

---

# Enhancing Treatment Effect Estimation: A Model Robust Approach Integrating Randomized Experiments and External Controls using the Double Penalty Integration Estimator (Supplementary Material)

---

Yuwen Cheng<sup>1</sup>

Lili Wu<sup>2</sup>

Shu Yang<sup>3</sup>

<sup>1</sup>Statistics Dept. North Carolina State University

<sup>2</sup>Microsoft Research NYC

<sup>3</sup>Statistics Dept. North Carolina State University

The supplementary material is structured as follows: Section S1 includes additional assumptions for the likelihood and density functions. Section S2 provides proofs for the main theorems. Section S3 presents the toy example mentioned in Section 2. Section S4 displays additional figures for the first simulation study in this section. The simulation codes are in the <https://github.com/yuwen997/simulation-codes>.

## S1 REGULARITY CONDITIONS

In this section, we provide the same assumptions on likelihood functions and density functions as those in Fan and Peng [2004], White [1982]. For a better understanding of these conditions, an informal summary of assumptions S1–S9 is provided here. Assumptions S1–S7 align with A1–A7 in White [1982], ensuring MLE estimator consistency and asymptotic normality in both misspecified and correct models. S8–S9 resemble F–G in Fan and Peng [2004], bounding  $f$  moments.

**Assumption S1** *The independent random vectors  $(p_i, Y_i)$ ,  $i = 1, \dots, N$ , have common joint distribution function  $G$  on  $\Upsilon$ , a measurable Euclidean space, with measurable Radon-Nikodym density  $g = dG/d\mu$ .*

**Assumption S2** *The family of distribution functions  $F(Y_1, p, \theta)$  has Radon-Nikodym densities  $f(y, p, \theta) = dF(y, p, \theta)/d\mu$  which are measurable in  $(y, p)$  for every  $\theta$  in  $\Theta$ , a compact subset of Euclidean space, and continuous in  $\theta$  for every  $(y, p)$  in  $\Upsilon$ .*

**Assumption S3** *(a)  $\mathbb{E}\{\log g(Y, p)\}$  exists and  $|\log f(y, p, \theta)| \leq m(y, p)$  for all  $\theta$  in  $\Theta$ , where  $m$  is integrable with respect to  $G$ ; (b)  $KLIC(g : f, \theta)$  has a unique minimum at  $\theta_*$  in  $\Theta$ .*

**Assumption S4**  *$\partial \log f(y, p, \theta)/\partial \theta_j$ ,  $j = 1, \dots, K$ , are measurable functions of  $(y, p)$  for each  $\theta$  in  $\Theta$  and continuously differentiable functions of  $\theta$  for each  $(y, p)$  in  $\Upsilon$ .*

**Assumption S5**  *$|\partial^2 \log f(y, p, \theta)/\partial \theta_i \partial \theta_j|$  and  $|\partial \log f(y, p, \theta)/\partial \theta_i \cdot \partial \log f(y, p, \theta)/\partial \theta_j|$ ,  $i, j = 1, \dots, K$  are dominated by functions integrable with respect to  $G$  for each  $\theta$  in  $\Theta$  and  $(y, p)$  in  $\Upsilon$ .*

**Assumption S6** *Define matrix*

$$A(\theta) = -\mathbb{E} \left\{ \frac{\partial^2 \log f(Y_1, p_1, \theta)}{\partial \theta_j \partial \theta_k} \right\} > 0,$$
$$B(\theta) = \mathbb{E} \left[ \left\{ \frac{\partial \log f(Y_1, p_1, \theta)}{\partial \theta} \right\} \left\{ \frac{\partial \log f(Y_1, p_1, \theta)}{\partial \theta} \right\}^T \right],$$

and (a)  $\theta_*$  is interior to  $\Theta$ ; (b)  $B(\theta_*)$  is nonsingular; (c)  $\theta_*$  is a regular point of  $A(\theta)$ .

**Assumption S7**  $|\partial \{ \partial f(y, p, \theta) / \partial \theta_i \cdot f(y, p, \theta) \} / \partial \theta_j|$ ,  $i, j = 1, \dots, K$  are dominated by functions integrable with respect to  $\mu$  for all  $\theta$  in  $\Theta$  and the minimal support of  $f(y, p, \theta)$  does not depend on  $\theta$ .

**Assumption S8** Define matrix

$$C(\theta) = A(\theta)B^{-1}(\theta)A(\theta).$$

Assume matrix  $A(\theta)$  and  $B(\theta)$  satisfy conditions

$$\begin{aligned} 0 < C_1 < \lambda_{\min} \{A(\theta)\} \leq \lambda_{\max} \{A(\theta)\} < C_2 < \infty \text{ for all } N, \\ 0 < C_1^* < \lambda_{\min} \{B(\theta)\} \leq \lambda_{\max} \{B(\theta)\} < C_2^* < \infty \text{ for all } N, \end{aligned}$$

and for  $j, k = 1, \dots, K$ ,

$$\mathbb{E}_\theta \left\{ \frac{\partial \log f(Y_1, p_1, \theta)}{\partial \theta_j} \frac{\partial \log f(Y_1, p_1, \theta)}{\partial \theta_k} \right\}^2 < C_3 < \infty$$

and

$$\mathbb{E}_\theta \left\{ \frac{\partial^2 \log f(Y_1, p_1, \theta)}{\partial \theta_j \partial \theta_k} \right\}^2 < C_4 < \infty.$$

**Assumption S9** There is a large enough open subset  $\omega_N$  of  $\Theta \in R^K$  which contains the parameter point  $\theta_*$ , such that for almost all  $(p_i, Y_i)$  the density admits all third derivatives  $\partial^3 f(Y_i, p_i, \theta) / \partial \theta_j \partial \theta_k \partial \theta_l$  for all  $\theta \in \omega_N$ . Furthermore, there are functions  $M_{jkl}$  such that

$$\left| \frac{\partial^3 \log f(Y_i, p_i, \theta)}{\partial \theta_j \partial \theta_k \partial \theta_l} \right| \leq M_{jkl}(Y_i, p_i)$$

for all  $\theta \in \omega_N$ , and

$$\mathbb{E}_\theta \{M_{jkl}^2(Y_i, p_i)\} < C_5 < \infty$$

for all  $K, N, j, k$  and  $l$ .

**Assumption S10** Define  $h(A, X, S | \beta, \delta) = \beta_{\text{int}} + \beta_A A + \beta_X^T p_\mu(X) + (1 - S)\delta^T p_b(X)$ .  $f, g$  satisfy that the minimizers  $\beta_*, \delta_*$  of  $\mathbb{E}\{Y - h(A, X, S | \beta, \delta)\}^2$  are also the minimizers  $\beta_*, \delta_*$  of  $KLIC(g : f, \beta, \delta) = \mathbb{E}[\log \{g(Y, p) / f(Y, p, \beta, \delta)\}]$ .

