL: Regularized calibrated estimation*. R package version 2.0.
- Van der Vaart, A. W. (2000). *Asymptotic Statistics*. Cambridge University Press.
- Zhang, Y., A. Giessing, and Y.-C. Chen (2023). Efficient inference on high-dimensional linear models with missing outcomes. *arXiv preprint arXiv:2309.06429*.