## S2 PROOF

In this section, we provide proof of Theorems 1-4. Define  $L(\theta) = L(\beta, \delta) = \sum_{i=1}^N \{\ln f(Y_i, p_i, \beta, \delta)\}$ . Subsequently, the penalized likelihood function is  $Q(\theta) = L(\theta) - N \sum_{i=1}^{K_1} P_{\lambda_1, i}(|\beta|) - N \sum_{i=1}^{K_2} P_{\lambda_2, i}(|\delta|)$ . Assume  $\beta_* = \begin{pmatrix} \beta_{*1} \\ \beta_{*2} \end{pmatrix}$ ,  $\delta_* = \begin{pmatrix} \delta_{*1} \\ \delta_{*2} \end{pmatrix}$  where  $\beta_{*1} \neq 0$  with  $s_1$  dimensions,  $\delta_{*1} \neq 0$  with  $s_2$  dimensions,  $\beta_{*2} = 0$  with  $K_1 - s_1$  dimensions and  $\delta_{*2} = 0$  with  $K_2 - s_2$  dimensions. Further, let  $\theta_* = \begin{pmatrix} \theta_{*1} \\ \theta_{*2} \end{pmatrix}$ , where  $\theta_{*1} = \begin{pmatrix} \beta_{*1} \\ \delta_{*1} \end{pmatrix} \neq 0$  with  $s = s_1 + s_2$  dimensions and  $\theta_{*2} = \begin{pmatrix} \beta_{*2} \\ \delta_{*2} \end{pmatrix} = 0$  with  $K - s$  dimensions.

### S2.1 THEOREM S1

Theorem S1 was previously demonstrated in Lorentz [1966] and Chen [2007], and we restate it here.

**Theorem S1** For any unknown function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , assuming function  $f(\cdot)$  is  $t$  times continuously differentiable. Let  $K = (q + 1)^d$  where  $x_1, \dots, x_d$  are at least up to power  $q$ , and let  $r^K(x)$  be the  $K$ -dimension power series basis function,  $R(x) = A_K r^K(x)$  where  $A_K$  is the matrix such that  $\mathbb{E}\{R(X) R^T(X)\} = \mathcal{I}$  where  $\mathcal{I}$  is the identity matrix. Then there is a  $K$ -vector  $\theta$  such that on the compact set  $\mathcal{X}$ ,  $\sup_{x \in \mathcal{X}} |f(x) - R^T(x)\theta| = O(K^{-t/d})$ .

## S2.2 PROOF FOR THEOREM 1

Under the combined dataset, the ANCOVA working model is

$$Y = \beta_{\text{int}} + \beta_A A + \beta_X^T p_\mu(X) + (1 - S)b_0(X).$$

We can rewrite the models as

$$\begin{aligned}\tilde{Y} &= \beta_{\text{int}} + \beta_A A + \beta_X^T p_\mu(X) \\ &= \beta^T p_\mu \\ &= \bar{\mu}_{A,1}(X; \beta),\end{aligned}$$

where  $\tilde{Y} := Y - (1 - S)b_0(X)$ . Note,  $\bar{\mu}_{A,1}(X; \beta) = \beta^T p_\mu$ . Remember  $\bar{\mu}_{A,1}(X) = \beta_0 + \beta_A A + \beta_X^T p_\mu(X)$ . The same in Proof of Theorem for Linear Models in Rosenblum and Van Der Laan [2009], the ordinary least squares estimate of  $\beta$  is asymptotically normal and converges in probability to the minimizer  $\tilde{\beta}_*$  of  $\mathbb{E}\{\tilde{Y} - \bar{\mu}_{A,1}(X; \beta)\}^2$ . Then

$$\begin{aligned}\mathbb{E}\{\tilde{Y} - \bar{\mu}_{A,1}(X; \beta)\}^2 &= \mathbb{E}\{\tilde{Y} - \mathbb{E}(\tilde{Y} | A, X, S) + \mathbb{E}(\tilde{Y} | A, X, S) - \bar{\mu}_{A,1}(X; \beta)\}^2 \\ &= \mathbb{E}\{\tilde{Y} - \mathbb{E}(\tilde{Y} | A, X, S)\}^2 + \mathbb{E}\left\{\mathbb{E}(\tilde{Y} | A, X, S) - \bar{\mu}_{A,1}(X; \beta)\right\}^2 \\ &= \mathbb{E}\{\tilde{Y} - \mathbb{E}(\tilde{Y} | A, X, S)\}^2 + \mathbb{E}\left[\mathbb{E}\left\{\mathbb{E}(\tilde{Y} | A, X, S) - \bar{\mu}_{A,1}(X; \beta)\right\}^2 \mid S\right] \\ &= \mathbb{E}\{\tilde{Y} - \mathbb{E}(\tilde{Y} | A, X, S)\}^2 + \mathbb{E}\left\{\mathbb{E}(\tilde{Y} | A, X, S = 1) - \bar{\mu}_{A,1}(X; \beta)\right\}^2 \mathbb{P}(S = 1) \\ &\quad + \mathbb{E}\left\{\mathbb{E}(\tilde{Y} | A = 0, X, S = 0) - \bar{\mu}_{0,1}(X; \beta)\right\}^2 \mathbb{P}(S = 0).\end{aligned}$$

Under HCs, by the definition of  $b_0(X)$ , we have

$$\begin{aligned}\mathbb{E}(\tilde{Y} | A = 0, X, S = 0) - \bar{\mu}_{0,1}(X; \beta) &= \mathbb{E}\{Y - b_0(X) | A = 0, X, S = 0\} - \bar{\mu}_{0,1}(X; \beta) \\ &= \mathbb{E}\{Y - \mathbb{E}(Y | A = 0, X, S = 0) + \bar{\mu}_{0,1}(X; \beta) | A = 0, X, S = 0\} - \bar{\mu}_{0,1}(X) \\ &= \mathbb{E}(Y | A = 0, X, S = 0) - \mathbb{E}(Y | A = 0, X, S = 0) + \bar{\mu}_{0,1}(X) - \bar{\mu}_{0,1}(X) \\ &= 0.\end{aligned}$$

Similar in [Wang et al., 2023], the  $\tilde{\beta}_*$  minimizing

$$\begin{aligned}\mathbb{E}\left\{\mathbb{E}(\tilde{Y} | A, X, S = 1) - \bar{\mu}_{A,1}(X; \beta)\right\}^2 &= \mathbb{E}\{\mathbb{E}(Y | A, X, S = 1) - \bar{\mu}_{A,1}(X; \beta)\}^2 \\ &= \mathbb{E}\{\mu_{A,1}(X) - \bar{\mu}_{A,1}(X; \beta)\}^2.\end{aligned}$$

On the other hand,  $\beta_*$  minimizes

$$\begin{aligned}\mathbb{E}\{[Y - \bar{\mu}_{A,1}(X; \beta)]^2 | S = 1\} &= \mathbb{E}\{[Y - \mu_{A,1}(X) + \mu_{A,1}(X) - \bar{\mu}_{A,1}(X; \beta)]^2 | S = 1\} \\ &= \mathbb{E}\{[Y - \mu_{A,1}(X)]^2 | S = 1\} + \mathbb{E}\{[\mu_{A,1}(X) - \bar{\mu}_{A,1}(X; \beta)]^2 | S = 1\} \\ &= \mathbb{E}\{[Y - \mu_{A,1}(X)]^2 | S = 1\} + \mathbb{E}\{\mu_{A,1}(X) - \bar{\mu}_{A,1}(X; \beta)\}^2.\end{aligned}$$

Therefore  $\tilde{\beta}_* = \beta_*$ , which both minimizes  $\mathbb{E}\{\mu_{A,1}(X) - \bar{\mu}_{A,1}(X; \beta)\}^2$ . Similar in [Wang et al., 2023], by the first formula of taking the first derivative of  $\mathbb{E}\{\mu_{A,1}(X) - \bar{\mu}_{A,1}(X; \beta)\}^2$ ,  $\beta_*$  satisfies  $\mathbb{E}\{\mu_{A,1}(X) - \bar{\mu}_{A,1}(X; \beta)\} = 0$ . That is,  $\beta_*$  satisfies  $\tau = \mathbb{E}\{\mu_{1,1}(X) - \mu_{0,1}(X)\} = \mathbb{E}\{\bar{\mu}_{1,1}(X) - \bar{\mu}_{0,1}(X)\} = \beta_{A*}$ . Similarly, in the REs, we have  $\tau = \beta_{A*}$ .

## S2.3 PROOF FOR THEOREM 2

Assume  $g(Y, p)$  is the true density function,  $f(Y, p, \theta)$  is our working density function. From Assumption S10, the selected  $f$  makes minimizing the KLIC equivalent to minimizing the least square to  $\mathbb{E}\{Y - h(A, X, S | \beta, \delta)\}^2$ , therefore,  $\theta_* = (\beta_*^T, \delta_*^T)^T$  is also the parameter which minimizes the Kullback-Leibler Information Criterion (KLIC),

$$KLIC(g : f, \theta) = \mathbb{E}[\log\{g(Y, p)/f(Y, p, \theta)\}].$$

We follow the similar proofs in Fan and Peng [2004], let  $a_N = \sqrt{K} (N^{-1/2} + \alpha_N)$  and set  $\|\mathbf{u}\|_2 = C$ , where  $C$  is a large enough constant, our aim is to show that for any given  $\epsilon$  there is a large constant  $C$  such that, for large  $N$  we have

$$\mathbb{P} \left\{ \sup_{\|\mathbf{u}\|_2=C} Q(\theta_* + a_N \mathbf{u}) < Q(\theta_*) \right\} \geq 1 - \epsilon.$$

This implies that with probability tending to 1 there is a local maximum  $\hat{\theta}$  in the call  $\{\theta_* + a_N \mathbf{u} : \|\mathbf{u}\|_2 \leq C\}$  such that  $\|\hat{\theta} - \theta_*\|_2 = O_p(a_N)$ . Because  $P_{\lambda_1}(0) = P_{\lambda_2}(0) = 0$ , We have

$$\begin{aligned} D(\mathbf{u}) &= Q(\theta_* + a_N \mathbf{u}) - Q(\theta_*) \\ &\leq \underbrace{L(\theta_* + a_N \mathbf{u}) - L(\theta_*)}_{(I)} \\ &\quad - \underbrace{N \sum_{j=1}^{s_1} \{P_{\lambda_1}(|\beta_{*,j} + a_N u_{1j}|) - P_{\lambda_1}(|\beta_{*,j}|)\}}_{(II)} \\ &\quad - \underbrace{N \sum_{j=1}^{s_2} \{P_{\lambda_2}(|\delta_{*,j} + a_N u_{2j}|) - P_{\lambda_2}(|\delta_{*,j}|)\}}_{(III)} \\ &:= (I) + (II) + (III), \end{aligned}$$

where  $\mathbf{u}^T = (\mathbf{u}_1^T, \mathbf{u}_2^T)$  with  $\mathbf{u}_1$  as  $K_1$  dimensions and  $\mathbf{u}_2$  as  $K_2$  dimensions. First for (II) we have

$$\begin{aligned} (II) &= - \sum_{j=1}^{s_1} \left[ N a_N P'_{\lambda_1}(|\beta_{*,j}|) \operatorname{sgn}(\beta_{*,j}) u_{1j} + N a_N^2 P''_{\lambda_1}(\beta_{*,j}) u_{1j}^2 \{1 + o(1)\} \right] \\ &:= I_1 + I_2, \\ |I_1| &\leq \sum_{j=1}^{s_1} |N a_N P'_{\lambda_1}(|\beta_{*,j}|) \operatorname{sgn}(\beta_{*,j}) u_{1j}| \leq \sqrt{s_1} N a_N \alpha_N \|\mathbf{u}_1\|_2 \leq N a_N^2 \|\mathbf{u}\|_2, \\ |I_2| &= \sum_{j=1}^{s_1} N a_N^2 P''_{\lambda_1}(|\beta_{*,j}|) u_{1j}^2 \{1 + o(1)\} \leq 2 \max_{1 \leq j \leq s_1} P''_{\lambda_1}(|\beta_{*,j}|) N a_N^2 \|\mathbf{u}\|_2^2. \end{aligned}$$

Similarly for (III). Then for (I) we have

$$\begin{aligned} (I) &= a_N \nabla^T L(\theta_*) \mathbf{u} + \frac{1}{2} \mathbf{u}^T \nabla^2 L(\theta_*) \mathbf{u} a_N^2 + \frac{1}{6} \nabla^T \{ \mathbf{u}^T \nabla^2 L(\theta_*) \mathbf{u} \} \mathbf{u} a_N^3 \\ &:= I_3 + I_4 + I_5, \end{aligned}$$

with the same proof in Theorem 1 in Fan and Peng [2004], by condition S8, we have

$$\begin{aligned} |I_3| &= |a_N \nabla^T L(\theta_*) \mathbf{u}| \leq a_N \|\nabla^T L(\theta_*)\|_2 \|\mathbf{u}\|_2 = O_p(a_N^2 N) \|\mathbf{u}\|_2. \\ I_4 &= \frac{1}{2} \mathbf{u}^T \left\{ \frac{1}{N} ([\nabla^2 L(\theta_*) - \mathbb{E} \{ \nabla^2 L(\theta_*) \}]) \right\} \mathbf{u} N a_N^2 \\ &\quad - \frac{1}{2} \mathbf{u}^T A(\theta_*) \mathbf{u} N a_N^2 \\ &= - \frac{N a_N^2}{2} \mathbf{u}^T A(\theta_*) \mathbf{u} + o_p(1) N a_N^2 \|\mathbf{u}\|_2^2. \end{aligned}$$

By condition S9 and  $K^4/N \rightarrow 0$  and  $K^2\alpha_N \rightarrow 0$  as  $N \rightarrow \infty$ , we have

$$\begin{aligned} |I_5| &= \left| \frac{1}{6} \sum_{i,j,k=1}^K \frac{\partial L(\theta^*)}{\partial \theta_i \partial \theta_j \partial \theta_k} u_i u_j u_k a_N^3 \right| \\ &\leq \frac{1}{6} \sum_{l=1}^N \left\{ \sum_{i,j,k=1}^K M_{ijk}^2(Y_i, p_i) \right\}^{1/2} \|\mathbf{u}\|_2^3 a_N^3 \\ &= o_p(N a_N^2) \|\mathbf{u}\|_2^2. \end{aligned}$$

Therefore, by Assumption 5 and allowing  $\|\mathbf{u}\|_2$  to be large enough, all  $I_1, I_2, I_3, I_5$  and (III) are dominated by  $I_4$ , which is negative, therefore proves  $\|\hat{\theta} - \theta_*\|_2 = O_p \left\{ \sqrt{K} (N^{-1/2} + a_N) \right\}$ . Further we have  $\max \left\{ \|\hat{\beta} - \beta_*\|_2, \|\hat{\delta} - \delta_*\|_2 \right\} \leq \|\hat{\theta} - \theta_*\|_2 = O_p \left\{ \sqrt{K} (N^{-1/2} + a_N) \right\}$ . For the SCAD penalty, it is clear that  $a_N = O_p(N^{-1/2})$ , therefore  $\hat{\beta}$  and  $\hat{\delta}$  are root- $(N/K)$ -consistent penalized likelihood estimators exist with probability tending to 1, and no requirements are imposed on the convergence rate of  $\lambda_1$  and  $\lambda_2$ .

## S2.4 PROOF FOR THEOREM 3

We follow the similar proofs in Fan and Peng [2004]. we first show that the nonconcave penalized estimator possesses the sparsity property  $\hat{\theta}_2 = 0$  by the following lemma.

**Lemma 1** *Assume Assumption 5, Assumption S1–S9 are satisfied, if  $\lambda_1, \lambda_2 \rightarrow 0$ ,  $\sqrt{N/K}\lambda_1 \rightarrow \infty$ ,  $\sqrt{N/K}\lambda_2 \rightarrow \infty$ , and  $K^5/N \rightarrow 0$  as  $N \rightarrow \infty$ , then first show that with probability tending to 1, for any given  $\theta_1$  satisfying  $\|\theta_1 - \theta_{*1}\|_2 = O_p(\sqrt{K/N})$  and any constant  $C$ ,*

$$Q \{(\theta_1^T, 0)^T\} = \max_{\|\theta_2\|_2 \leq C(K/N)^{1/2}} Q \{(\theta_1^T, \theta_2^T)^T\}.$$

Proof: Let  $\epsilon = C\sqrt{K/N}$ . It is sufficient to show that with probability tending to 1 as  $N \rightarrow \infty$ , for any  $\theta_1 - \theta_{*1} = O_p(\sqrt{K/N})$  we have for  $j = s+1, \dots, K$ ,

$$\begin{aligned} \frac{\partial Q(\theta)}{\partial \theta_j} &< 0 \quad \text{for } 0 < \theta_j < \epsilon, \\ \frac{\partial Q(\theta)}{\partial \theta_j} &> 0 \quad \text{for } -\epsilon < \theta_j < 0. \end{aligned}$$

By Taylor expansion,

$$\begin{aligned} \frac{\partial Q(\theta)}{\partial \theta_j} &= \frac{\partial L(\theta)}{\partial \theta_j} - NP'_\lambda(|\theta_j|)\text{sgn}(\theta_j) \\ &= \frac{\partial L(\theta_*)}{\partial \theta_j} + \sum_{l=1}^K \frac{\partial^2 L(\theta_*)}{\partial \theta_j \partial \theta_l} (\theta_l - \theta_{*,l}) \\ &\quad + \sum_{l,k=1}^K \frac{\partial^3 L(\theta^*)}{\partial \theta_j \partial \theta_l \partial \theta_k} (\theta_l - \theta_{*,l}) (\theta_k - \theta_{*,k}) \\ &\quad - NP'_\lambda(|\theta_j|)\text{sgn}(\theta_j) \\ &:= I_1 + I_2 + I_3 + I_4, \end{aligned}$$

where  $\theta^*$  lies between  $\theta$  and  $\theta_*$ , and  $P'_\lambda(|\theta_j|) = P'_{\lambda_1}(|\beta_j|)$  for  $j = s+1, \dots, K_1 - s_1 + s$ , and  $P'_\lambda(|\theta_j|) = P'_{\lambda_2}(|\delta_j|)$  for  $j = K_1 - s_1 + s + 1, \dots, K$ .

Following the same proof in Fan and Peng [2004], we prove  $I_1 + I_2 + I_3 = O_p(\sqrt{NK})$ . First,  $I_1 = O_p(\sqrt{N}) = O_p(\sqrt{NK})$ . Also,

$$\begin{aligned} I_2 &= \sum_{l=1}^K \left( \frac{\partial^2 L(\theta_*)}{\partial \theta_j \partial \theta_l} - \mathbb{E} \left\{ \frac{\partial^2 L(\theta_*)}{\partial \theta_j \partial \theta_l} \right\} \right) (\theta_l - \theta_{*,l}) \\ &\quad + \sum_{l=1}^K \mathbb{E} \left\{ \frac{\partial^2 L(\theta_*)}{\partial \theta_j \partial \theta_l} \right\} (\theta_l - \theta_{*,l}) \\ &:= S_1 + S_2. \end{aligned}$$

Using the Cauchy-Schwarz inequality and  $\|\theta - \theta_*\|_2 = O_p(K/N)$ , we have

$$\begin{aligned} |S_2| &= |N \sum_{l=1}^K A(\theta_*)(j, l) (\theta_l - \theta_{*,l})| \\ &\leq N O_p \left( \sqrt{\frac{K}{N}} \right) \left\{ \sum_{l=1}^K A^2(\theta_*)(j, l) \right\}^{1/2}. \end{aligned}$$

By Assumption S8, as the eigenvalues of the  $A(\theta)$  are bounded, we have  $S_2 = O_p(\sqrt{NK})$ . On the other hand,

$$|S_1| \leq \|\theta - \theta_*\|_2 \left( \sum_{l=1}^K \left[ \frac{\partial^2 L(\theta_*)}{\partial \theta_j \partial \theta_l} - \mathbb{E} \left\{ \frac{\partial^2 L(\theta_*)}{\partial \theta_j \partial \theta_l} \right\} \right]^2 \right)^{1/2}.$$

By Assumption S8, we have

$$\left( \sum_{l=1}^K \left[ \frac{\partial^2 L(\theta_*)}{\partial \theta_j \partial \theta_l} - \mathbb{E} \left\{ \frac{\partial^2 L(\theta_*)}{\partial \theta_j \partial \theta_l} \right\} \right]^2 \right)^{1/2} = O_p(\sqrt{NK}).$$

Therefore  $S_1 = O_p(\sqrt{NK})$  and  $I_2 = O_p(\sqrt{NK})$ . Further,

$$\begin{aligned} I_3 &= \sum_{l,k=1}^K \left[ \frac{\partial^3 L(\theta_*)}{\partial \theta_j \partial \theta_l \partial \theta_k} - \mathbb{E} \left\{ \frac{\partial^3 L(\theta_*)}{\partial \theta_j \partial \theta_l \partial \theta_k} \right\} \right] (\theta_l - \theta_{*,l}) (\theta_k - \theta_{*,k}) \\ &\quad + \sum_{l,k=1}^K \mathbb{E} \left\{ \frac{\partial^3 L(\theta_*)}{\partial \theta_j \partial \theta_l \partial \theta_k} \right\} (\theta_l - \theta_{*,l}) (\theta_k - \theta_{*,k}) \\ &:= S_3 + S_4. \end{aligned}$$

By Assumption S9,  $|S_4| \leq C_5^{1/2} NK \|\theta - \theta_*\|_2^2 = O_p(K^2) = o_p(\sqrt{NK})$ . Further,

$$S_3 \leq \sum_{l,k=1}^K \left[ \frac{\partial^3 L(\theta_*)}{\partial \theta_j \partial \theta_l \partial \theta_k} - \mathbb{E} \left\{ \frac{\partial^3 L(\theta_*)}{\partial \theta_j \partial \theta_l \partial \theta_k} \right\} \right]^2 \|\theta - \theta_*\|_2^4,$$

where under the Assumption S9 and Assumption 5,  $S_3 = O_p \left\{ \left( NK^2 \frac{K^2}{N^2} \right)^{1/2} \right\} = o_p(\sqrt{NK})$ . Then

$$I_1 + I_2 + I_3 = O_p(\sqrt{NK}).$$

Because we focus on the SCAD penalty, Fan and Peng [2004] illustrates that under Assumption 5, the SCAD penalty satisfies that

$$\begin{aligned} \liminf_{N \rightarrow +\infty} \liminf_{\beta \rightarrow 0+} P'_{\lambda_1}(\beta) / \lambda_1 &> 0 \\ \liminf_{N \rightarrow +\infty} \liminf_{\delta \rightarrow 0+} P'_{\lambda_2}(\delta) / \lambda_2 &> 0, \end{aligned}$$

therefore from

$$\frac{\partial Q(\theta)}{\partial \theta_j} = N\lambda \left\{ -\frac{P'_\lambda(|\theta_j|)}{\lambda} \text{sgn}(\theta_j) + O_p\left(\sqrt{\frac{K}{N}}/\lambda\right) \right\},$$

where  $\lambda = \lambda_1$  if  $j = s+1, \dots, K_1 - s_1 + s$  and  $\lambda = \lambda_2$  if  $j = K_1 - s_1 + s + 1, \dots, K$ , and  $P'_\lambda(|\theta_j|) = P'_{\lambda_1}(|\beta_j|)$  for  $j = s+1, \dots, K_1 - s_1 + s$ , and  $P'_\lambda(|\theta_j|) = P'_{\lambda_2}(|\delta_j|)$  for  $j = K_1 - s_1 + s + 1, \dots, K$ , the sign of  $\theta_j$  completely determines the sign of  $\partial Q(\theta)/\partial \theta_j$ . We complete the proof of Lemma 1.

By Lemma 1 we prove  $\hat{\theta}_2 = \begin{pmatrix} \hat{\beta}_2 \\ \hat{\delta}_2 \end{pmatrix} = 0$ . Then we prove the part 2.

Let

$$\begin{aligned} \Sigma &= \text{diag} \left\{ P''_\lambda(\theta_{*,1}), \dots, P''_\lambda(\theta_{*,s}) \right\} \\ &= \text{diag} \left\{ P''_{\lambda_1}(\beta_{*,1}), \dots, P''_{\lambda_1}(\beta_{*,s_1}), P''_{\lambda_2}(\delta_{*,1}), \dots, P''_{\lambda_2}(\delta_{*,s_2}) \right\} \end{aligned}$$

and

$$\begin{aligned} b &= \left\{ P'_\lambda(|\theta_{*,1}|) \text{sgn}(\theta_{*,1}), \dots, P'_\lambda(|\theta_{*,s}|) \text{sgn}(\theta_{*,s}) \right\}^\top \\ &= \left\{ P'_{\lambda_1}(|\beta_{*,1}|) \text{sgn}(\beta_{*,1}), \dots, P'_{\lambda_1}(|\beta_{*,s_1}|) \text{sgn}(\beta_{*,s_1}), P'_{\lambda_2}(|\delta_{*,1}|) \text{sgn}(\delta_{*,1}), \dots, P'_{\lambda_2}(|\delta_{*,s_2}|) \text{sgn}(\delta_{*,s_2}) \right\}^\top. \end{aligned}$$

If we can show that

$$\{A(\theta_{*1}) + \Sigma\} (\hat{\theta}_1 - \theta_{*1}) + b = \frac{1}{N} \nabla L(\theta_{*1}) + o_p(N^{-1/2}),$$

then

$$\begin{aligned} &\sqrt{N} W A^{-1/2}(\theta_{*1}) \{A(\theta_{*1}) + \Sigma\} \left[ \hat{\theta}_1 - \theta_{*1} + \{A(\theta_{*1}) + \Sigma\}^{-1} b \right] \\ &= \frac{1}{\sqrt{N}} W A^{-1/2}(\theta_{*1}) \nabla L(\theta_{*1}) + o_p \left\{ W A^{-1/2}(\theta_{*1}) \right\} \\ &= \frac{1}{\sqrt{N}} W A^{-1/2}(\theta_{*1}) \nabla L(\theta_{*1}) + o_p(1). \end{aligned}$$

Let  $R_i = \frac{1}{\sqrt{N}} W A^{-1/2}(\theta_{*1}) \nabla L_i(\theta_{*1})$ ,  $i = 1, \dots, N$ . Following the same proof in Fan and Peng [2004], for any  $\epsilon$ , we have

$$\begin{aligned} \sum_{i=1}^N \mathbb{E} \|R_i\|_2^2 \mathbf{1} \{\|R_i\|_2 > \epsilon\} &= N \mathbb{E} \|R_1\|_2^2 \mathbf{1} \{\|R_1\|_2 > \epsilon\}, \\ &\leq N (\mathbb{E} \|R_1\|_2^4)^{1/2} \{\mathbb{P}(\|R_1\|_2 > \epsilon)\}^{1/2}. \end{aligned}$$

By Assumption S8 and  $W W^\top \rightarrow G$ , we obtain

$$\mathbb{P}(\|R_1\|_2 > \epsilon) \leq \frac{\mathbb{E} \|W A^{-1/2}(\theta_{*1}) \nabla L_1(\theta_{*1})\|_2^2}{N \epsilon^2} = O(N^{-1})$$

and

$$\begin{aligned} \mathbb{E} \|R_1\|_2^4 &= \frac{1}{N^2} \mathbb{E} \|W A^{-1/2}(\theta_{*1}) \nabla L_1(\theta_{*1})\|_2^4 \\ &\leq \frac{1}{N^2} \lambda_{\max}(W W^\top) \lambda_{\max} \{A^{-1}(\theta_{*1})\} \mathbb{E} \|\nabla^\top L_1(\theta_{*1}) \nabla L_1(\theta_{*1})\|_2^2 \\ &\leq O\left(\frac{K^2}{N^2}\right). \end{aligned}$$

Thus, we have

$$\sum_{i=1}^N \mathbb{E} \|R_i\|_2^2 \mathbf{1} \{\|R_i\|_2 > \epsilon\} = O\left(N \frac{K}{N} \frac{1}{\sqrt{N}}\right) = o(1).$$

and

$$\sum_{i=1}^N \text{cov}(R_i) = \text{cov} \left\{ W A^{-1/2}(\theta_{*1}) \nabla L_1(\theta_{*1}) \right\} = W A^{-1/2}(\theta_{*1}) B(\theta_{*1}) A^{-1/2}(\theta_{*1}) W^T,$$

so that the  $R_i$  satisfies the conditions of the Lindeberg-Feller central limit theorem. Further, using the Taylor expansion on  $\nabla Q(\hat{\theta}_1)$  at the point  $\theta_{*1}$ , we have

$$\begin{aligned} & \frac{1}{N} \left[ \{ \nabla^2 L(\theta_{*1}) - \nabla^2 P_\lambda(\theta_1^{**}) \} (\hat{\theta}_1 - \theta_{*1}) - \nabla P_\lambda(\theta_{*1}) \right] \\ &= -\frac{1}{N} \left[ \nabla L(\theta_{*1}) + \frac{1}{2} (\hat{\theta}_1 - \theta_{*1})^T \nabla^2 \{ \nabla L(\theta_1^*) (\hat{\theta}_1 - \theta_{*1}) \} \right], \end{aligned}$$

where  $\theta_1^*$  and  $\theta_1^{**}$  lie between  $\hat{\theta}_1$  and  $\theta_{*1}$ . Now define

$$\mathcal{L} := \nabla^2 L(\theta_{*1}) - \nabla^2 P_\lambda(\theta_1^{**})$$

and

$$C := \frac{1}{2} (\hat{\theta}_1 - \theta_{*1})^T \nabla^2 \{ \nabla L(\theta_1^*) (\hat{\theta}_1 - \theta_{*1}) \}.$$

Following the proof in Fan and Peng [2004], under Assumption S9 and Assumption 5 and by the Cauchy-Schwarz inequality, we have  $\|1/NC\|_2^2 = o_p(1/N)$ . Further, we have

$$\lambda_i \left\{ \frac{1}{N} \mathcal{L} + A(\theta_{*1}) + \Sigma \right\} = o_p \left( \frac{1}{\sqrt{K}} \right), i = 1, \dots, s,$$

where  $\lambda_i(M)$  is the  $i$ th eigenvalue of a symmetric matrix  $M$ . Therefore,

$$\left\{ \frac{1}{N} \mathcal{L} + A(\theta_{*1}) + \Sigma \right\} (\hat{\theta}_1 - \theta_{*1}) = o_p \left( \frac{1}{\sqrt{N}} \right).$$

Then, we have  $\{A(\theta_{*1}) + \Sigma\} (\hat{\theta}_1 - \theta_{*1}) + b = \frac{1}{N} \nabla L(\theta_{*1}) + o_p(N^{-1/2})$ , and finally we have

$$\begin{aligned} & \sqrt{N} W A^{-1/2}(\theta_{*1}) \{A(\theta_{*1}) + \Sigma\} \left[ \hat{\theta}_1 - \theta_{*1} + \{A(\theta_{*1}) + \Sigma\}^{-1} b \right] \\ & \rightarrow \mathcal{N}(0, W A^{-1/2}(\theta_{*1}) [B(\theta_{*1}) - \mathbb{E} \{ \nabla L_1(\theta_{*1}) \}] A^{-1/2}(\theta_{*1}) W^T). \end{aligned}$$

Further, based on the SCAD penalty,  $\Sigma = 0$  and  $b = 0$ , therefore, we have

$$\sqrt{N} W A^{1/2}(\theta_{*1}) (\hat{\theta}_1 - \theta_{*1}) \rightarrow \mathcal{N} \left( 0, W A^{-1/2}(\theta_{*1}) B(\theta_{*1}) A^{-1/2}(\theta_{*1}) W^T \right).$$

If the model is correctly specified, i.e.,  $g(Y, p_i) = f(Y, p_i, \theta)$  for some  $\theta \in \Theta$ , then  $\theta_0 = \theta_*$ , and

$$\sqrt{N} W I^{1/2}(\theta_{01}) (\hat{\theta}_1 - \theta_{01}) \rightarrow \mathcal{N}(0, W W^T).$$

We finish the second part.

## S2.5 COMMENTS ON THEOREM 2 AND THEOREM 3

It's important to clarify that if our focus isn't primarily on the least square loss, we can adjust the assumptions on  $f$  and  $g$  from S1-S10 to only satisfy Assumptions S1-S9 for Theorem 2 and Theorem 3. The critical point is that for any other  $f$  and  $g$  fulfilling Assumptions S1-S9, a local maximizer  $\hat{\theta}$  of the  $Q(\theta)$  exists such that  $\hat{\theta}$  converges to  $\theta_*$ , where  $\theta_*$  minimizes the KLIC between  $f$  and  $g$ , and also exhibits the oracle property and asymptotic normality. However, in the main paper, our primary concern is the least square loss, and we further constrain  $f$  to meet assumption S10 in addition to assumptions S1-S9. Assumption S10 implies that the parameters minimizing the KLIC also minimize a specific form: the least square form. This is because in the main paper, we are particularly interested in ANCOVA least square estimates.



## S2.6 PROOF FOR THEOREM 4

The consistency and the asymptotic normality of the estimator  $\hat{\tau}$  are from Theorem 2 and Theorem 3. Here we only prove the calculation of the variance. If  $\mathbb{V}(\epsilon) = \sigma^2$ , then the influence function of  $\theta$  is

$$\phi_{\theta}(x, y, s) = \mathbb{E} \left\{ \left( \begin{array}{c} p_{\mu} \\ (1-S)p_b(X) \end{array} \right) (p_{\mu}^T, (1-S)p_b^T(X)) \right\}^{-1} \left( \begin{array}{c} p_{\mu} \\ (1-s)p_b(x) \end{array} \right) (y - p^T\theta).$$

Then the asymptotic variance of  $\hat{\beta}$  is

$$\begin{aligned} \mathbb{V}(\hat{\beta}) &= \mathbb{E} \{ \phi_{\theta}(X, Y, S) \phi_{\theta}^T(X, Y, S) \}_{11} \\ &= \sigma^2 \mathbb{E} \left\{ \left( \begin{array}{cc} p_{\mu} p_{\mu}^T & (1-S)p_{\mu} p_b^T(X) \\ (1-S)p_b(X) p_{\mu}^T & (1-S)p_b(X) p_b^T(X) \end{array} \right) \right\}_{11}^{-1} \\ &= \sigma^2 \mathbb{E} \left[ S p_{\mu} p_{\mu}^T + (1-S) p_{\mu} p_{\mu}^T - (1-S) \{ p_{\mu} p_b^T(X) \} \{ p_b(X) p_b^T(X) \}^{-1} \{ p_b(X) p_{\mu}^T \} \right]^{-1}. \end{aligned}$$

On the other hand, under the RE data, the ANCOVA working model is

$$Y = \beta_{\text{int}} + \beta_A A + \beta_X^T p_{\mu}(X) + \epsilon, \quad \epsilon \sim (0, \sigma^2),$$

Similarly, the asymptotic variance of  $\hat{\beta}_{\text{RE}}$  is

$$\mathbb{V}(\hat{\beta}_{\text{RE}}) = \sigma^2 \mathbb{E} (S p_{\mu} p_{\mu}^T)^{-1}.$$

Then by Holder inequality,

$$\begin{aligned} &\mathbb{E} \{ (1-S) p_{\mu} p_{\mu}^T \} \\ &- \mathbb{E} \left[ (1-S) \{ p_{\mu} p_b^T(X) \} \{ p_b(X) p_b^T(X) \}^{-1} \{ p_b(X) p_{\mu}^T \} \right] \geq 0, \end{aligned}$$

where the inequality holds iff  $p_{\mu} = M p_b(X)$  for some matrix  $M$ . Therefore, we have  $\mathbb{V} \{ \hat{\beta} \} \leq \mathbb{V} \{ \hat{\beta}_{\text{RE}} \}$ , and the inequality holds iff  $p_{\mu} = M p_b(X)$  for some matrix  $M$ . Therefore  $\mathbb{V}(\hat{\tau}_{\text{RE}}) \geq \mathbb{V}(\hat{\tau})$ .

## S3 TOY EXAMPLE

Consider the case  $y = x^T \beta + (1-s)x^T \delta + \epsilon$ ,  $x = (x_1, \dots, x_{50})^T \in \mathbb{R}^{50}$ ,  $\beta^T = (1, \dots, 50)/10$ ,  $\delta^T = (1, \dots, 50) \times 30$ ,  $\epsilon \sim \mathcal{N}(0, 1)$ , where  $s$  denotes the zero-one indicator variable that determines whether the observation belongs to the REs, for simplicity, we assume  $x^T \beta$  is the correct outcome mean function of the REs, and  $x^T \delta$  is the bias function reflecting the difference between ECs and REs, i.e., if  $\delta = 0$ , then the observed covariates capture all confounders in the ECs and REs and thus the exchangeability assumption is valid. For didactic purposes, the magnitude of  $\delta$  is much larger than the magnitude of  $\beta$ . Using the same regularization parameter appears to assign the same weight for  $\beta$  and  $\delta$ , thus any penalty regularization methods tend to omit small signals  $\beta$  and only pick up big signals  $\delta$ . Therefore, in order to make penalizations between different parameters comparable, it is crucial to add regularizations to  $\beta$  and  $\delta$  separately. Figure S1 shows the smoothed linear regression between  $\hat{\beta}$  and  $\beta$  after applying the single-penalty regularization (denoted as ‘‘Single’’) and the double-penalty regularizations (denoted as ‘‘Double’’) to select variables and refitting the model using selected variables, where double penalties make  $\hat{\beta}$  more accurate than the single penalties.

## S4 PLOTS

We present the extra figures for the first simulation study in this section. Figure S2 shows the results for the case  $\|\beta_0\|_1 > \|\delta_0\|_1$  when setting half of parameters in  $\delta_0$  equal to zero:  $\|\beta_0\|_1 = c \|\delta_0\|_1$ ,  $c = 1, 3, 5, 7, 9$ ; and Figure S3 shows the results for varying the sparsity level of  $\delta_0$  when setting  $\|\beta_0\|_1 = \|\delta_0\|_1$  with the x axis as the the number of variables in  $\delta_0$  equal to zero. Each figure shows the MSE results and the percentage of Under-select and Over-select. Figure S2 shows the SPIE has a larger MSE compared to the DPIE, which is consistent with the theoretical results. On the other hand, the changes in the sparsity of  $\delta_0$  affect little on the results. This finding also consists of the theoretical result, where we only need to restrain the magnitude of different parameters to guarantee the consistency and oracle properties.

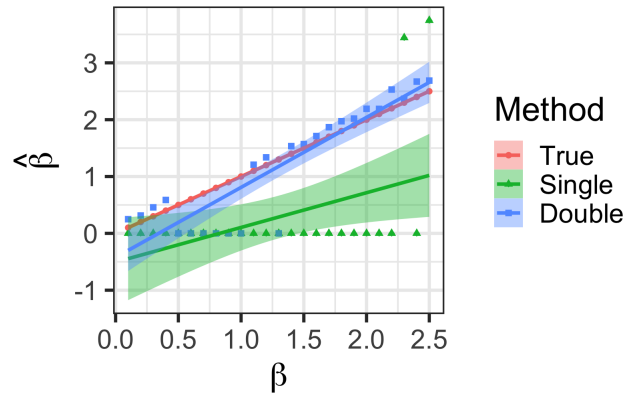


Figure S1: The smoothed linear regression between  $\hat{\beta}$  and  $\beta$  with the 95% confidence intervals as the shaded area and  $(\beta, \hat{\beta})$  as the points.

## References

- Xiaohong Chen. Large sample sieve estimation of semi-nonparametric models. *Handbook of Econometrics*, 6:5549–5632, 2007.
- Jianqing Fan and Heng Peng. Nonconcave penalized likelihood with a diverging number of parameters. *The Annals of Statistics*, 32(3):928–961, 2004.
- G.G. Lorentz. *Approximation of Functions*. Holt, Rinehart and Winston, 1966.
- Michael Rosenblum and Mark J Van Der Laan. Using regression models to analyze randomized trials: Asymptotically valid hypothesis tests despite incorrectly specified models. *Biometrics*, 65(3):937–945, 2009.
- Bingkai Wang, Ryoko Susukida, Ramin Mojtabei, Masoumeh Amin-Esmaeili, and Michael Rosenblum. Model-robust inference for clinical trials that improve precision by stratified randomization and covariate adjustment. *Journal of the American Statistical Association*, 118(542):1152–1163, 2023.
- Halbert White. Maximum likelihood estimation of misspecified models. *Econometrica*, 50(1):1–25, 1982.

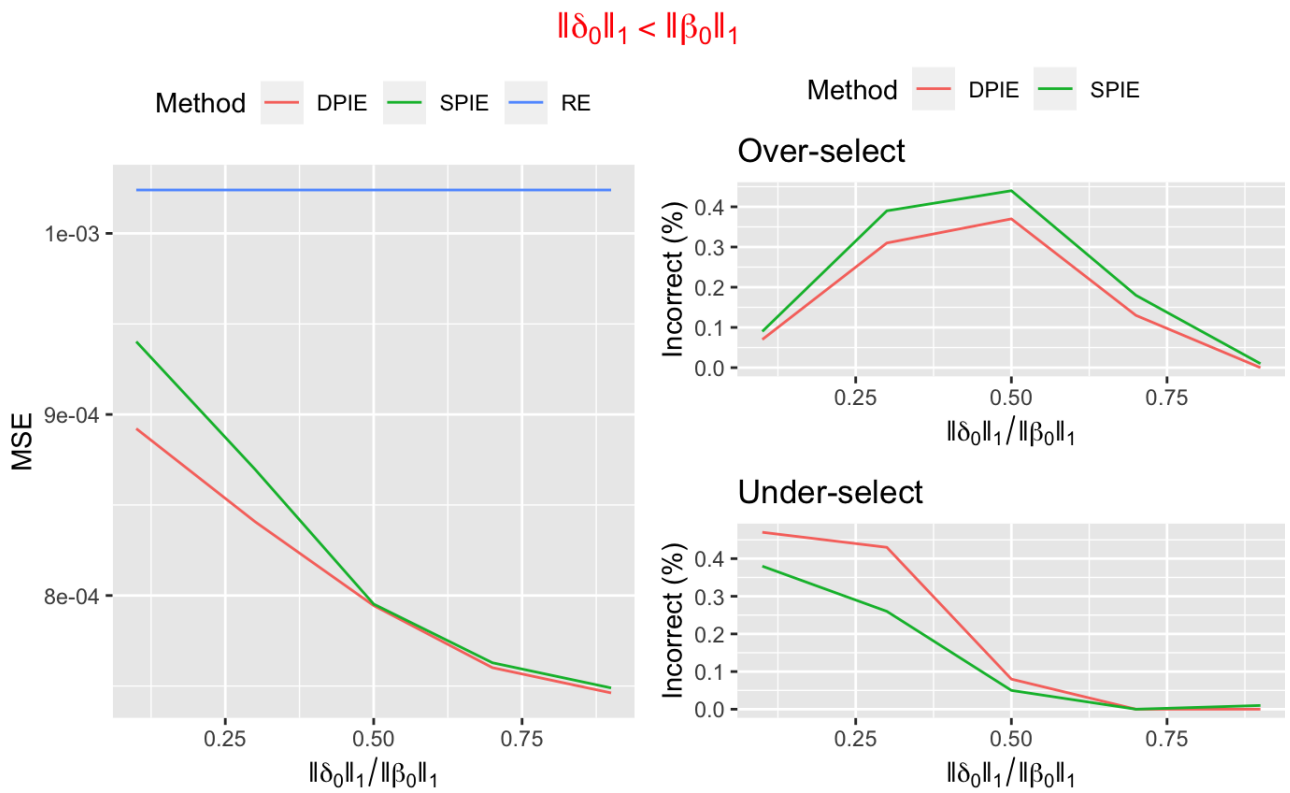


Figure S2: Simulation results based on 100 Monte Carlo times. The left panel shows the MSE versus the magnitude ratio between  $\delta$  and  $\beta$ . The right panel shows the percentage of wrongly choosing more and less parameters, separately.

### Varing Sparsity Level

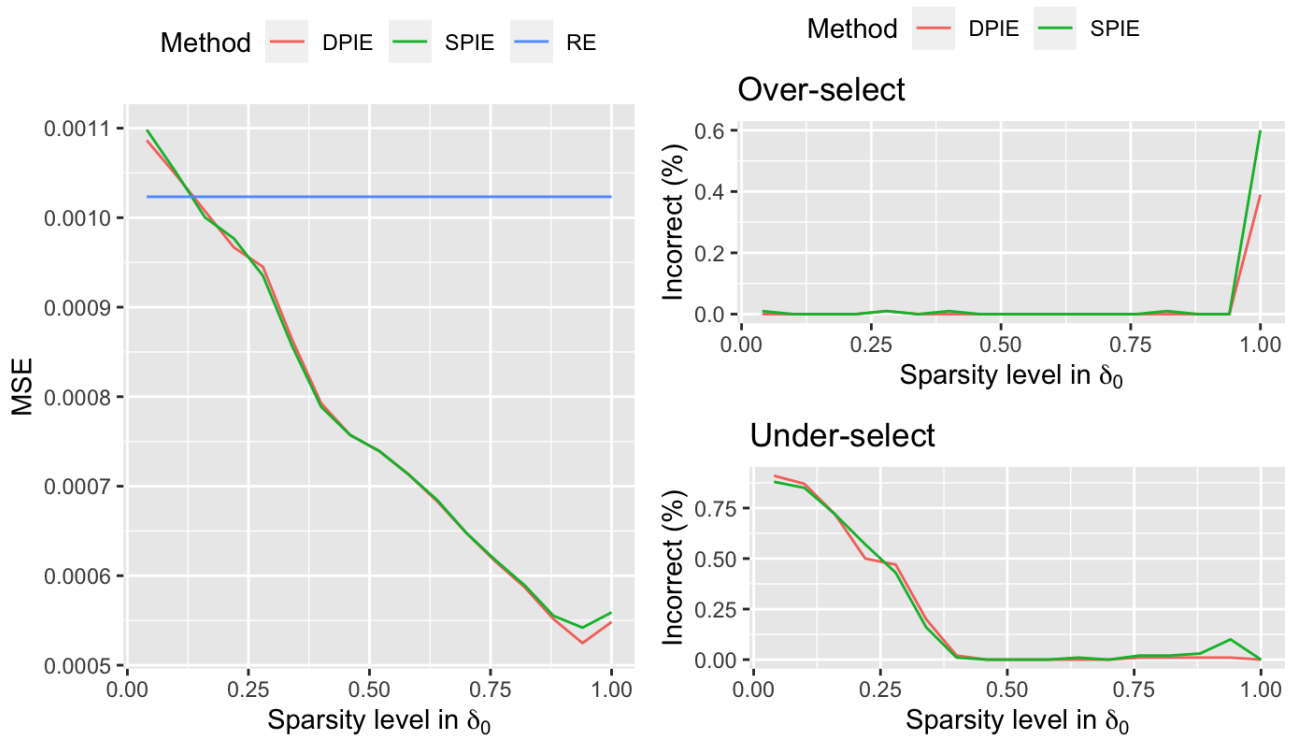


Figure S3: Simulation results based on 100 Monte Carlo times. The left panel shows the MSE versus the sparsity level in  $\delta$ . The right panel shows the percentage of wrongly choosing more and less parameters, separately.